# Journal Pre-proof

Data sampling strategies for accurate fault analyses: A scale-independent test based on a Machine Learning approach

Tiago M. Alves, Joshua Taylor, Padraig Corcorant, Tao Ze

Please cite this article as: Alves, T.M., Taylor, J., Corcorant, P., Ze, T., Data sampling strategies for accurate fault analyses: A scale-independent test based on a Machine Learning approach, *Journal of Structural Geology*, https://doi.org/10.1016/j.jsg.2025.105342.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

1    **Data sampling strategies for accurate fault analyses: A scale-independent**

2    **test based on a Machine Learning approach**

3    **Tiago M. Alves[1], Joshua Taylor[1,2], Padraig Corcorant[2]; Tao Ze[3]**

4

5    [1]3D Seismic Lab, School of Earth and Environmental Sciences, Cardiff University, Main

6    Building, Park Place, Cardiff, CF10 3AT, United Kingdom (email: alvest@cardiff.ac.uk)

7    [2]School of Computer Science and Informatics, Cardiff University, Abacws Building,

8    Senghennydd Road, Cathays, Cardiff, CF24 4AG, United Kingdom

9    [e]Key Laboratory of Tectonics and Petroleum Resources, Ministry of Education, China

10    University of Geosciences, Wuhan, 430074, China

11

12    **Abstract**

13    Seismic and outcrop data from SE Brazil, Greece and SW England are used to develop

14    a new method to correctly identify tectonic fault segments – either active or quiescent - using

15    a machine learning approach. Three-dimensional (3D) analyses of tectonic faults are often

16    based on the mapping of throw values (T) along their full length (D) or depth (Z) using a

17    wide range of data. Yet, the collection of these throw values using geophysical or outcrop

18    data is often time-consuming and onerous. In contrast to many empirical measurements of

19    T/D and T/Z, our new method supports the mapping of active (or potentially active) fault

20    segments and limits data undersampling, a caveat that results in the grouping of faults as

21    single zones, systematically overlooking their natural segmentation. The new method is

22    scale-independent and resulted in the definition of a minimum sampling ratio necessary for

23    accurate fault segment mapping. Determined through the gradual downsampling of T/D and

24    T/Z data to a critical point of information loss, the minimum sampling interval (δ) in T/D and

25    T/Z data, expressed as a percentage of fault length, or height, is: a) 1.02% ± 0.02 for faults

26    that are longer or higher than 3.5 km; b) 4.167% ± 0.18 for isolated faults that are shorter

27    than 3.5 km in either length or height. This work is therefore important as it shows that one

28    should never acquire T/D and T/Z data above a threshold δ value of 4% to identify

29    successive, linked fault segments, whatever their scale. Total accuracy in fault-segment

30    detection is only assured for δ values of 1% when in the presence of fault zones with

31    segments longer than 3.5 km. As a corollary, we confirm that T/D and T/Z data are often

32    undersampled in the published literature, leading to a significant bias of subsequent

33    interpretations towards coherent constant-length growth models when analyzing both active

34    and old, quiescent fault systems.

35

36    **Keywords:** Data sampling; Machine Learning; Tectonic faults; fault growth; Sampling

37    errors; Fault propagation models

38

39    **1. Introduction**

40        The mapping and geometrical characterization of faults and joints at varied scales of

41    observation are vital to geological, structural and earthquake-risk analyses. Recognizing

42    faults and joints is also important in hydrocarbon and geothermal energy production, in

43    engineering works, and to the implementation of sub-surface storage solutions

44    (Gudmundsson et al., 2002; Misra and Mukherjee, 2018; Trippetta et al., 2019; Torabi et al.,

45    2023). Measurements of both active and quiescent tectonic faults need to be accurate

46    because: a) the trapping and accumulation of subsurface fluid often depend on the geometry

47    and interaction styles of faults and joints (Yielding, 2015), b) drilling-related hazards are

48    frequent in highly-faulted areas, as well as in prospects where reservoir quality is much

49    reduced by joint systems (Saeidi et al., 2014; Kozłowska et al., 2017), c) the migration and

50    preservation of sub-surface fluid is, in many a prospect, associated with the timing of

51    formation, growth, and sizes of tectonic faults and joints (Ferrill et al., 2017; 2020). The size

52    of tectonic faults, their inherent geometry, and the location of their intersection (linkage)

53    points are important for a safe and sustainable production of geological resources (Jentsch et

54    al., 2020; Purba et al. 2019; Moska et al., 2021; Huenges et al., 2013).

55        Another characteristic of tectonic faults is that their size is a predictor of earthquake

56    magnitude, i.e. faults over a certain length are capable of generating destructive earthquakes

57    and associated geohazards (Trippeta et al., 2019). Importantly, large magnitude earthquakes

58    can be generated in seemingly discrete fault segments that are connected to form a single

59    large fault zone, at depth, whereas relatively isolated, smaller fault segments present a much

60    lower seismic risk (Cloetingh et al. 2010; He et al. 2019; Alves, 2024). Generally speaking,

61    fault intersections, geotechnically unstable fault zones and active faults capable of generating

62    earthquakes must be avoided in engineering projects. The systematic undersampling of

63    tectonic faults' geometries can result in a rapid degradation of infrastructure after the

64    completion of construction works, or in unexpected cost increases (Aydin et al., 2004; Wang

65    et al. 2022).

66        Notwithstanding the fact that accurate fault analyses are crucial in structural geology,

67    producing a high-resolution image of fault structures is time-consuming, expensive and not

68    always practical. Automatic methods to extract faults using remote sensing data have been

69    developed by authors such as Gloaguen et al. (2007), but these types of data are not always

70    available or concern the scales of observation necessary for a particular aim, or analysis. In

71    most instances, the solution followed by both academia and industry is to reduce data

72    collection to the minimum required level, just enough to understand where the loci of fault

73   interactions are. Unfortunately, such an approach results in coarse and often random data

74   sampling ratios and techniques, as one can easily verify in most scientific articles published

75   in the past 20-30 years. Purposely anonymous examples from the published literature,

76   include: a) fault-throw values measured every 600 m for a single fault segment that is 4 km-

77   long, in which only 6.6 data points were acquired for such a segment (for an average of 15

78   measurements per fault in the same article), b) fault throws measured every 3 cm for faults

79   that are 1.0 m-long, returning 33 measurements per fault, on average, c) a third example in

80   which fault throws are measured every 100-200 m for a fault zone that is 4 km-long,

81   returning an average of less than 20 data points per fault segment. The coarse, and often

82   random samplings of throw/distance (T/D) and throw/depth (T/Z) data are common in the

83   literature, and surprisingly not always deriving from the use of seismic and geophysical data

84   of relatively poor resolution.

