



Oei, L., Koromani, F., Breda, S. J., Schousboe, J. T., Clark, E. M., van Meurs, J. B. J., ... Rivadeneira, F. (2018). Osteoporotic Vertebral Fracture Prevalence Varies Widely Between Qualitative and Quantitative Radiological Assessment Methods: The Rotterdam Study. *Journal of Bone and Mineral Research*, 33(4), 560-568. <https://doi.org/10.1002/jbmr.3220>

Peer reviewed version

License (if available):
Other

Link to published version (if available):
[10.1002/jbmr.3220](https://doi.org/10.1002/jbmr.3220)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Wiley at <http://onlinelibrary.wiley.com/doi/10.1002/jbmr.3220/abstract>. Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/pure/about/ebr-terms>

1 **Osteoporotic Vertebral Fracture Prevalence Varies Widely Between Qualitative and**
2 **Quantitative Radiological Assessment Methods: The Rotterdam Study**

3 Oei L^{1,2*}, Koromani F^{1,2,3*}, Breda SJ³, Schousboe JT⁴, Clark EM⁵, van Meurs JBJ¹, Ikram
4 MA², Waarsing JH⁶, van Rooij FJA², Zillikens MC¹, Krestin GP³,
5 Oei EHG^{3**}, Rivadeneira F^{1,2**}

6 ¹ *Department of Internal Medicine, Erasmus MC, Rotterdam, The Netherlands*

7 ² *Department of Epidemiology, Erasmus MC, Rotterdam, The Netherlands*

8 ³ *Department of Radiology, Erasmus MC, Rotterdam, The Netherlands*

9 ⁴ *Park Nicollet Clinic and HealthPartners Institute, HealthPartners Inc., Minneapolis, MN, USA*

10 ⁵ *Musculoskeletal Research Unit, School of Clinical Science, University of Bristol, Southmead*

11 *Hospital, Bristol, UK*

12 ⁶ *Department of Orthopedics, Erasmus MC,*

13

14 Keywords: osteoporosis; fracture; vertebral; diagnosis

15 Word count: **Text: 4,117 Abstract: 281 Tables: 5 Figures: 4**

16 Osteoporosis, Epidemiology, Screening, Radiology

17 Corresponding author:

18 Fernando Rivadeneira MD, PhD

19 Departments of Internal Medicine

20 and Epidemiology

21 P.O. Box 2040 Ee5-59b

22 3000CA Rotterdam

23 The Netherlands

24 Email: f.rivadeneira@erasmusmc.nl

25 Phone number: +31 10 7044015

26 Supplemental data: Tables 2, Figures 2

27

28

29

30 ABSTRACT

31

32 Background: Accurate diagnosis of vertebral osteoporotic fractures is crucial for the
33 identification of individuals at high risk of future fractures. Different methods for radiological
34 assessment of vertebral fractures exist, but a gold standard is lacking.

35 The aim of our study was to estimate statistical measures of agreement and prevalence of

36 osteoporotic vertebral fractures in the population-based Rotterdam Study, across two

37 assessment methods. Methods: The quantitative morphometry assisted by SpineAnalyzer®

38 (QM SA) method, evaluates vertebral height loss that affects vertebral shape whereas the

39 algorithm based qualitative (ABQ) method judges endplate integrity and includes guidelines

40 for the differentiation of vertebral fracture and non- fracture deformities.

41 Results: Cross-sectional radiographs were assessed for 7,582 participants aged 45-95 years.

42 With QM SA, the prevalence was 14.2% (95% CI: 13.4% to 15.0%), compared to 4.0% (95%

43 CI: 3.6% to 4.5%) with ABQ. Inter-method agreement according to kappa (κ) was 0.24. The

44 highest agreement between methods was among females ($\kappa=0.31$), participants aged above 80

45 ($\kappa =0.40$) and at the L1 level ($\kappa =0.40$). With ABQ, most fractures were found at the thoraco-

46 lumbar junction (T12-L1) followed by the T7-T8 level, whereas with QM SA, most

47 deformities were in the mid (T7-T8) and lower thoracic spine (T11-T12) with similar number

48 of fractures in both peaks. Excluding mild deformities (grade 1 with QM) from the analysis

49 increased the agreement between the methods from $\kappa=0.24$ to $\kappa=0.40$, whereas re-examining

50 mild deformities based on endplate depression increased agreement from $\kappa =0.24$ to 0.50 (p-

51 value< 0.001).

52 Conclusion: Vertebral fracture prevalence differs significantly between QM SA and ABQ; re-

53 examining QM mild deformities based on endplate depression would increase the agreement

54 between methods. More wide-spread and consistent application of an optimal method may
55 improve clinical care.

56 **Introduction**

57 Of all osteoporotic fractures, vertebral fractures are the most common type.⁽¹⁾
58 Vertebral fractures have been synonymous with the diagnosis of osteoporosis since its earliest
59 description as a metabolic bone disorder.⁽²⁾ Furthermore, osteoporotic vertebral fractures are a
60 major health problem worldwide. Given the ageing of populations, osteoporotic vertebral
61 fractures are likely to become an even increasingly important health issue. The costs of
62 osteoporotic vertebral fractures were estimated to be € 1.5 billion in Europe in 2010⁽³⁾ and are
63 expected to have increased by more than 50% by 2025.⁽⁴⁾

64 Vertebral fractures may occur in the absence of trauma or after normal activities
65 involving bending, lifting or turning.⁽¹⁾ Although, two thirds of vertebral fractures are not
66 clinically detected, they are associated with decreased quality of life, back pain, functional
67 limitations⁽⁵⁾ and mortality⁽⁶⁾ and can only be detected by formal screening. Vertebral
68 fractures are often a first presentation of osteoporosis, therefore, accurate diagnosis is
69 important to identify patients at high risk for future fractures. It has been shown that women
70 with preexisting vertebral fractures have four times greater risk of subsequent vertebral
71 fractures and 1.5 to 2 times greater risk of non-vertebral fractures than those without prior
72 fractures, and this risk increases with the number and severity of prior vertebral fractures.⁽⁷⁻⁹⁾
73 It is important to detect these fractures, since anti-osteoporotic therapy has been proven
74 highly effective in reducing the risk of both non-vertebral and vertebral fractures.

75 Several methods for radiological assessment of vertebral fractures exist, but a gold
76 standard is lacking.⁽¹⁰⁾ The most commonly applied assessment methods include (semi-)
77 quantitative morphometry (QM) and the algorithm based qualitative (ABQ) method. In
78 contrast to semi-quantitative methods relying on expert visual inspection of height reduction,
79 actual QM-based methods determine relative vertebral height loss by calculating ratios of the
80 measured vertebral heights. Rather than only placing morphometry points manually on a

81 vertebral body, software packages such as Spine Analyzer®⁽¹¹⁾ apply Genant's
82 classification⁽¹²⁾ to define vertebral deformities. Finally, the algorithm based qualitative
83 (ABQ) method by Jiang et al.⁽¹³⁾ mainly judges endplate integrity, regardless of vertebral
84 height reduction, and includes defined guidelines for the differentiation of vertebral fracture
85 and non-fracture deformities. The key assumption is that the endplate is always deformed in
86 vertebral fractures, and therefore endplate depression has perfect specificity for vertebral
87 fracture. Vertebral height may appear to be decreased as a result of oblique image projection,
88 specific diseases, and anatomical variants that can mimic vertebral fractures.⁽¹²⁻¹⁵⁾ To deal
89 with this misclassification, ABQ uses an algorithm to systematically rule out non-fracture
90 deformities.

91 The aim of our study was to analyze differences in prevalence and fracture location
92 between methods. We applied two methods, i.e., ABQ and SpineAnalyzer® software-assisted
93 QM, for assessing vertebral fractures in the population-based Rotterdam Study, an ongoing
94 prospective cohort study in elderly persons.

95

96 **Materials and Methods**

97 *The Rotterdam Study:* The Rotterdam Study is a prospective population-based cohort
98 studying the determinants of chronic diseases and disability in Dutch men and women. Both
99 the objectives and the study design have been described previously.⁽¹⁶⁾ The study targets
100 investigations on endocrine diseases like osteoporosis amongst others. It includes 14,926
101 inhabitants aged 45 years and over of Rotterdam city's Ommoord district in The Netherlands.
102 *Vertebral fracture assessment:* Radiographic examinations of the spine were obtained by a
103 digitized Fuji FCR system (FUJIFILM Medical Systems). All radiographs were acquired
104 according to a standardized protocol with a focus film distance of 120 cm. In some instances
105 evaluability was suboptimal, mostly in the upper spine levels (supplementary Fig 1). In the

106 current report we have included participants with sufficient evaluability from T4-L4. Two
107 teams, each composed of seven trained research assistants assessed lateral spine radiographs
108 (T4-L4) independent of each other, using either ABQ or software-assisted QM
109 (SpineAnalyzer[®], Optasia Medical Ltd, Cheadle, UK). The mean inter- observer agreement
110 for ABQ according to kappa statistic (κ) was moderate for both QM SA and ABQ ($\kappa= 0.51$
111 and $\kappa=0.53$ respectively). A subset of 76 radiographs were scored by two independent
112 external readers; one reader with ABQ and one reader with QM SA; the agreement was poor,
113 $\kappa= 0.19$. With ABQ, radiographs were triaged as normal, uncertain or definite fracture, based
114 on integrity of the endplates. Definite and uncertain vertebral fractures were re-assessed by a
115 musculoskeletal radiologist. SpineAnalyzer[®] software automatically identifies vertebral shape
116 to calculate the exact heights of the vertebrae. After labeling the vertebrae of interest by
117 placing thirteen points at the center of each vertebral body from L4 to T4, SpineAnalyzer[®]
118 will place six morphometry points for each labeled vertebra, corresponding to the four
119 corners and the middle of the vertebral body. The analyst can make manual adjustments to
120 these six morphometry points to fine-tune their exact locations. The morphometry points are
121 used to assess reductions in anterior, middle and posterior heights of the vertebrae by
122 determining if one height measure is “reduced” in relation to another height (e.g., anterior
123 height/posterior height <1 for a wedge shaped deformity). The SpineAnalyzer[®] software
124 output provides a classification for deformities of shape (wedge, biconcave, crush) and
125 severity (mild, moderate, severe). The wedge ratio is calculated by dividing anterior height by
126 posterior height (hA/hP). Biconcavity is calculated by dividing mid height by posterior height
127 (hM/hP). The calculation of crush fractures makes use of adjacent vertebral heights. Height
128 loss less than 20% is considered normal. Mild fracture (grade one) is defined as height loss
129 between $\geq 20\%$ and $<25\%$, moderate fracture (grade two) between $\geq 25\%$ and $<40\%$ and

130 severe fracture (grade three) $\geq 40\%$ according to Genant's classification scheme for
131 osteoporotic vertebral fractures.⁽¹²⁾

132 *Incident fracture* were new fractures identified and reported by the general practitioners
133 (GPs) or assessed from hospital records that occurred after baseline assessment. All events
134 were then reviewed and coded by a research physician. For the current study we examined
135 incident non-vertebral, hip and clinical-vertebral fractures.

136 *Statistical analysis:* We compared fracture prevalence and distribution according to vertebral
137 level for QM SA and ABQ. Since there is no consensus whether most of the grade 1 or mild
138 deformities are true osteoporotic vertebral fractures or not⁽¹⁴⁾, we performed secondary
139 analyses by excluding those fractures from the analysis. Agreement between the diagnostic
140 approaches (inter- method agreement) and between raters (inter-rater agreement) for the
141 identification of prevalent vertebral fractures was analyzed using kappa . The kappa value
142 takes into account the proportion of agreement attributable to chance alone and can range
143 from 0 (no agreement) to 1 (complete agreement); values greater than 0.8 are considered
144 strong and values lower than 0.6 moderate⁽¹⁷⁾. Given that kappa is influenced by the
145 imbalances in the distribution of marginal totals in the 2x2 table ^(18,19), together with kappa
146 we have reported Bias Index (BI) which estimates the difference in proportions of "Yes" for
147 the two raters, Prevalence Index which estimates the difference between the probability of
148 "Yes" and the probability of "No" , observed agreement (p_o); proportion of positive
149 agreement (p_{pos}) which estimates the conditional probability, given that one of the raters/
150 method, randomly selected, makes a positive rating, the other rater/ method will also do so;
151 proportion of negative agreement (p_{neg}) which estimates the conditional probability, given
152 that one of the raters/ method, randomly selected, makes a negative rating, the other rater/
153 method will also do so. We also calculated PABAK which is an index developed to account

154 for the effect that low prevalence and the difference in observer assessment of the frequency
 155 occurrence, have on kappa. All these statistics are derived from a 2x2 table as follows⁽¹⁹⁾.

		ABQ/ Rater 1	
		+	-
QM SA/ Rater 2	+	a	b
	-	c	d

156

157 $P_o=(a+d)/N$ where N denotes total sample size

158 $P_e=(((a+b)(a+c))/N)+(((c+d)(b+d))/N))/N$

159 $P_{pos}=2a/(2a+b+c)$

160 $P_{neg}=2d/(2d+b+c)$

161 $BI=(b-c)/N$

162 $PI=(a-d)/N$

163 $PABAK=2P_o-1$

164 We calculated the above mentioned statistics per a) subject level; where prevalent cases were
 165 defined as subjects having at least one vertebra fractured from T4 to L4 and controls as
 166 having none of the vertebrae from T4 to L4 fractured, and per b) vertebral level; we counted
 167 as cases any fracture from T4 to L4; furthermore we calculated agreements of the methods
 168 between cohorts, sexes, age categories and vertebral level. We used four age categories: ≥ 45
 169 and <60 ; ≥ 60 and <70 ; ≥ 70 and <80 ; ≥ 80 . We separated vertebral level into three categories:
 170 T4-T9, T10-T12 and L1-L4. Additionally we assessed differences in baseline characteristics
 171 between cases and non-cases defined by either method and also differences between
 172 concordant and discordant cases defined as follows: QMSA + ABQ-, QM SA- ABQ+, QM
 173 SA+ ABQ+ against the reference group QM SA- ABQ- . The future incident fracture
 174 prediction ability by prevalent vertebral fractures scored by either method was estimated

175 using a Cox regression model adjusted for Age, Sex, BMI, cohort effect and FN-BMD with a
176 mean follow up of 12 years. All analyses were performed using SPSS 21.0 (IBM Corp. NY,
177 USA).

178 **Results**

179 *Per subject analyses*

180 Radiographs were assessed for 7,582 participants of which 61.7% (n=4,672) were
181 from RS I, 21.8% (n=1,655) from RS II and 16.5% (n=1,255) from RS III. 60% of our study
182 participants were females and age ranged from 46 to 95 years (mean 65.3, Fig. 1). QM SA
183 scored vertebral fracture prevalence was 14.2% (95% CI: 13.4%-15.0%), compared to 4.0%
184 (95% CI: 3.6%-4.5%) scored by ABQ. Participants who had sustained a fracture were
185 significantly older according to both QM (67.4 vs. 64.9, p-value <0.001) and ABQ (70.4 vs.
186 65.1, p-value <0.001) compared to non-fractured participants. 54.5% of QM SA cases were
187 females vs. 45.5 % males (p-value <0.001) and 74.0 % of ABQ cases were females against
188 26% males (p-value<0.001). Both QM SA and ABQ fractured participants had lower FN-
189 BMD; 0.864 g/cm² vs. 0.890 g/cm² and 0.827 g/cm² vs. 0.894 g/cm², p-value <0.001
190 respectively. Fractured cases defined by ABQ were significantly shorter and lighter compared
191 to the healthy participants 163.5 vs. 167.5 and 72.6 kg vs. 75.4 kg (p-value <0.001). No
192 differences were seen between QM SA cases and controls in height and weight (p-value>
193 0.05) (Table 1a). When comparing (QM SA+) (ABQ-) participants vs. (QM SA-) (ABQ+),
194 the latter had lower FN-BMD (0.846 vs. 0.877, p-value<0.001), were lighter (74.1 vs. 76.9,
195 p-value<0.001), shorter (164.8 vs. 168.6) and comprised a higher number of females (74.3%
196 vs. 50.1%, p-value<0.001) (Table 1b). According to QM SA, the prevalence of vertebral
197 fractures was higher among males compared to females (16.0% vs. 13.0%), whereas
198 according to ABQ it was higher among females compared to males (5.0% vs. 2.6%) (Table
199 2). According to both methods the prevalence increased with increasing age (Table 3).

200 According to QM SA, 10% of the participants had only one spinal fracture, 2.6% had two
201 fractures, 1.0% had three and 0.5% more than three fractures, whereas according to ABQ the
202 estimates were lower with 2.9% of participants having only one fracture, 0.7% having two
203 fractures, 0.2% three and close to 0% more than three.

204 The estimated concordance between ABQ and QM SA was $\kappa = 0.24$. When assessing
205 agreement across sexes, it was significantly higher among females compared to males;
206 $\kappa = 0.31$ vs. $\kappa = 0.14$, $p\text{-value} < 0.001$ (Table 2). The agreement across age categories increased
207 with increasing age; the highest kappa was among those aged above 80 and was significantly
208 higher compared to the youngest group $\kappa = 0.40$ vs. 0.12 , $p\text{-value} < 0.001$ (Table 2).

209 Participants with a QM SA prevalent fracture had an increased risk for future non-vertebral
210 fractures compared to those with absent prevalent vertebral fracture (HR= 1.15, 95% CI
211 1.007; 1.32) and also an increased risk of future clinical vertebral fracture (HR= 2.70, 95% CI
212 2.18; 3.35) but not for incident hip fracture (HR= 1.49, 95% CI 0.92; 1.71). The same trend
213 was observed for participants with prevalent ABQ fractures although with higher estimates;
214 participants with prevalent ABQ fracture had an increased risk to sustain a future non-
215 vertebral fracture (HR= 1.30, 95% CI 1.06; 1.60), hip (HR= 1.47, 95% CI 1.05; 2.05) also an
216 increased risk of incident clinical fractures (HR= 5.27, 95% CI 4.00; 6.77) compared to those
217 with absent prevalent vertebral fracture (Fig 3).