85        Such erratic sampling strategies can lead to a generic failure in recognizing that fault

86   segments are components of a larger fault zone (Walsh et al., 2003). In fact, Tao and Alves

87   (2019) have shown the systematic undersampling of fault throws in seismic, remote sensing

88   and outcrop data will inevitably lead to an over-reliance of models reflecting coherent 'fast

89   propagation' styles of fault growth (Walsh et al., 2003; Nicol et al., 2020). In other words,

90   naturally segmented faults, or fault zones, will appear as single long structures if fault throws

91   are undersampled. This caveat is compounded when interpreters overlook map-view

92   geometries and concentrate only on collecting throw values without an accurate structural

93   mapping accompanying their workflows.

94        The aim of this work is to produce reliable predictions of fault segmentation in an

95   automated manner, without human bias, and avoiding any under- or overfitting of data to

96   emphasize a particular fault growth model. Overfitting in this case would involve finding

97   more faults than exist through misinterpretation of signal noise and height undulation caused

4

98    by erosion or poor exposure, for instance. Underfitting would be to exaggerate fault throw so

99    that multiple segments appear coherent in their growth and part of a single fault zone (Torabi

100   and Berg, 2011; Tao and Alves, 2019). Our approach is scale-independent and works for both

101   active and quiescent (ancient) tectonic faults that may or may not reactivate by anthropogenic

102   means. In summary, the main research questions addressed in the work include:

103

104      a)   What mathematical methods can be applied to Machine Learning tools to avoid

105           interpretative errors when identifying tectonic faults?

106      b)   What are the implications of misrepresenting fault segmentation in terms of

107           understanding their growth modes?

108      c)   What are the threshold fault-throw (or displacement) values necessary for a correct

109           identification of fault segmentation in nature?

110

111   **2. Theoretical aspects concerning fault-segment recognition**

112   *2.1 Coherent vs. isolated growth modes and scale variance in structural observations*

113        Tectonic faults and joints, universally named as 'rock fractures' in the published

114   literature, comprise sets of related segments, or strands, that can be kinematically and

115   spatially related (Pollard and Segall, 1987; Gudmundsson, 2012) (Fig. 1). They represent

116   continuous, brittle breaks in rocks, be it crustal-scale stress in the case of tectonic faults or

117   smaller localized stresses that hardly offset rocks in the case of polygonal faults and joints

118   (Peacock et al., 2017; Laubach et al., 2018). The largest of faults, those documenting a clear

119   vertical or horizontal offset in strata or rocks, are often part of system of related fault

120   segments that interact and link - and are restricted to a relatively narrow band or volume -

121    also called a Fault Zone (Peacock et al., 2000; 2017). Fault zones are formed by the 3D

122    linkage of multiple segments in a broad region of deformation, leaving behind fault segments

123    not frequently affected by such a strain (Rotevatn et al., 2019).

124        Fault segments may show geometries that are indicative of a 'fault-linkage' and

125    'coherent' growth (Kim and Sanderson, 2005) or, instead, develop individually to obey an

126    'isolated' growth mode (Walsh et al., 2003) (**Fig. 2**). In practice, many 'linked' or 'coherent'

127    faults are part of a larger zone of deformation, while isolated faults show growth histories and

128    throw distributions that are independent or disparate from nearby faults segments (Nicol et

129    al., 2020) (**Fig. 2**). The recognition of such fault growth modes in geophysical or outcrop data

130    relies on the correct mapping of fault throws (T) against fault zone length (D) and depth (Z)

131    to produce T/D and T/Z plots (Cartwright et al., 1998; Baudon et al., 2008) (**Figs. 1 and 2**).

132    Multiple examples of how throw data can be used to understand fault growth modes are given

133    in the literature for Norway (Tvedt et al., 2013; King and Cartwright, 2020), SE Brazil

134    (Varela and Mohriak, 2013; Plawiak et al., 2024), Gulf of Mexico (Cartwright et al., 1998;

135    Shen et al., 2018) and for tectonically active areas in the Gulf of Corinth (Fernández-Blanco

136    et al., 2019; Robertson et al., 2020; Nixon et al., 2024), offshore Crete (Caputo et al., 2010;

137    Nicol et al., 2020; Mechernich et al., 2023) or the Basin and Range, where topographic

138    information has been combined with local tectonic analyses (Lee et al., 2023). Whenever

139    available, fault displacement should be used instead of throw (see **Fig. 3**), but the acquisition

140    of such data is time-consuming in practice when analysing outcrop or geophysical data - as a

141    result, fault throw (T) is more frequently measured (e.g. Cartwright et al., 1998). Fault throw

142    is a measure of the vertical distance between the footwall tip of a fault and its corresponding

143    hanging-wall tip (Mukherjee, 2019) (**Fig. 3**). Fault displacement concerns the total movement

144    of two fault blocks along a fault plane, measured in any specified direction. It represents the

145    distance between two separated pieces of a marker layer on both sides of a fault. The time

6

146     and effort needed to collect such fault data is often the source of 'censorship' and 'truncation'

147     in data (Torabi and Berg, 2011), leading to incorrect assumptions regarding the relative

148     timing of fault activity.

149         A caveat often overlooked by structural interpreters is that recognizing fault segments

150     depends on the distinction of meaningful throw gradients that represent segment linkages on

151     T/D (or $D_{max}$/L) plots, accompanied by their analysis on vertical sections and map view

152     (Walsh and Watterson, 1991; Walsh et al., 2002, 2003; Kim and Sanderson, 2005) (**Fig. 1**).

153     With lower resolution images, or remote-sensing data of lower quality, comes a high level of

154     uncertainty over the linkage points of discrete fault segments when acquiring such T/D or T/Z

155     data. The lack of chronostratigraphic markers can also result in the misinterpretation of

156     important gaps between faults, and multiple small segments may appear as a single large fault

157     when a low-resolution dataset obscures lows, or minima, in throw (Tao and Alves, 2019).

158

159     *2.2 Use of T/D and T/Z data in fault-segment recognition*

160         Fault throw/distance (T/D) and throw/depth (T/Z) data are often measured on a seismic

161     section, or exhumed fault plane, in order to identify distinct fault segments and interpret a

162     fault propagation mode (Torabi and Berg, 2011). Fault throw (T) is often used instead of

163     displacement as it is an easier variable to define, and quantify, in geophysical and field data,

164     regardless if a fault is planar or listric. Throw measurements in listric faults will overlook

165     their horizontal component (heave) but can still be used to identify discrete fault segments.

166     In parallel, throw/distance (T/D) plots measure throw distributions along a fault's length and

167     can be complemented by Throw-Depth (T/Z) measurements. While T/D data help

168     recognizing distinct, linked fault segments, T/Z data indicate the areas where the mechanical

169   properties of rocks may vary across a fault, at the same highlighting any evidence for vertical

170   fault linkage (Cartwright et al., 1998; Baudon et al., 2008).

171       Distinct faults, and also their constituting segments, show distinct orientations and

172   curvatures in map view (Kim and Sanderson, 2005). On T/D profiles, steep decreases in

173   throw values relate to the existence of an intersection (a 'hard' or 'soft' linkage point)

174   between two fault segments or, instead, points out to a fault's lateral tip (**Figs. 1 and 2**). Two

175   linking fault segments will also be recorded as sudden gradient changes in T/D and T/Z plots

176   (**Figs. 1 and 2**). Conversely, variations in fault height caused by erosion and local sediment

177   deposition will be seen as high-frequency, low-magnitude undulations that resemble a noise-

178   like pattern of throw distributions (Torabi et al., 2019). Throw and displacement can be

179   particularly affected by erosion of a fault scarp, as both are measured from a defined height at

180   the immediate footwall block of a fault (**Fig. 3**).

181

182   **3. Data and Machine Learning methods**

183

184   *3.1 Fault-throw data*

185       Measurements of fault throw used in this work were taken from distinct parts of the

186   world (**Fig. 4a**). T/D and T/Z measurements for 415 faults were used in our analysis - they

187   were collected at regular intervals and used to test the sampling distance necessary to

188   correctly interpreted fault linkages and their growth modes. The primary source of data

189   comes from the Southern North Sea and Southeast Brazil (Alves et al. 2022; Zhang et al.,

190   2022; Tao and Alves, 2016; 2019). Outcrop data were gathered in various locations in Crete

191   and Somerset (Tao and Alves 2019; Gaki-Papanastassiou et al. 2009; Caputo et al. 2010;

192   Alves and Cupkovic 2018).

193

194   *3.1.1 3D Seismic data*

195       Seismic data in this work comprises two high-quality seismic volumes from the SE

196   Brazil (**Figs. 4a,b and 5a**). The volume was stacked with a bin (or trace) spacing of 12.5 m

197   and a vertical sampling rate of 2 ms. The vertical resolution of the investigated seismic data

198   varies from 5 to 8 m near the seafloor, and c. 12 m at the maximum depth of faults

199   investigated in this work (**Fig. 5a**). Fifty-nine (59) faults, including crestal faults, radial faults

200   and low-angle normal faults flanking salt diapirs were interpreted every 1, 3, 5, 10 and 20

201   inlines and crosslines (**Fig. 5a**). Composite lines were also used, when needed, to collect data

202   perpendicularly to fault-plane dip. Interpreted faults are 225 m to 5,000 m long and show

203   throw values varying from 6 ms to 73 ms two-way time (twt). These faults are still active at

204   present as some offset strata that are very close to the modern seafloor due to on-going salt

205   tectonics in SE Brazil (**Fig. 5a**).

206

207   *3.1.2 Ierapetra Fault Zone (SE Crete)*

208       The modern Ierapetra Fault Zone is located in SE Crete and is >25 km long (**Fig. 4a,c

209   and 5b**). It has been active since, at least, the Late Miocene and is one of the most prominent

210   structures on the island (Caputo et al., 2010; Gaki-Papanastassiou et al., 2009). Several fault

211   segments striking NNE–SSW and dipping to the WNW played a crucial role in the evolution

212   of the fault zone, namely the Kavousi, Ha and Ierapetra segments (Gaki-Papanastassiou et al.,

213   2009) (**Fig. 5b**). Each of these segments has its own characteristic geometry (**Fig. 5b**). Due to

9

214    its activity, thick sediments cover the fault zone's hanging-wall, while the immediate

215    footwalls are barren of marine sediment and feed adjacent basins at present (**Fig. 5b**).

216        Throw/distance (T/D) data reveal that the fault segments are 0.5 to 7.1 km long, show

217    maximum throw values between 250 and 1000 m. Nevertheless, a synchronous Holocene

218    reference horizon was identified in the study area and used as a marker to compile T/D plots

219    for outcropping fault segments (**Fig. 3**). During the collection of fault-throw data, the

220    following were performed:

221        (i)    Fault scarps were mapped in detail in the field and projected on 1:50,000 maps

222    from the Hellenic Mapping and Cadastral Organization – the maps with the highest resolution

223    in the region. The present-day height of footwall tips and any associated erosional and

224    depositional features were taken into consideration in our throw measurements of active

225    tectonic faults,

226    (ii)    Throw data were collected at a regular interval of 50 m along the fault segments

227    observed in the field. Throw measurements were gathered where the geometry of the faults is

228    clear on the maps and in panoramic photos (**Fig. 5b**).