218 *Per vertebral body analyses*

219 Among 7,582 participants, there were 1,574 (20.7%) vertebrae fractured according to
220 QM SA and 447 (5.8%) according to ABQ. Figure 2 shows the distribution of osteoporotic
221 vertebral fractures at each level assessed according to ABQ and QM SA. Both methods show
222 a bimodal distribution, but according to ABQ, most fractures were found at the thoraco-
223 lumbar junction (T12-L1) region, whereas according to QM SA, most deformities were at the
224 middle (T7-T8) and lower thoracic regions (T11-T12), showing a more prominent bimodal

225 pattern (Fig. 2). The frequencies for QM SA deformities' classification of severity was 49.2%
226 mild, 30.8% moderate and 4.7% severe; 53.5% of the deformities were wedge shaped, 11.9%
227 were biconcave and 19.3% were crush (supplementary Table 1 and supplementary Figure 2).
228 The agreement statistics per vertebral level could not be calculated for T4 since according to
229 ABQ there were no T4 vertebrae fractured in any of the participants. The kappa statistic in
230 the other vertebrae varied from 0.04 at T5 to 0.40 at L1. When assessing the agreement per
231 region of the spine the highest agreement was in the L1-L4 region $\kappa=0.37$ (p-value<0.001)
232 and when further stratifying by sex it reached $\kappa=0.41$ (p-value<0.001) among females (Table
233 4).

234 *Excluding mild fractures from the study*

235 We observed an increase in the net agreement between methods, mostly because the
236 deformities with height loss but intact endplates were excluded. Out of 1,075 participants that
237 were classified as fractured by QM SA, 614 of them had mild fractures. When excluding
238 these subjects from the analysis, according to QM SA the prevalence decreased from 14.1%
239 to 6.6%. Excluding these participants slightly affected the prevalence of ABQ scored
240 fractures with a decrease from 4.0% to 3.8%. On the other hand the kappa statistic increased
241 from 0.24 to 0.40 (p-value<0.001) and reached its maximum among participants aged above
242 80, $\kappa=0.47$ among females $\kappa=0.48$ and at the L1 level $\kappa=0.53$ (Table 5). The prevalence of
243 fractured vertebrae by grading of QM SA deformities is displayed by vertebral level
244 distribution in Figure 4. According to QM SA, the highest concentration of fractured
245 vertebrae was at T7-T8 and T11-T12-L1, showing again a bimodal distribution with almost
246 the same number of fractured vertebrae for both peaks. A bimodal distribution was observed
247 for ABQ as well, but with the highest peak at T12-L1.

248 **Discussion**

249 In this large population based study where we compared two assessment methods,
250 osteoporotic vertebral fracture prevalence was four times higher when applying
251 SpineAnalyzer[®] software-assisted QM compared to ABQ. Each method classified a
252 considerable number of deformities that were assessed as normal by the other, reflected by
253 poor between-method agreement statistics. Our study is the first to compare SpineAnalyzer[®]
254 software-assisted QM and ABQ. According to ABQ, vertebral fracture prevalence was higher
255 among females than males, whereas according to QM SA prevalence was higher among
256 males. Differences in baseline characteristics were also observed; the difference in age,
257 height, weight, FN-BMD and over-representation of females among cases compared to
258 controls were stronger when they were defined by ABQ then when they were defined by QM
259 SA. Also differences in BMD levels were observed among participants with discordant
260 assessment of vertebral fractures, where participants with (ABQ+) (QM SA-) deformities had
261 lower FN-BMD, weight and height compared to participants with (QM SA+) (ABQ-)
262 deformities. We also observed difference in the ability to predict future non-vertebral and
263 clinical vertebral fracture by prevalent vertebral fractures scored by either method with ABQ
264 being more strongly associated with future fractures. The vertebral fracture prevalence
265 estimate in our population for the ABQ method is similar to previous findings in other
266 populations^(13,20) mostly consisting of elderly females in a clinical setting; and also taking
267 into account that we included subjects of both genders and even a subset comprising a
268 relatively young population (RS-III). In previous work of the Rotterdam Study⁽²¹⁾, including
269 a sample of RS-I subjects assessed with the McCloskey-Kanis method⁽²²⁾, the prevalence was
270 found to be 6.3%. This prevalence is intermediate between the prevalence of ABQ (~4.0%)
271 and QMSA (~14.1%) and very similar to the prevalence of QM SA after excluding Grade 1
272 (~6.6%). The agreement was significantly higher in females compared to males, L1-L4 level
273 and older age. The bimodal fracture distribution over the vertebral column was obvious for

274 the QM SA method in our cohort, with maxima at the mid-thoracic and lower thoracic
275 regions including the thoraco-lumbar junction and less pronounced in ABQ. This pattern has
276 been reported previously using other assessment methods. However, some argue that the
277 more pronounced mid-thoracic peak with QM is to a great extent due to degenerative
278 changes, normal anatomical variation (i.e., short vertebral height) and old traumatic fractures
279 ⁽²³⁾. It has been put forward that ABQ would be able to differentiate these entities⁽¹⁵⁾
280 compatible with our findings (Fig 3). When assessing QM SA morphometry, the far majority
281 of deformities were classified as mild wedges located mostly at the T7-T8 level. By
282 excluding QM SA-mild deformities, the difference in prevalence between the methods
283 decreased and all agreement statistics increased.

284 We have assessed vertebral levels T4 to L4, as T1-T3 has poor evaluability and L5 is
285 usually not affected by osteoporotic fractures. Several studies have compared assessment
286 methods, but only a few have evaluated SpineAnalyzer[®] software or ABQ, and none have
287 directly compared these two methods. SpineAnalyzer[®] software-assisted QM reading by a
288 non-radiologist has been found to agree relatively well with conventional semi-quantitative
289 (SQ) grading, i.e., visual estimation of vertebral body heights performed by experienced
290 radiologists, with a kappa for agreement of 0.78.⁽²⁴⁾ ABQ comparisons with QM (Eastell-
291 Melton and McCloskey definitions) have yielded kappa statistics between 0.39 and 0.64.⁽¹³⁾
292 Most notably, the lowest agreement found to date is between ABQ and Genant's SQ
293 methods, observing kappa statistics of 0.30 to 0.58.^(15,25,26) The agreement between
294 SpineAnalyzer software-assisted QM and ABQ in this study was even lower than the
295 agreement between ABQ and Genant's SQ methods. This could have been further amplified
296 because we have examined a relatively young and generally healthy population in RSIII, in
297 which there might be many mild non-fracture deformities. This is also sustained by the
298 results where kappa tended to increase with the increase of age. The kappa statistic is

299 associated with two paradoxes described by Feinstein et al.^(18,19) These paradoxes arise from
300 the chance-adjustment applied to kappa; adjustment that also helps to “standardize” and allow
301 comparison across different studies. Kappa is estimated as the difference between *observed*
302 and *expected agreement* divided by [$1 - \text{expected agreement}$]. Indeed in our study we observe
303 a tendency towards *Paradox 1*, where there is high *expected agreement* (p_e) as well as high
304 *observed agreement* which still results in a low kappa (Table 2). In addition, *Paradox 2* is
305 also present given the population-based setting of our study, resulting in a large number of
306 individuals without events, which creates unbalance of the marginal totals reflected in a high
307 PI. The marginal totals are already determined by the (relatively low) prevalence of VFs and
308 (healthy) population we studied and they can explain only partly the low kappa values. The
309 remaining explanation of low kappa will arise from the method’s separate performances for
310 P_{pos} and P_{neg} . While kappa helps to compare agreement across studies, positive and negative
311 agreement statistics help to better understand the individual study. In the present study, QM
312 SA and ABQ agreed excellently to identify controls, but poorly to identify cases. Having said
313 this and given that vertebral fracture diagnosis requires adaptation of current approaches to
314 conciliate the differences between methods, we propose that one way would be by re-
315 examining QM mild deformities for endplate depression. We simulated in our data a
316 redistribution of the 2x2 table when reconsidering mild QM fractures for endplate depression
317 and we saw that all agreement statistics increase significantly (supplementary Table 2c).

318 Nonetheless, it should be noted that agreement statistics concern precision of a study
319 and may not necessarily relate to its validity. QM SA would not diagnose vertebral fractures
320 in the case of endplate depression without reduced vertebral height, and conversely, ABQ
321 would not diagnose a QM SA -based vertebral deformity with reduced height but intact
322 endplates. More research is needed to clarify which of these discordant cases are clinically
323 relevant vertebral fractures and which are false-positives.

324 It is important to recognize that although Spine Analyzer[®] software uses the Genant height
325 criteria to judge severity of deformities defined by QM, QM methods on Spine Analyzer[®]
326 software are *not* the same as the Genant semi-quantitative method⁽¹²⁾. While the Genant SQ
327 method⁽¹²⁾ unlike ABQ, does not specifically state how to differentiate non-fracture
328 deformities from true fractures, it relies on the expertise of the evaluator⁽²⁷⁾ to discriminate
329 them from vertebral height loss due to other causes such as degenerative remodeling and
330 Scheuermann's disease⁽²⁸⁾. In an accompanying article in this issue, Lentle et al.⁽²⁹⁾ employed
331 the standard Genant methodology and draw similar conclusions with regard to the drastic
332 differences in fracture prevalence and low concordance with a modified ABQ methodology.