229

230    *3.1.3 Sub-seismic scale faults from SW England (Kilve)*

231        The Bristol Channel Basin records four distinct stages of faulting: 1) N-S extension and

232    associated normal faulting in the Mesozoic, accompanying the development of the Bristol

233    Channel Basin, 2) reactivation of some of the normal faults formed during the first stage, 3)

234    reverse reactivation of Mesozoic and older structures during the Alpine orogenic pulses

235    (Underhill and Paterson, 1998), 4) reverse-reactivation of normal faults that were

236    subsequently cut by conjugate strike-slip faults (Dart et al., 1995), 5) jointing of strata after

237    Alpine-related fault reactivation (Rawnsley et al., 1998).

238         A certain degree of tectonic reactivation thus occurred in the Bristol Channel Basin

239    during the Cenozoic and was of an enough magnitude to generate: a) structures formed by N-

240    S contraction - chiefly reverse reactivated planar normal faults, b) structures formed by east–

241    west contraction, c) intersecting N- to NNW-trending and NE-trending faults (Glen et al.,

242    2005). Importantly, the faults analysed in this paper were formed by N-S extension, record no

243    apparent tectonic reactivation, and only occur in Liassic limestones and shales (Peacock et

244    al., 2017).

245         Thirteen (13) faults with lengths varying from 1.65 m to 7.55 m, and maximum throw

246    values ranging from 3 cm to 29 cm, were measured and interpreted in the field (**Figs. 4 and

247    5c,d**). Fault-throw measurements in the field depended on how clear they were exposed at the

248    surface. Throw values were measured where the hanging-wall and footwall were totally

249    exposed on the two sides of the fault trace. The throw-distance data were acquired along the

250    exposed fault trace every 5 cm. T/D plots were also computed and analyzed for these faults

251    considering different sampling spacings as exemplified in the Supplementary Materials in

252    Tao and Alves (2019).

253

254    *3.2 Machine Learning and mathematical algorithms*

255         Machine Learning algorithms were implemented using the Python programming

256    language applied on NumPy (Harris et al., 2020), PyWavelets/Pywt V1.4.1 (Lee et al., 2023)

257    and SciPy 1.0 (Virtanen et al., 2020) software libraries.

258

259     *3.2.1 Wavelet transforms for fault-segment detection*

260     The main advantage of using Wavelet Transforms to detect discrete fault segments is

261     that they permit the analysis of features that vary in character over different scales

262     (Kalbermatten et al., 2012; Shen et al., 2022). For acoustic or optical signals, such features

263     are often frequencies varying over time. In image data, features of interest include edges and

264     textures, as is the case of throw maxima and minima in T/D and T/Z curves (Shen et al.,

265     2022), or object-based classes of images recorded after segmenting remote sensing data into

266     homogeneous regions (Gloaguen et al., 2007).

267     In mathematical terms, Wavelet Transforms allow for the decomposition of an input

268     signal into the intensity of individual frequency bands. The advantage of the Wavelet

269     Transform method over the Fast Fourier Transform is the former's ability to identify both the

270     frequency and spatial position of frequencies in the data. Fast Fourier Transforms only

271     provide frequency information over a fixed range, with no location value along that range

272     (Sifuzzaman et al., 2009). A wavelet can thus be convolved with a signal and the resulting

273     signal gives the intensity of the wavelet at each point along the signal. The wavelet size can

274     be changed to give an intensity for each frequency band.

275     In order to have a successful Wavelet Transform, a wavelet must follow a set of

276     criteria, namely the wavelet function $\psi$ needs to return a zero average:

277

278     $\int_{-\infty}^{+\infty} \psi(t)dt = 0.$                    (Eq. 1)

279

280     The wavelet is then multiplied by a scale parameter *s* and translated by *u* such as:

12

281

282    $\psi_{u,s}(t) = \frac{1}{\sqrt{s}} \psi\left(\frac{t-u}{s}\right)$                        (Eq. 2)

283

284    The Wavelet Transform of $f$, at a scale $s$ and position $u$, is finally computed by correlating $f$

285    with a wavelet atom:

286

287    $Wf(u,s) = \int_{-\infty}^{+\infty} f(t) \frac{1}{\sqrt{s}} \psi^*\left(\frac{t-u}{s}\right) dt$                        (Eq. 3)

288

289          The above equations use $t$ as the measure of displacement across the signal, as wavelets

290    are most often related to signals represented as a function of time. In our particular case the

291    Wavelet Transform will not be reprocessed as a function of time; it will be estimated along a

292    measured distance, making no difference to the mathematics used. Time (t) will be replaced

293    by distance (D) in Wavelet Transforms, this parameter D being the distance along a fault

294    plane considered earlier in the paper, with frequency and wavelength being also be processed

295    in reference to distance. By convention, t is used in signal processing, but for our case study

296    distance (D) is used where t is seen in Equations 1 to 3.

297          Particular wavelet types are more often used in signal processing, and is thus best to

298    choose one of these common types when performing a Wavelet Transform. A wavelet that

299    follows a similar shape to the expected signal is required to get the best results (Mallat,

300    2009). In this work we used the so-called *Ricker wavelet* (see **Fig. 6**). Such a wave shape

301    allows for the isolation of peaks, or throw maxima, in a fault segment, with throw minima

302    being mathematically defined as the wavelet boundaries – as with distinct fault segments that

13

303    are part of a fault zone (Mallat, 2009). Hence, the *Ricker wavelet* closely matches the linkage

304    behaviour of fault segments (see Wang, 2015a; 2015b), i.e. it better identifies sharp throw

305    minima, which are known to indicate the places where distinct fault segments were originally

306    linked (**Fig. 6**). Such an approach results in the strongest correlation possible whenever the

307    *Ricker wavelet* is equal to the fault size.

308

309    *3.2.2 Polynomial regressions as a complementary method for fault-segment detection*

310        Polynomial regressions follow a similar process to linear regressions whereby a line

311    with a minimum average distance to the data points is found. Such a distance is quantified by

312    a Sum of Square Errors (Heiberger and Neuwirth, 2009). The advantage of such a polynomial

313    regression relates to the ability of using a higher order equation to define the line of best fit to

314    T/D and T/Z data, respectively where the x-axis is distance (D) and the y-axis corresponds to

315    depth (Z). In the case of a third order polynomial, the key values are the coefficients, and a

316    polynomial regression model can be simplified to these coefficient terms, i.e. the terms can

317    be used as predictors for the values in the 'real' field data. This simplification to a single

318    equation is important in our work, as it allows for data comparisons for the same fault

319    whenever the T/D and T/Z measurements are correctly sampled vs. when data are

320    downsampled.

321        To avoid data overfitting, we used a lower degree polynomial of degree 3. In practical

322    terms, a three-parameter polynomial equation is first generated for each of the identified fault

323    segments. The absolute minimum number of sample points used to generate this first model

324    of fault shape is three (3), so the sampling space is so low that only the two tips and the point

325    of maximum throw of a fault are identified (e.g. **Fig. 7**). The purpose of this method is to

326    allow a comparison of fault detection approaches, using different sampling ratios, by

14

327   reducing them all to the same dimensions. The low complexity of this method also helps to

328   ensure that the model does not overfit the T/D and T/Z data in this work. We verified that the

329   above findings could be generalized to our specific data by verifying that the error between

330   the model and the T/D and T/Z data in question was small. A more detailed explanation of

331   the polynomial regression process can be found in Ostertagová (2012) and James et al. (2013)

332   and Section 4.4 in this work.

333      The modelling of faults via polynomials works well due to the process of fault creation,

334   itself the result of forces, or stresses, developing and growing fractures in a volume of rock.

335   Over geological periods of time, such forces change in terms of their direction and

336   magnitude, and multiple factors can cause local variations in space and time (4D) in stress-

337   strain relationships (Kim and Sanderson, 2005). At a single point in time, a skewed

338   polynomial shape can accurately follow the shapes of faults and joints in nature, as the forces

339   acting on a volume of rock result in a fault following a path of least resistance. This promotes

340   the formation, in nature, of Gaussian T/D and T/Z curve shapes in faults and joints. The

341   various (unpredictable) factors acting on these same structures, and altering their T/D and T/Z

342   profiles, can thus be simplified as skewed Gaussian curves. Goff (1991) found that a skewed

343   Gaussian curve provides a model of low complexity that accurately fits our type of data.

344

345   **4. Results**

346      As a summary, the workflow used in this work is shown in **Fig. 8** to highlight the

347   different steps of the proposed machine learning methodology.

348

15

349   *4.1 Step 1 – Application of discrete Continuous Wavelet Transforms (CWT) to resolve faults*

350   *at different scales*

351   Theoretically, Continuous Wavelet Transforms (CWTs) can produce a 2D plot of

352   frequency band strength. For the purposes of this work, these band strengths correspond to

353   variations in fault throw (T) when this throw is interpreted as a part of a wave. Hence, fault

354   segments in the field or in seismic data can be represented as wavelets.

355   In this work, computed CWTs were visualized against T/D plots, with a clear

356   correlation being observed between frequency band strength and the throw maxima recorded

357   for each fault segment (**Figs. 7 and 8**). In fault zones containing multiple throw maxima, the

358   largest fault segments correlate with a peak in low frequency wavelet amplitude (**Fig. 7**).

359   When performing a CWT, the wavelets of various frequencies are compared across the

360   input signal. The correlation of the signal with that wavelet is measured at each point.

361   Therefore, when reaching the throw maxima of fault segments with a similar frequency, the

362   accuracy behind correlating wavelets with T/D (and T/Z) plots (i.e. *correlation strength*)

363   reaches a maximum. Such an approach simplifies the analysis of fault segments by splitting

364   the throw measurements made in the field, or in seismic data, into frequency bands. This

365   allows a computer algorithm to pick out certain frequencies that are likely to correspond to

366   fault segments. Wavelets that are most similar in shape and size to fault segments, will result

367   in a higher correlation between the CWTs and real T/D and T/Z data after convolution. This

368   means the peak in convolution output will give the 'best match' from possible wavelet sizes

369   and locations along a fault. Peaks in the CWTs' output can then be assumed to be the 'top'

370   (i.e. the point of maximum throw) of a fault of a particular size.

371   Modifications were made in our analysis to the CWT technique so that discrete fault

372   segments could be found. The main modification consisted in changing how the frequency of

16

373    a fault is decided. Initially the frequency of discrete faults was found by identifying the point

374    of greatest throw amplitude across a fault, or a fault zone, for each frequency. We then

375    visually confirmed which of these amplitude maxima coincides with the throw maxima of

376    fault segments by comparing then with acquired throw data, acquired at maximum resolution.

377    However, such an approach was deemed unreliable when: a) multiple segments in a fault

378    zone show similar lengths, and b) an entire fault zone follows a shape similar to the wavelet,

379    in which case a very low frequency will be used, spanning the entire fault zone. This caveat

380    results in the smallest segments being ignored by the algorithm. An example can be found in

381    **Fig. 7**, where some obvious local fault segments were missed.

382    A successful solution was found using an approach that required the application of a

383    computational step to remove the highest frequencies representing a 'noisy' signal. Peaks and

384    Troughs in the computed wavelets were found by a comparison of points to their immediate

385    neighbours (see Section 4.2). This is called mathematically as calculating *prominence*.

386    Prominence is calculated by finding the minimum between a peak and the next higher peak,

387    so the comparison happens over a range around a peak, not just the immediate neighbours - a

388    full mathematical explanation of *prominence* is given in

389    https://www.mathworks.com/help/signal/ug/prominence.html. If a point was found to be

390    higher in value than its adjacent points and had sufficient *prominence* in the whole of the

391    fault zone, it was taken as a Peak by our algorithm. Troughs were found in the same way,

392    using an inverse algorithm so the same function can be used.

393    In a second stage of this process, a wavelet band was chosen by removing wavelets that

394    are not considered relevant, as they mostly represent noise (**Fig. 8**). The highest wavelet

395    frequency band remaining in the dataset was then considered to be ready for fault scanning.