333 Our overall aim was to objectively compare radiological assessment methods for
334 osteoporotic vertebral fractures. Strengths of our study are that we systematically applied two
335 very different assessment methods by two independent teams of trained readers which
336 eliminates the risk of ascertainment bias. Applying two methods in a very large setting with
337 two independent teams, proved to be very labor-intensive, requiring extra consensus
338 meetings, supervision by musculoskeletal radiologists and double readings. Although
339 radiographs were assessed by well-trained reader teams, it was not feasible to have all
340 radiographs assessed by musculoskeletal radiologists. We are aware that more subtle endplate
341 depression fractures could have been missed. As the Rotterdam Study is deemed
342 representative of the general Dutch middle-aged to elderly population, we believe that our
343 results may be extrapolated to other settings as well.

344 The semi-automated SpineAnalyzer[®] software-assisted QM method proved to be an
345 excellent recording tool for research purposes, providing a standardized data output.⁽³⁰⁾
346 Surprisingly, ABQ was in our experience even more time-efficient, but this method requires
347 more intensive initial training. Quantitative assessment is based on morphometry alone,
348 which may result in the inclusion of deformities that are not truly vertebral fractures. For this

349 reason it might be better to refer to “deformities” instead of “fractures” for cases defined by
350 QM. Yet, we experienced that further triage for both methods requires a lot of extra effort
351 involving extra double-reading of up to thousands of participants. Further standardization and
352 automation of this triage procedure with clear-cut classification criteria would be very
353 helpful.

354 Vertebral fractures are often a first presentation of osteoporosis and should be
355 regarded as an opportunity to trace individuals at high-risk for additional fractures and other
356 related adverse health outcomes. To accomplish this, accurate vertebral fracture diagnosis is
357 needed to identify these patients at high risk, as many effective treatment options are
358 available. Conversely, individuals without true vertebral fractures should not be unnecessarily
359 treated with medication, which is associated with unnecessary costs and potential adverse
360 effects.⁽³¹⁾ Improvement of radiological vertebral fracture definition, clearer criteria for non-
361 fracture deformities differential diagnosis⁽³²⁾ and more wide-spread and consistent application
362 of an optimal method may improve clinical care.

363 We have undertaken meticulous phenotyping on our ABQ and SpineAnalyzer®
364 morphometric raw data. With these data, different cut-offs and vertebral fracture definitions
365 could be linked to various clinically relevant outcomes. Furthermore, the remaining
366 Rotterdam Study cohorts, which in total will yield ~11,000 subjects aged 45 years and over,
367 will be assessed for the presence of osteoporotic vertebral fractures. In addition, our
368 measurements could serve as population reference data.

369 In conclusion, we procured an impartial comparison of osteoporotic vertebral fracture
370 assessment methods in the large population-based Rotterdam Study, with extensive recording
371 of vertebral fracture distribution according to sex, age, deformity shape, severity and location.
372 Osteoporotic vertebral fracture prevalence is significantly different when applying either
373 software-assisted QM or ABQ. Further work is needed to reveal which of the discordant

374 cases are actually clinically relevant true vertebral fractures and which are not. We propose
375 that mild deformities should be assessed for endplate depression, decreasing this way the
376 false-positive QM fractures and conciliating the two methods.

377 **Acknowledgements**

378 We would like to thank Dr. Guirong Jiang for the training in the algorithm-based qualitative
379 assessment method (ABQ). We are thankful to the employees from Optasia Medical Ltd who
380 familiarized us with the use of the SpineAnalyzer[®] software. The Rotterdam Study is funded
381 by Erasmus Medical Center and Erasmus University, Rotterdam, Netherlands Organization
382 for the Health Research and Development (ZonMw), the Research Institute for Diseases in
383 the Elderly (014-93-015; RIDE2), RIDE), the Ministry of Education, Culture and Science, the
384 Ministry for Health, Welfare and Sports, the European Commission (DG XII), and the
385 Municipality of Rotterdam. The authors are grateful to the study participants, the staff from
386 the Rotterdam Study (particularly Lydia Buist and Hannie van den Boogert for acquisition of
387 the radiographs) and the participating general practitioners and pharmacists. We thank René
388 Vermeren. Nano Suwarno and Mart Rentmeester for their technical support. Last but not
389 least, we acknowledge the tremendous efforts from our team of radiographic readers.

390

391 **Authors' contributions**

392 LO, FK, SJB, MAI, EHGO, and FR designed the study. LO, FK, SJB, JBJvM, JHW, FJAvR
393 collected and processed the data. LO, FK, SJB, JTS, EMC, JHW, FJAvR assessed and
394 (statistically) analyzed the data. LO, FK, SJB, JTS, EMC, JBJvM, MCZ, GPK, EHGO, FR
395 interpreted the results. LO, FK, SJB created the figures and tables. LO, FK, SJB, EHGO, FR
396 drafted the manuscript. All authors (LO, FK, SJB, JTS, EMC, JBJvM, MAI, JHW, FJAvR,
397 MCZ, GPK, EHGO, FR) read and revised the manuscript, and approved the final submitted
398 version. LO and FK, EHGO and FR contributed equally. EHGO and FR assume

399 responsibility for the completeness and accuracy of the data and analyses, and for adherence
400 to the study protocol.

401

402 **References**

- 403 1. Szulc P, Bouxsein ML. Overview of osteoporosis: Epidemiology and clinical
404 management. Vertebral Fracture Initiative Resource Document. 2011;PART I:1-65.
- 405 2. Cooper C. Epidemiology and public health impact of osteoporosis. *Baillieres Clin*
406 *Rheumatol.* 1993;7(3):459-77.
- 407 3. Ström O, Borgström F, Kanis JA, et al. Osteoporosis: burden, health care provision
408 and opportunities in the EU. *Arch Osteoporos.* 2011;DOI 10.1007/s11657-011-0060-
409 1.
- 410 4. Burge R, Dawson-Hughes B, Solomon DH, Wong JB, King A, Tosteson A. Incidence
411 and economic burden of osteoporosis-related fractures in the United States, 2005-
412 2025. *J Bone Miner Res.* 2007;22(3):465-75.
- 413 5. Nevitt MC, Ettinger B, Black DM, et al. The association of radiographically detected
414 vertebral fractures with back pain and function: a prospective study. *Ann Intern Med.*
415 1998;128(10):793-800.
- 416 6. Bliuc D, Nguyen ND, Milch VE, Nguyen TV, Eisman JA, Center JR. Mortality risk
417 associated with low-trauma osteoporotic fracture and subsequent fracture in men and
418 women. *Jama.* 2009;301(5):513-21.
- 419 7. Klotzbuecher CM, Ross PD, Landsman PB, Abbott TA, 3rd, Berger M. Patients with
420 prior fractures have an increased risk of future fractures: a summary of the literature
421 and statistical synthesis. *J Bone Miner Res.* 2000;15(4):721-39.
- 422 8. Burger H, van Daele PL, Algra D, et al. Vertebral deformities as predictors of non-
423 vertebral fractures. *Bmj.* 1994;309(6960):991-2.
- 424 9. Black DM, Arden NK, Palermo L, Pearson J, Cummings SR. Prevalent vertebral
425 deformities predict hip fractures and new vertebral deformities but not wrist fractures.
426 Study of Osteoporotic Fractures Research Group. *J Bone Miner Res.* 1999;14(5):821-
427 8.
- 428 10. Oei L, Rivadeneira F, Ly F, et al. Review of radiological scoring methods of
429 osteoporotic vertebral fractures for clinical and research settings. *Eur Radiol.*
430 2013;23(2):476-86.
- 431 11. Brett A, Miller CG, Hayes CW, et al. Development of a clinical workflow tool to
432 enhance the detection of vertebral fractures: accuracy and precision evaluation. *Spine*
433 *(Phila Pa 1976).* 2009;34(22):2437-43.
- 434 12. Genant HK, Wu CY, van Kuijk C, Nevitt MC. Vertebral fracture assessment using a
435 semiquantitative technique. *J Bone Miner Res.* 1993;8(9):1137-48.
- 436 13. Jiang G, Eastell R, Barrington NA, Ferrar L. Comparison of methods for the visual
437 identification of prevalent vertebral fracture in osteoporosis. *Osteoporos Int.*
438 2004;15(11):887-96.
- 439 14. Ferrar L, Jiang G, Adams J, Eastell R. Identification of vertebral fractures: an update.
440 *Osteoporos Int.* 2005;16(7):717-28.
- 441 15. Ferrar L, Jiang G, Cawthon PM, et al. Identification of vertebral fracture and non-
442 osteoporotic short vertebral height in men: the MrOS study. *J Bone Miner Res.*
443 2007;22(9):1434-41.
- 444 16. Hofman A, Brusselle GG, Darwish Murad S, et al. The Rotterdam Study: 2016
445 objectives and design update. *Eur J Epidemiol.* 2015;30(8):661-708.
- 446 17. Cohen J. A Coefficient of Agreement for Nominal Scales. *Educational and*
447 *Psychological Measurement.* 1960;20(1):37-46.
- 448 18. Feinstein AR, Cicchetti DV. High agreement but low kappa: I. The problems of two
449 paradoxes. *J Clin Epidemiol.* 1990;43(6):543-9.

- 450 19. Cicchetti DV, Feinstein AR. High agreement but low kappa: II. Resolving the
451 paradoxes. *J Clin Epidemiol*. 1990;43(6):551-8.
- 452 20. Ferrar L, Roux C, Felsenberg D, Gluer CC, Eastell R. Association between incident
453 and baseline vertebral fractures in European women: vertebral fracture assessment in
454 the Osteoporosis and Ultrasound Study (OPUS). *Osteoporos Int*. 2012;23(1):59-65.
- 455 21. Van der Klift M, De Laet CE, McCloskey EV, Hofman A, Pols HA. The incidence of
456 vertebral fractures in men and women: the Rotterdam Study. *J Bone Miner Res*.
457 2002;17(6):1051-6.
- 458 22. McCloskey EV, Spector TD, Eyres KS, et al. The assessment of vertebral deformity:
459 a method for use in population studies and clinical trials. *Osteoporos Int*.
460 1993;3(3):138-47.
- 461 23. Adams JE, Lenchik L, Roux C, Genant HK. Radiological assessment of vertebral
462 fracture. *Vertebral Fracture Initiative Resource Document*. 2011;PART II:1-48.
- 463 24. Kim YM, Demissie S, Genant HK, et al. Identification of prevalent vertebral fractures
464 using CT lateral scout views: a comparison of semi-automated quantitative vertebral
465 morphometry and radiologist semi-quantitative grading. *Osteoporos Int*.
466 2012;23(3):1007-16.
- 467 25. Ferrar L, Jiang G, Clowes JA, Peel NF, Eastell R. Comparison of densitometric and
468 radiographic vertebral fracture assessment using the algorithm-based qualitative
469 (ABQ) method in postmenopausal women at low and high risk of fracture. *J Bone
470 Miner Res*. 2008;23(1):103-11.
- 471 26. Ferrar L, Jiang G, Schousboe JT, DeBold CR, Eastell R. Algorithm-based qualitative
472 and semiquantitative identification of prevalent vertebral fracture: agreement between
473 different readers, imaging modalities, and diagnostic approaches. *J Bone Miner Res*.
474 2008;23(3):417-24.
- 475 27. Grados F, Fechtenbaum J, Flipon E, Kolta S, Roux C, Fardellone P. Radiographic
476 methods for evaluating osteoporotic vertebral fractures. *Joint Bone Spine*.
477 2009;76(3):241-7.
- 478 28. Armbrecht G, Felsenberg D, Ganswindt M, et al. Vertebral Scheuermann's disease in
479 Europe: prevalence, geographic variation and radiological correlates in men and
480 women aged 50 and over. *Osteoporos Int*. 2015;26(10):2509-19.
- 481 29. Lentle BC, Berger C, Probyn L, et al. Comparative Analysis of the Radiology of
482 Osteoporotic Vertebral Fractures in Women and Men: Cross-Sectional and
483 Longitudinal Observations from the Canadian Multicentre Osteoporosis Study
484 (CaMos). In Press. 2017.
- 485 30. Oei L, Ly F, El Saddy S, et al. Multi-functionality of computer-aided quantitative
486 vertebral fracture morphometry analyses. *Quant Imaging Med Surg*. 2013;3(5):249-
487 55.
- 488 31. Breda SJ, Oei L, Oei EH, Zillikens MC. [Osteoporotic vertebral fractures or
489 Scheuermann's disease?]. *Ned Tijdschr Geneesk*. 2013;157(45):A6479.
- 490 32. Makurthou AA, Oei L, Saddy SE, et al. Scheuermann's Disease: Evaluation of
491 Radiological Criteria and Population Prevalence. *Spine (Phila Pa 1976)*. 2013.
492

493
494 **Tables**

495
496 **Table 1a. Baseline characteristics of study participants shown by vertebral fracture**
497 **status as scored by each definition.** Fractured participants according to both QM SA and
498 ABQ were significantly older, had lower FN-BMD and an over-representation of females.
499 According to ABQ they were also shorter and lighter. Among QM SA cases, 57% were
500 classified as grade 1, 37 % as grade 2 and 6% grade 3. Among ABQ defined cases, 39 were
501 also scored as grade 1 by QM SA, 111 as grade 2 and 49 as grade 3.
502

		QM SA		ABQ	
	Overall N=7,582	Controls n=6,506	Cases n=1,076	Controls n=7,278	Cases n=304
Age	65.3 (8.8)	64.9 (8.6)	67.4 (9.7)	65.1 (8.7)	70.4 (9.9)
Sex (female)	4,516 (59.6)	3,930 (60.4)	586 (54.5)	4,291 (59.0)	225 (74.0)
Height	167.4 (9.1)	167.4 (9.0)	167.5 (9.3)	167.6 (9.0)	163.5 (8.5)
Weight	75.3 (12.9)	75.2 (12.8)	76.0 (13.8)	75.4 (12.9)	72.6 (13.4)
BMI	26.8 (3.9)	26.8 (3.9)	27.0 (4.1)	26.8 (3.9)	27.1 (4.3)
FN- BMD*	0.890 (0.14)	0.895 (0.14)	0.864 (0.14)	0.894 (0.14)	0.827 (0.14)
QM SA Grade					
1			614 (57.0)		39
2			399 (37.0)		111
3			63 (6.0)		49

503 *adjusted for age, sex, height, weight

504
 505
 506
 507
 508
 509
 510
 511
 512

Table 1b. Baseline characteristics among participants with discordant and concordant assessment of vertebral fractures. Participants classified as cases according to QM but not according to ABQ were used as reference group for comparisons. Participants classified as cases according to ABQ but not to QM, were lighter, shorter , had lower FN-BMD and a higher representation of females.

N=7,582	(QM SA-) (ABQ-) (ref)	(QM SA+) (ABQ-)	(QM SA-) (ABQ+)	(QM SA+) (ABQ+)	(QM SA G2 or G3+) (ABQ+)
	N=(6,401)	N=(877)	N=(105)	N= (199)	N= (160)
Age	64.9 (8.5)	66.4 (9.4)	67.6 (10.1)	71.9 (9.5)	72.4 (9.4)
Sex (female)	3852 (60.2)	439 (50.1)	78 (74.3)	143 (73.9)	121 (75.6)
Height	167.4 (9.0)	168.6 (9.1)	164.8 (8.0)	162.8 (8.7)	161.9 (8.4)
Weight	75.27 (12.8)	76.9 (13.7)	74.13 (13.2)	71.8 (13.5)	71.1 (13.0)
BMI	26.8 (3.9)	27.0 (4.1)	27.2 (4.4)	27.0 (4.2)	27.0 (4.2)
FN- BMD*	0.896 (0.14)	0.877 (0.14)	0.846 (0.14)	0.820 (0.14)	0.763 (0.14)
QM SA Grade					
1		575		39	
2		288		111	111
3		14		49	49

513 **Table 2. Participants with prevalent vertebral fractures and agreement statistics**
514 **between QM SA and ABQ, stratified by cohort and sex.** The prevalence of VFs is the
515 highest in RS III according to both QM SA and ABQ. The agreement statistics are the highest
516 in RS I. According to ABQ, the prevalence of VFs is higher among females but not according
517 to QM SA
518

	Cohort			Sex		Pooled (N=7,582)
	RS I (N=4,672)	RS II (N=1,655)	RS III (N=1,255)	Males (N=3,066)	Females (N=4,516)	
QM SA (%)	578 (12.4)	249 (15.0)	249 (19.8)	490 (16.0)	586 (12.9)	1076 (14.1)
ABQ (%)	190 (4.1)	59 (3.6)	55 (4.4)	79 (2.6)	225 (5.0)	304 (4.0)
Kappa	0.28	0.20	0.16	0.14	0.31	0.24
Observed agreement	0.89	0.86	0.81	0.85	0.89	0.87
Expected Agreement	0.85	0.82	0.77	0.82	0.83	0.83
Bias Index	0.08	0.11	0.15	0.13	0.08	0.10
Prevalence Index	-0.83	-0.81	-0.75	-0.81	-0.82	-0.81
Positive agreement	0.33	0.25	0.22	0.18	0.36	0.29
Negative agreement	0.94	0.92	0.89	0.91	0.94	0.93
PABAK	0.78	0.72	0.62	0.70	0.78	0.74

519
520

521 **Table 3. Participants with prevalent vertebral fractures and agreement statistics**
 522 **between QM SA and ABQ stratified by age categories.** The prevalence increases as age
 523 increases according to both methods. The highest prevalence is , as expected, among
 524 participants above 80 years old and kappa statistic is the highest in the same category.
 525
 526

	Age category			
	45-60 (N=2,396)	60-70 (N=2,932)	70 -80 (N=1,745)	>80 (N=509)
QM SA (%)	269 (11.2)	375 (12.8)	315 (18.1)	117 (23.0)
ABQ (%)	53 (2.2)	85 (2.9)	113 (6.5)	53 (10.4)
Kappa	0.12	0.20	0.30	0.40
Observed agreement	0.89	0.88	0.84	0.83
Expected agreement	0.87	0.85	0.77	0.71
Bias Index	0.09	0.10	0.11	0.12
Prevalence Index	-0.86	-0.84	-0.75	-0.66
Positive agreement	0.15	0.23	0.37	0.48
Negative agreement	0.94	0.93	0.91	0.90
PABAK	0.78	0.76	0.68	0.66

527

Table 4. Agreement statistics regarding number of fractured vertebrae by regions in the spine and by sex; note is per vertebral level. The lower in the spine is the fracture located, the higher is the agreement between methods.