396    The use of the highest wavelet frequency avoided locating the longest faults early in the

397    process, as the scale (and wavelet range) of these long faults usually overprint the smallest

17

398    fault segments before these are found. For instance, in **Fig. 7** we can identify discrete throw

399    maxima relating to the presence of small fault segments that were overlooked by the

400    algorithm that, in Step 1, was focused on picking the greatest throw maxima. This means that

401    the identification of relatively small throw maxima need to be prioritized in a Machine

402    Learning approach.

403        In summary, the maximum value in the wavelet band that is not interpreted by the

404    algorithm as a discrete fault segment was defined as the maximum throw value of a new

405    segment. Conversely, the throw minima on each side of this maximum were taken as

406    comprising the lateral tips of a fault segment. Such a method could be applied to a map all

407    maxima and minima in the produced CWT matrix. This method allows for a rigorous

408    definition of fault segment distribution and their linkage points. To avoid errors in our

409    analysis a cross-validation was used to select the most suitable frequency. We split the data

410    into training and test cases. The frequency was selected using the training cases and this was

411    determined to be a suitable frequency through evaluation of the test cases. A subset of the

412    dataset was chosen randomly to use for validation of the frequency constant. A ground-

413    truthed set of fault locations was then marked on the dataset. The constant that came closest

414    to this ground-truthed data was taken and modified by smaller amounts for a different subset

415    of the data, repeating the same process to address any bias introduced.

416

417    *4.2 Step 2 – Detection of throw gradients from the point of throw minima*

418        Step 2 in this work consisted in the application of a gradient descent from the point of

419    frequency minima. The aim was to find the nearest throw minimum representing the linkage

420    point between two fault segments. If no frequency minima are found before reaching the end

421 of the dataset, the last value picked by the algorithm is taken as the end of the segment (**Figs.**

422 **8 and 9**).

423       The method consisted in the scanning of every wavelet frequency for their Peaks and

424 Troughs, which are then reduced down to frequencies that contain enough Peaks and Troughs

425 to form at least one discrete fault segment. A second reduction is completed by removing the

426 frequencies that result in too many Peaks per meter. Such a step is important for removing

427 frequencies that reflect irrelevant, spurious throw maxima, usually comprising measurement

428 errors and resolution issues when measuring throw data in seismic and at outcrop (truncation

429 and censoring cf. Torabi and Berg, 2011). The threshold Peak values can be changed, with a

430 stricter threshold resulting in the identification of only the larger fault segments (see **Fig. 9**),

431 and a looser threshold resulting in multiple fault segments being found. Naturally, if it is set

432 too loose, unwanted segments may appear in one's fault tracing.

433

434 *4.3 Step 3 – Integration of Continuous Wavelet Transforms (CWTs) with a threshold Peak*

435 *rate*

436       To improve the accuracy of our results, a re-sampling was applied as a third step before

437 undertaking a CWT. The sample count was scaled to 1,000 times the longest wavelet length,

438 which resulted in a less unusual behavior whereby segments are too large to be detected by

439 any of the wavelets. This allows the wavelet bands to be kept the same for all tests, even

440 while the dataset sample sizes vary. After all processing was done, an optional process

441 allowed for the joining of the segments, to remove gaps between them. Step 3 returned more

442 accurate results when undertaken on a series of faults where no gaps are expected, i.e. the

443 approach also meant the lowest throw between any two segments was always considered as

444 the linkage point of successive fault segments, regardless of their scale in nature.

19

445      A value between 0.03 and 0.04 for the threshold Peak rate (number of peaks per

446     sample) was found to provide good results in the datasets tested in this work. This value was

447     decided by plotting Peak rate values against frequency and finding where the graph **in Fig. 10**

448     begins to level out. In this graph, the rapid descent recorded with increasing wavelength sizes

449     represents the reduction in noise occurring as the small changes in throw are filtered out by

450     the algorithm. Once this noise is filtered out, and the curve approaches a flat, we can be

451     confident that the remaining data is accurate. Cross-validation can then be used to select the

452     above values for the threshold Peak rate.

453      Selecting a threshold Peak rate must be consistently applied across all tests to make

454     one's results comparable later on. However, in practical terms, the threshold can be changed

455     based on the smallest fault sizes one has to find in a dataset, although using a threshold too

456     high results in the detection of throw maxima that are the result of random noise or constitute

457     irrelevant changes in fault height in a discrete segment. The adoption of a 0.04 Peak rate

458     returned positive results in this work - all Peaks that are clearly not part of faults were

459     ignored, without overlooking any possible faults, examples of which can be seen in **Fig. 11**.

460     Cross-validation against the ground-truthed throw data was again used to obtain a value of

461     0.04. If a different dataset with different properties is used, then cross-validation is also

462     performed with respect to that dataset to select the most appropriate value. In practice, 0.04

463     was chosen by validating it across a large dataset and should be considered a 'default' value

464     to use but can also be changed depending on whether its use causes false positive or false

465     negative faults. This approach allowed us to use the previously defined method of Wavelet-

466     Transform scanning described in Section 4.1, starting with the highest frequency, as we have

467     now removed noisy wavelet bands that could hinder such a Machine Learning approach.

20

468     *4.4 Step 4 – Throw-profile fitting via a cubic model*

469     The computational steps so far described are successful in identifying the tips and

470     throw maxima for each fault segment, but fault shape is often not accurately depicted. To best

471     represent fault segment shape, a third order polynomial regression needs to be applied

472     individually to each fault segment (**Fig. 12**).

473     In our database, fault shape approaches a cubic equation in almost all cases; the

474     evolution of fault shape is a result of stress and ruptures in the lithosphere that can be

475     interpreted using the same models that dictate the geometry of failure in the smaller scale and

476     in other materials (Scholz and Aviles, 2013). A discrete fracture developing in a rheological

477     uniform material usually produces a parabolic fault in 2D (Walsh et al., 2002; 2003; Kim and

478     Sanderson, 2005). However, in nature the interaction with varying rock types, adjacent faults

479     and other irregularities within the crust add an order of complexity to fault shapes. This is

480     correctly accounted for with the use of a cubic model (Goff, 1991; Ostertagová, 2012). A

481     second order polynomial is only capable of modelling a curve with a single peak or a single

482     trough. Since the dataset used in this work contains multiple peaks and troughs, such a model

483     is unsuitable; using a third order polynomial overcomes this limitation whereby it can model

484     curves with multiple peaks and troughs.