	Spine Level								
	T4-T9			T10-T12			L1-L4		
	Males	Females	Pooled	Males	Females	Pooled	Males	Females	Pooled
N=7,582									
Males									
n=3,066									
Females									
n=4,516									
QM (%)	335(10.9)	339(7.5)	674(8.9)	156(5.1)	187(4.1)	343(4.5)	87(2.8)	129(2.9)	216(2.8)
ABQ (%)	29 (0.9)	51 (1.1)	80 (1.1)	24(0.8)	92 (2.0)	116(1.5)	43(1.4)	125(2.8)	168(2.2)
Kappa	0.10	0.17	0.14	0.14	0.39	0.29	0.28	0.41	0.37
Observed agreement	0.90	0.93	0.92	0.95	0.97	0.96	0.97	0.97	0.97
Expected Agreement	0.88	0.91	0.90	0.94	0.94	0.94	0.96	0.94	0.95
Bias Index	0.09	0.06	0.07	0.04	0.02	0.03	0.01	0.00	0.006
Prevalence Index	-0.88	-0.91	-0.90	-0.94	-0.94	-0.94	0.96	-0.94	-0.95
Positive agreement	0.12	0.18	0.15	0.16	0.40	0.31	0.29	0.43	0.38
Negative agreement	0.94	0.96	0.96	0.97	0.98	0.98	0.98	0.98	0.98
PABAK	0.80	0.86	0.84	0.90	0.92	0.92	0.94	0.94	0.94

Table 5. Agreement statistics regarding fractured subjects after excluding from the study those who had a mild fracture. After excluding participant with mild fractures from the study, all agreement statistics increase and the difference in prevalence between QM and ABQ decreases.

	Age Category				Sex		Pooled
	45-60 (N=2,217)	60-70 (N=2,698)	70 -80 (N=1,590)	>80 (N=463)	Males (N=2,768)	Females (N=4,200)	
QM SA (%)	90 (4.0)	141 (5.2)	160 (10.0)	71 (15.3)	192 (6.9)	270 (11.2)	462 (6.6)
ABQ (%)	46 (2.0)	71 (2.6)	101 (6.3)	47 (10.1)	66 (2.4)	199 (4.7)	265 (3.8)
Kappa	0.25	0.35	0.47	0.53	0.28	0.49	0.41
Observed agreement	0.95	0.95	0.92	0.90	0.93	0.95	0.94
Expected Agreement	0.94	0.92	0.85	0.78	0.91	0.89	0.90
Bias Index	0.02	0.03	0.04	0.05	0.04	0.02	0.03
Prevalence Index	-0.94	-0.92	-0.83	-0.74	-0.90	-0.89	-0.89
Positive agreement	0.26	0.38	0.51	0.60	0.30	0.52	0.44
Negative agreement	0.98	0.97	0.96	0.94	0.97	0.97	0.97
PABAK	0.90	0.90	0.84	0.80	0.86	0.90	0.88

Figures

Fig. 1. Age at baseline distribution within the Rotterdam Study population, stratified by sex and cohort. RS III is the youngest cohort and RS I the oldest. Mean age among both sexes is 65.1 years but the study population is made up by approximately 60% females and 40% males.

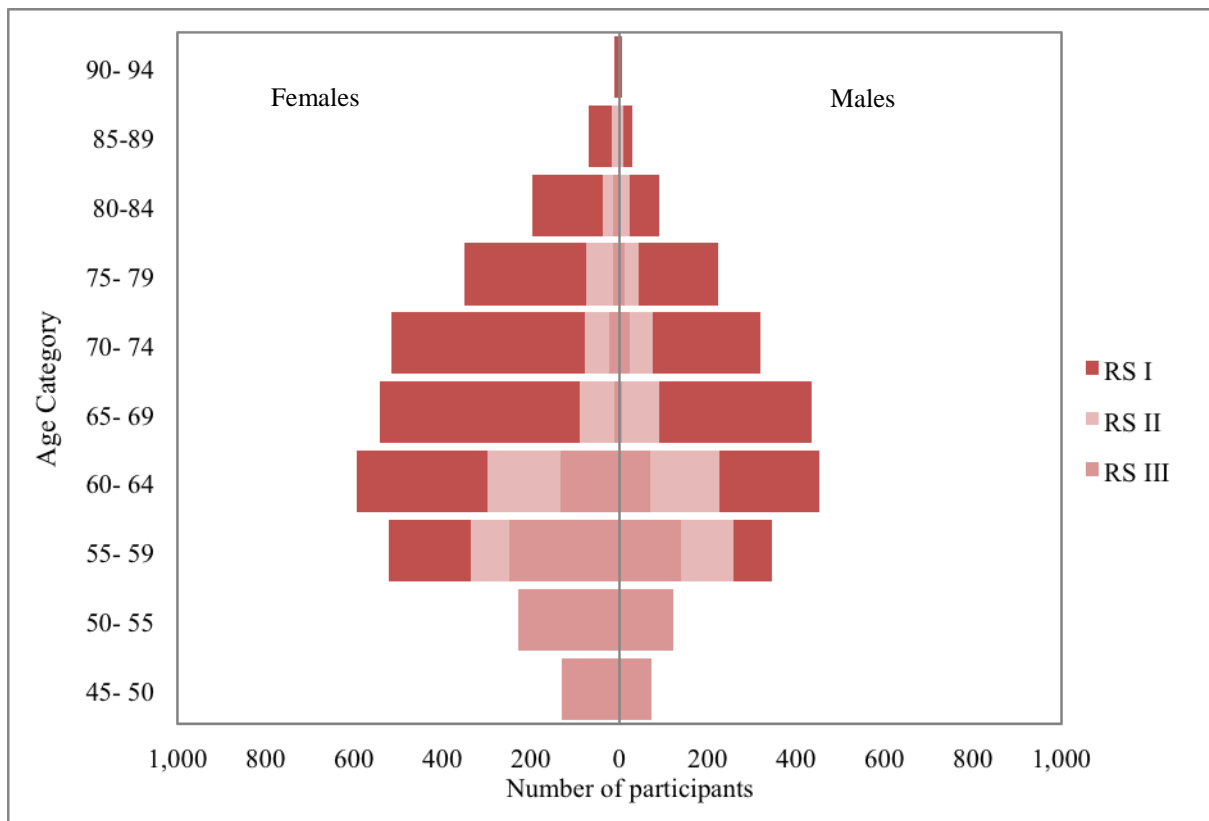


Fig. 2. Distribution of osteoporotic vertebral fractures across the thoracic and lumbar spine assessed according to the algorithm-based qualitative (ABQ) method and quantitative morphometry (QM) performed by SpineAnalyzer[®] software-assisted quantitative morphometry (vertebral height loss $\geq 20\%$). For both methods a bi-modal distribution can be seen but it is more pronounced for QM. According to QM the peaks are located at T7-T8 and T11-T12, whereas according to ABQ the highest peak is at T12-L1 and second highest at T7-T8.

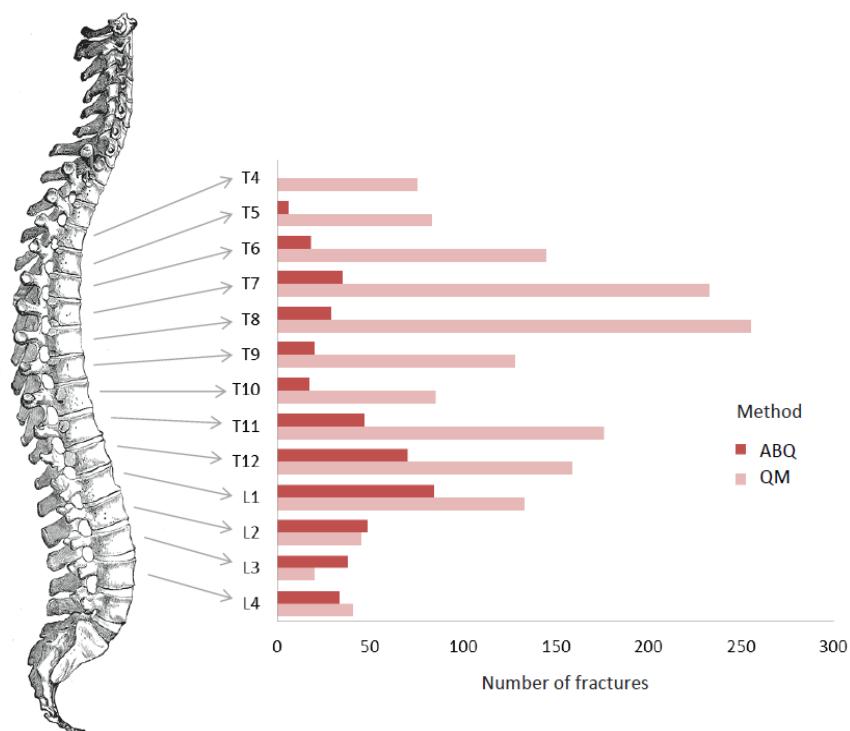
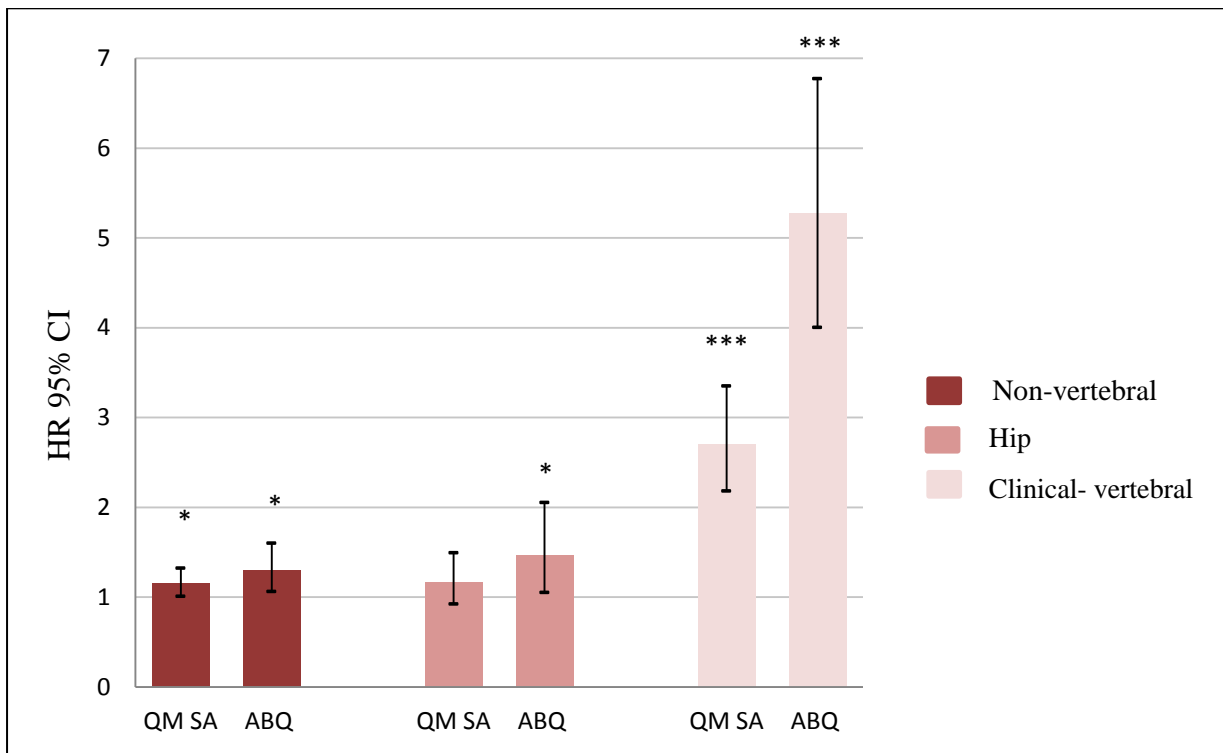


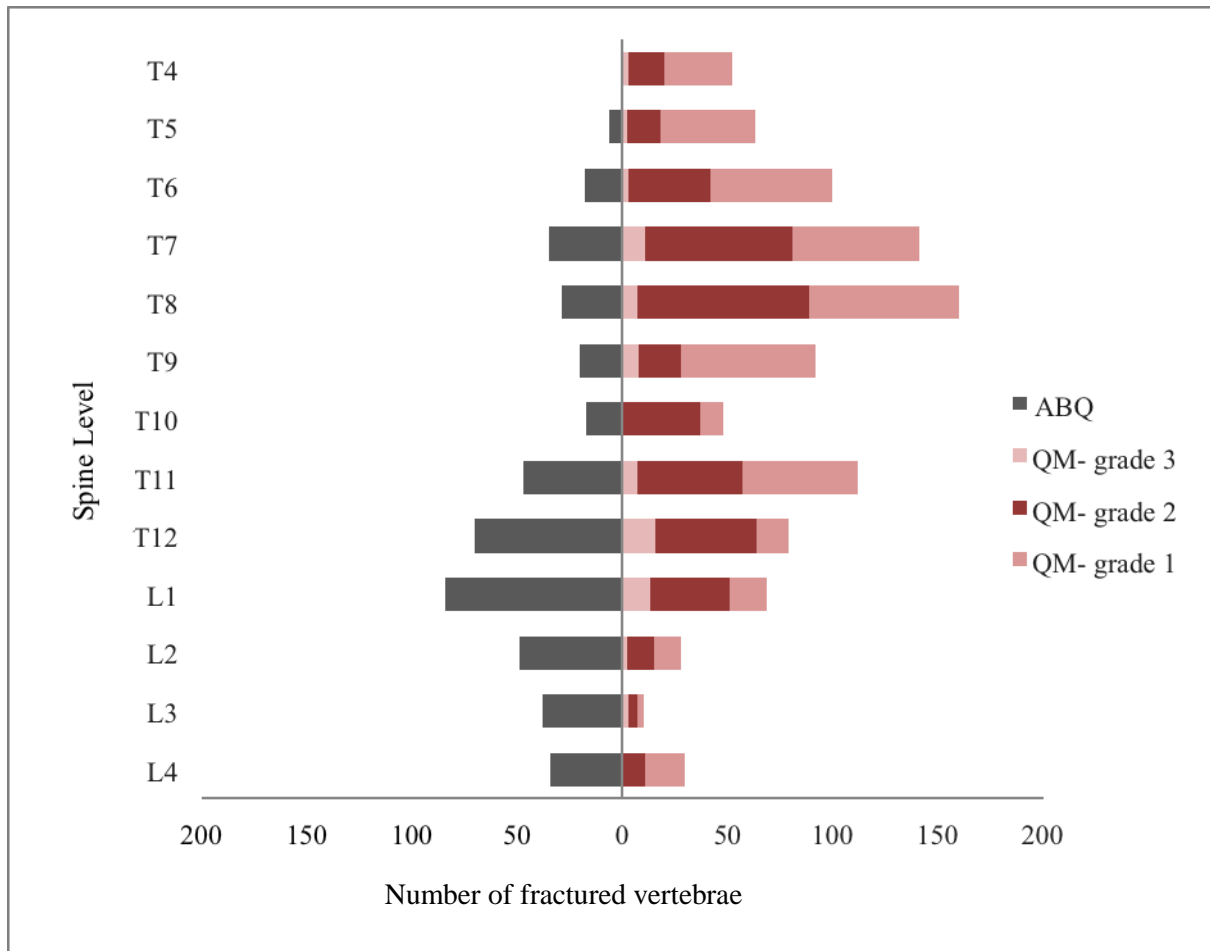
Fig. 3. The association between prevalent vertebral fractures scored by either method and incident non-vertebral and clinical vertebral fractures. During a mean follow-up time of 12 years, the 7,582 participants of this study sustained 1700 new non-vertebral fractures, 459 hip and 444 clinical-vertebral fractures. Participants with either prevalent QM or prevalent ABQ had increased risk of incident non-vertebral or clinical-vertebral fractures compared to participants who had not sustained either a QM or ABQ (respectively) fracture at baseline. Participants with an ABQ prevalent vertebral fracture at baseline were slightly more strongly associated with future non-vertebral fractures and significantly more strongly associated with incident clinical vertebral fractures compared to QM SA.