485     In this fourth step, the regression model developed for fault segment detection is

486     provided with throw data at the maximum resolution possible. However, a set weighting was

487     added for the minima, maxima and Peak throw values, thus ensuring the final curve passes

488     through each of these points. In addition, a lower weighting is given to the peak to prevent

489     the detection of unusual shapes due to other points being ignored by the model. A regression

490     was then applied through the implementation of the python software library Scikit-learn

491     (Grisel et al., 2023), which implements a simple and effective regression algorithm that

492     allows for quick implementation into the code used in previous steps (Grisel et al., 2023;

493     Raschka and Mirjalili, 2018). An advantage of this step is that the resulting curves can model

494     each fault using a single equation, and the computation of such equation simplifies any

495     further analysis needed for a fault (**Fig. 12**). It should be noted that the resulting equation will

496     only give an accurate model of fault shape within the range of the fault's predicted length.

497     Outside this range the cubic equation does not fit with the real fault shape (**Fig. 12**).

498     **5. Critical mathematical tests of minimum sampling rates for T/D and T/Z analyses**

499     A minimum sampling rate for T/D and T/Z analyses was previously estimated by Tao

500     and Alves (2019) as a percentage of the smallest fault segment. They approached the

501     detection of fault linkage points to the mapping of a fault's total area, and geometry, in the

502     2D space. Hence, a downsampling method was gradually applied by Tao and Alves (2019) to

503     data collected at maximum resolutions so to highlight fault linkage points in T/D and T/Z

504     plots. Fault-segment linkage points were detected for each iteration. The number of fault

505     segments was then measured and, once this number was reduced, fault segments could not be

506     detected below a specific sampling rate $\delta$.

507     Mathematically speaking, the standard approach to downsampling a dataset is through

508     decimation, which involves the application of an integer decimation factor $M$. The new

509     decimated data are obtained by simply selecting every $M^{\text{th}}$ value of a signal $x(n)$, a step that

510     returns a new sample rate of:

511

512     $n' = \frac{n}{M}$    (Eq. 4)

513

514    Decimation methods most commonly used involve the application of a low-pass filter prior to

515    decimation so that aliasing is avoided. However, to most accurately simulate the degradation

516    of data that derives from a lower sampling of field measurements, we decided to avoid the

517    application of a low-pass filter to our data. In fact, the decimation approach in Equation 4 is

518    the most basic and allows only for quick tests of the effect of sampling reductions on the

519    shape of T/D and T/Z data.

520        In our analysis, decimation was found to introduce a bias to downsampled data. The

521    results were often determined by the locations of the decimated samples relative to the 'real'

522    linkage points of discrete fault segments. Hence, to allow for a consistent approach to

523    downsampling, an interpolation algorithm was used whereby an interpolation function was

524    generated and followed the input data. This interpolation function was then applied to a new

525    set of sample points, an approach most closely following what happens when acquiring T/D

526    and T/Z data in the field, or in seismic data, as we can choose – by using this interpolation

527    function - a completely new set of sampling points independently of how the original data

528    was acquired.

529        A linear interpolation was therefore followed in our approach by computing two

530    adjacent samples, with the desired sampling point falling between these adjacent sample

531    locations along the throw axis. The normalized spacing between these two samples is *1/U*. If

532    the distance of the first sample comes before the desired sample distance by $x_m$, then the

533    sampling distance of the second sample leads the desired sampling distance by $(1/U) - x_m$.

534    If we designate the two samples as $y_1(m)$ and $y_2(m)$, and use a linear interpolation, the

535    approximation of the desired sample becomes (Proakis, 1992):

536

537

$$y(m) = (1 - a_m)y_1(m) + a_m y_2(m)$$
$$\text{where} \quad a_m = U x_m \quad \text{and} \quad 0 \le a_m \le 1$$

538 (Eq. 5)

539

540 Through this resampling approach, a sampling ratio can be increased and decreased while

541 retaining the original fault shape. The position of sampling points can also be tweaked to find

542 possible sampling intervals that cause information loss.

543

544 *5.1 Integral Error test*

545   The main result of performing a polynomial regression fit is that an interpreted can

546 obtain a discrete equation for each perceived fault. Building upon the method of Modulus

547 Error analysis in Tao and Alves (2019), we created a measure of the scale of changes caused

548 by a reduction in throw sampling. We subtracted the equations of faults measured at different

549 sampling spaces and took the absolute value of the resultant equation, where x is distance:

550

$$\text{Total error} = \frac{\sum_{i=0}^{n} \int_{p_i}^{p_{i+1}} |f_i(x) - g_i(x)|\, dx}{\sum_{i=0}^{n} \int_{p_i}^{p_{i+1}} f_i(x) dx}$$
$$p = \text{Intersection points}$$

551 (Eq. 6)

552

553 For a single fault, Equation 6 can be simplified to:

554

555
$$\text{Fault error} \quad = \frac{\int_p^q |f(x) - g(x)| \, dx}{\int_p^q f(x) dx}$$
$$p = \text{Fault start} \quad \text{and} \quad q = \text{Fault end}$$

556     (Eq. 7)

557

558         Performing an Integral Error calculation on each step of a sampling reduction test

559     reveals some of the effects imposed on the identification of fault-linkage points when one

560     randomizes data (throw) sampling (**Fig. 13**). As the sampling is reduced, the Integral Error

561     increases, responding to the fact that the sampled locations may miss the fault linkage points

562     if the sampling is too coarse. The error will reach a maximum value and then decrease over

563     smaller changes in sampling. This means that coarse and random sampling techniques can

564     drastically change the results, leading to erroneous estimations of fault segments' shape,

565     hindering their subsequent identification. In other words, it is certain that one is overlooking

566     the presence of discrete fault segments when the error starts to decrease in its magnitude (**Fig.**

567     **13**). In addition, when the sampled points are being incrementally reduced along a fault, the

568     distance to the nearest sample may also vary with some degree of randomness. An interpreter

569     may thus be fortunate enough (or not) to collect data near a point where fault segments are

570     linked solely by chance. The influence this has on the error value means that sometimes, but

571     also randomly, error will decrease for a lower number of samples.