*p-value < 0.05

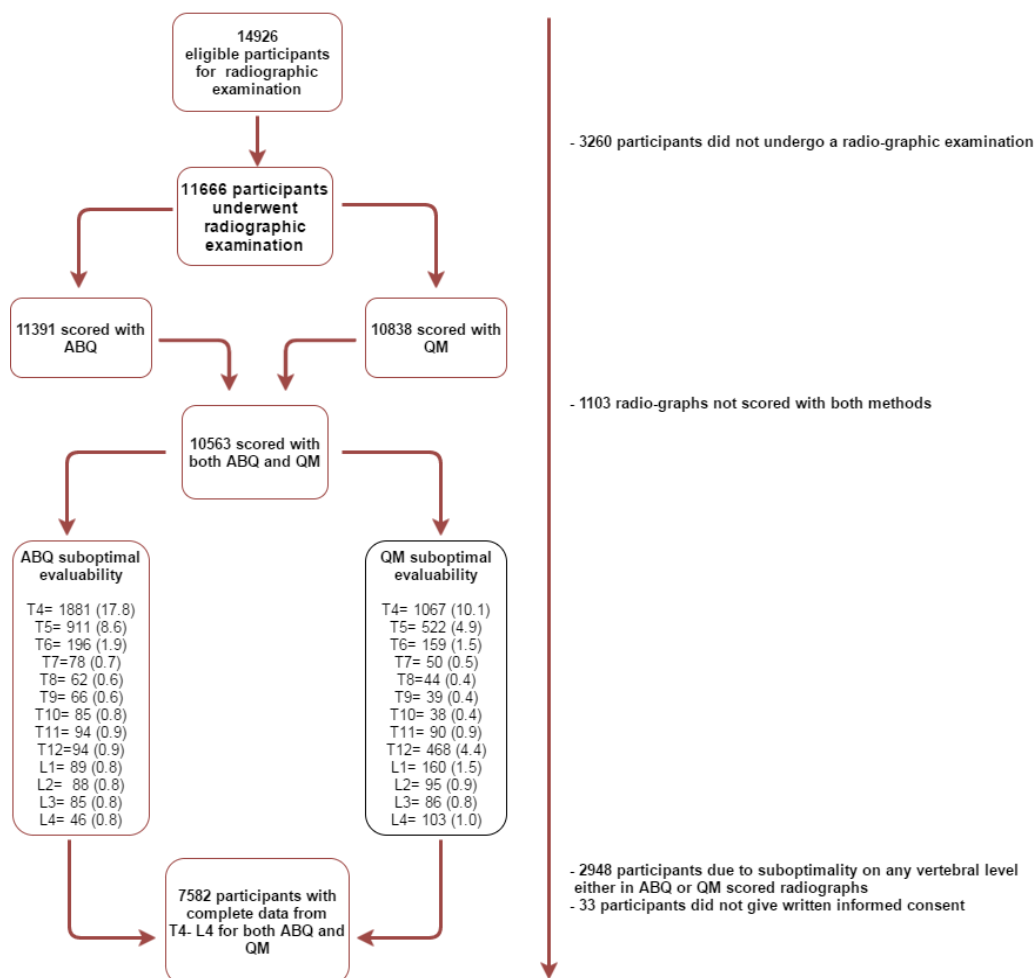
***p-value < 0.001

Fig. 4. Distribution of osteoporotic vertebral fractures per vertebral level assessed with the algorithm-based qualitative (ABQ) method and quantitative morphometry (QM) performed by SpineAnalyzer[®] software-assisted quantitative morphometry. Mild deformities constitute around 62% of QM vertebral fractures, followed by grade two , 33% and the least common, grade three with 5%



Supplementary Figures and tables

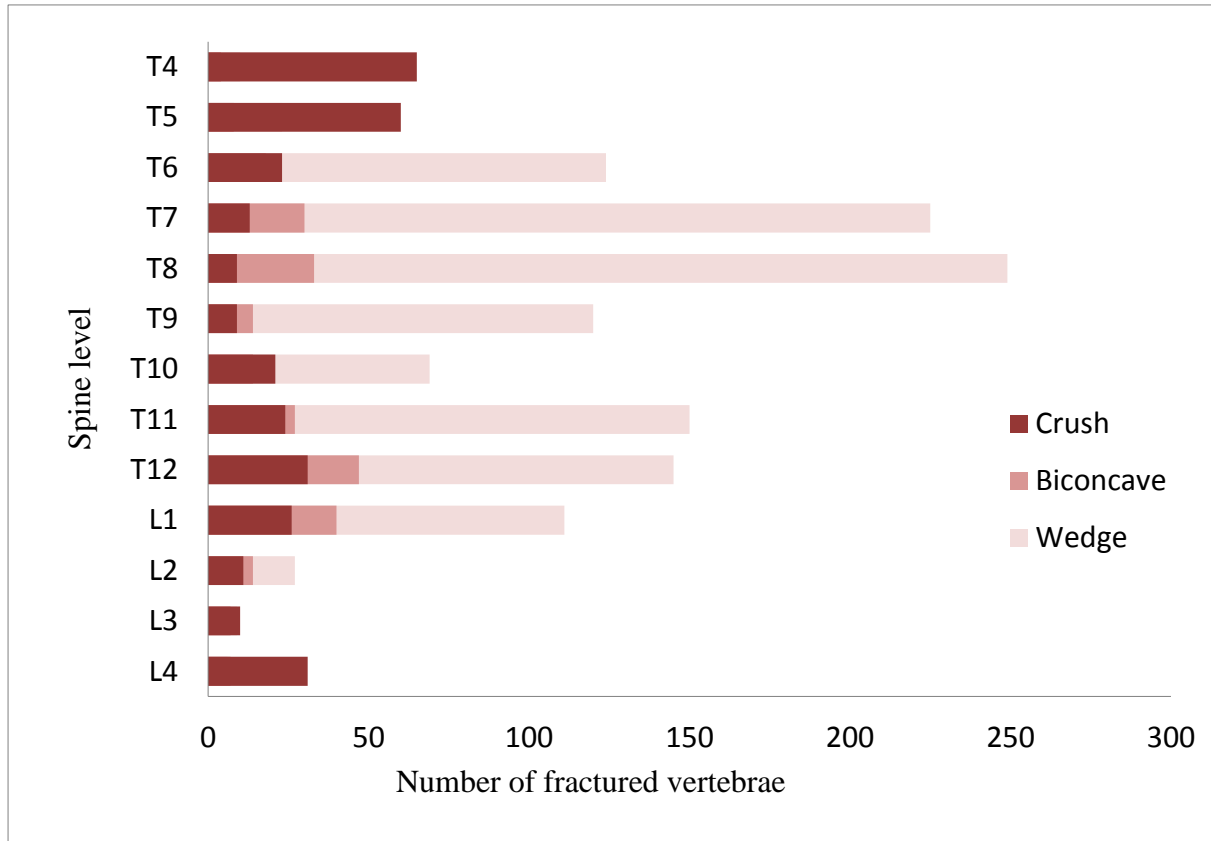
Supplementary Fig 1. Flowchart of the study participants. There were 14,926 participants that were eligible for radiographic examination and 3,260 did not undergo the exam. Out of 11,666 participants with radiography data, 828 were scored with ABQ but not with QM SA and 275 were scored by QM SA but not by ABQ, reducing the number of participants with radiographs scored by both methods to 10,563. Since we decided to perform analyses not only per subject level but also per vertebral body, we excluded participants that had any missing data from T4 to L4 level. Those missing were due to suboptimal visibility and no informed consent; this filter reduced the study population to 7,582 participants



Supplementary Fig 2. Distribution of QM fractures in the spine by morphometry. Crush

fractures are mostly located at the upper thoracic level at T4-T5, biconcave at T7-T8 and

T12-L1 and Wedge at T7-T8 and T11-T12.



Supplementary Table 1: Frequencies of QM SA vertebral fractures by shape and severity. In a population of 7,582 subjects, there were 1,574 vertebral bodies fractured of which 54.0% were wedge, 11.9% biconcave and 19.3% crush. On the other hand, 49.2% were classified as mild deformities, 30.8% as moderate deformities and 4.7% as severe

N=7,582 n=1,574	Wedge (n=842)	Biconcave (n=188)	Crush (n=304)
Mild- Grade 1 (n=775)	441	97	237
Moderate-Grade 2 (n=485)	348	73	64
Severe-Grade 3 (n=74)	53	18	3

Supplementary Table 2. The agreement between QM SA and ABQ and distribution in 2x2 tables in different scenarios; when applying the standard QM definition to QM SA, excluding mild deformities from the definition or assessing mild deformities based on endplate depression.

a) Agreement statistics for the study population when using the standard definition

for QM SA

		ABQ	
		+	-
QM SA	+	199	877
	-	105	6,401

	(N=7,582)
QM SA (%)	1,076 (14.2)
ABQ (%)	304 (4.0)
Kappa	0.24
Observed agreement	0.87
Expected Agreement	0.83
BI	0.10
PI	-0.81
Positive agreement	0.29
Negative agreement	0.93
PABAK	0.74

b) Agreement statistics for the study population when excluding subjects who had mild QM SA deformities

		ABQ	
		+	-
QM SA	+	160	302
	-	105	6,401

	(N=6,968)
QM SA (%)	462 (6.6)
ABQ (%)	265 (3.8)
Kappa	0.41
Observed agreement	0.94
Expected agreement	0.90
Bias Index	0.03
Prevalence Index	-0.89
Positive agreement	0.44
Negative agreement	0.97
PABAK	0.88

1 **c) Agreement statistics for the study population if we re-examine mild fractured**
 2 **subjects based on presence of endplate depression.** Out of 614 subjects that had a
 3 mild fracture, 39 were classified as fractured also by ABQ. If we classify those 39
 4 mild+ and ABQ + as true positives and the 575 remaining we classify as true
 5 negatives, the redistributed 2x2 table would look like the one below. Calculating
 6 agreement statistics for that 2x2 table, produces even higher agreement than just
 7 excluding those deformities from the study analysis.

	ABQ	
	+	-
QM SA	+	-
	199	302
	-	6,976

	(N=7,582)
QM SA (%)	501 (6.6)
ABQ (%)	304 (4.0)
Kappa	0.50
Observed agreement	0.95
Expected Agreement	0.90
Bias Index	0.03
Prevalence Index	-0.89
Positive agreement	0.50
Negative agreement	0.97
PABAK	0.90