572

573     *5.2 Modulus Error test*

574         The approach in Section 5.1 resulted in the calculation of a ratio resolving the size of

575     the error relative to the size of the fault $f(x)$. As we mostly recorded an increase in Integral

576     Error up to the point where a fault is no longer detected, the variation in Integral Error

25

577    became a good indication of the reliability of predictions made at decreasing sampling space.

578    The similarity of this equation to the Modulus Error equation in Tao and Alves (2019) allows

579    for a direct comparison between different error-calculation methods as a function of sampling

580    space:

581

582    $$\text{Modulus Error} = \frac{\sum_1^n |A_m - A'_m|}{\sum_1^n A_m}$$

583    (Eq. 8)

584

585    Taking the integral of Equation 8 will give a value for the area between the two faults, which

586    can be used as a measure of error between two measurements of the same fault zone. In our

587    case it was used to compare the downsampled datasets to the original ground-truth ones as

588    the sampling space is being tested.

589

590    *5.3 Intersection Error test*

591        The lateral tips of discrete fault segments can sometimes change in their relative

592    position (as identified by our algorithm) if data decimation is too 'coarse'. Once again, an

593    interpreter may be fortunate enough (or not) to collect data near a point where fault segments

594    are linked solely by chance. As a result, information is lost; when fault linkage points are not

595    identified in their accurate location, any resulting interpretations of a fault's geometry may be

596    inaccurate. Small changes in the location of fault segments' linkage points may not indicate

597    issues with their identification, so a threshold value needs to be defined if a particular sample

598    strategy is inaccurate.

599    We devised a way to measure the change in intersection points, i.e. the difference

600    between the lateral tip of a fault in one case is compared to the closest lateral tip of a fault in

601    another measurement of the same fault zone. This distance is divided by the length of the

602    fault to give an error value. The average of all the faults' errors gives a final *Intersection*

603    *Error* for the comparison.

604
$$\text{Intersection Error} = (n_{max} - n_{min}) + \sum_{i=1}^{n_{min}} \frac{\min(|a_{i,0} - b_{(1 \to n_{max}),0}|)}{a_{i,2} - a_{i,0}}$$

605
$$n = \text{number of faults} \quad a_{c,d}, b_{c,d} = \text{sequence of faults}$$
$$c = \text{fault number} \quad d = \text{fault start, peak and end}$$

606    (Eq. 9)

607    The Intersection Error returns similar values to the Integral Error. However, it will

608    more clearly identify situations where a fault segment has been overlooked. Other measures

609    of error also prioritize changes in the general shape of faults, while in many cases the more

610    important aspect of the faults we analyzed is where they lateral tips are, i.e. where they begin

611    and end laterally.

612    **6. Discussion**

613

614    *6.1 Downsampling techniques to highlight interpretation errors*

615    A comparison of error percentages when reducing the sampling spacing in T/D and T/Z

616    data reveals some interesting trends (**Fig. 14**). In most cases the error gradually increases

617    when sampling decreases, but there are some examples of minima in Integral and Intersect

618    errors occurring due to a sample coinciding exactly with a fault segment linkage point (see

619    low error percentages in **Fig. 14**). In other words, by simple coincidence, one can select a

27

620    sample that coincides exactly with, or be very close to, a fault intersection point. This finding

621    constitutes an important addition to the analysis of Tao and Alves (2019); it provides further

622    confirmation that obeying a minimum threshold sampling ratio is paramount to analyzing

623    fault segmentation in nature.

624    We applied an iterative downsampling approach to all the fault data available to find a

625    minimum sample ratio as a percentage of fault length. Three (3) approaches were followed to

626    measure minimum sampling ratios from the strictest to the most lenient:

627    a) Strict - sampling considers a percentage of the total data input range, i.e. the total

628    length of a fault zone that is composed of multiple segments,

629    b) Moderate - sampling is calculated considering the longest segment found in a fault

630    zone, and,

631    c) Lenient - sampling only considers the very first fault lost as a result of reducing

632    throw sampling rate.

633    The use of multiple minimum sampling definitions allowed us to identify what are the

634    upper and lower limits for the required sampling ratio in order to map discrete fault segments

635    with accuracy. Our datasets often include a wider range of fault geometries, with faults

636    varying in size along a fault zone. Results are shown in **Fig. 15**.

637    The results show that, with relatively short fault zones, in which only a few faults need

638    to be found and modelled, relatively lenient sample ratios are sufficient when compared to

639    long fault zones. The main caveat of analyzing fault zones is that they may contain long and

640    short fault segments, and the shortest segments need to be accurately identified using strict

641    sampling ratios. This means some fault zone geometries require a much higher sample ratio

642    than, for instance, two-three linked segments with relatively constant sizes.

28

643

*6.2 Minimum sampling ratios in T/D and T/Z analyses*

**Figure 16** illustrates the relationship between each error-testing approach and the

critical sampling ratio, with detailed information being provided in **Table 1**. The purpose of

quantifying error is to understand how much information is lost by a reduction in the

sampling ratio, or distance. The larger the percentage error observed in **Fig. 16**, at a critical

sample ratio, the better the measure of the accuracy of fault predictions is. The critical sample

ratio is the point at which important fault information is lost.

Modulus Error works independently of any fault shape data, so it results in a smaller

distribution error - it cannot reliably tell an interpreter how much information is lost in terms

of fault shapes and their linkage points. In comparison, Integral Error reflects a compromise

between the Modulus and the Intersection errors, though it only returns information on the

accuracy of lateral tips (start and end points) of fault segments. In spite of this, Integral Error

has a much higher average result for error, meaning the changes in fault shape are relatively

greater than the change in position of faults' linkage points.

From these results, and also via the successful visualization of fault shape, we

demonstrate that Integral Error is a superior tool to gauge the loss in information when

comparing variable sampling ratios for faults. The high correlation with Intersection Error

also tells us that there is little use for combining the two error-defining methods (Intersection

and Interval errors) in individual cases, as they are heavily dependent. The Intersection Error

can therefore be used separately to Interval Error as a good indication for how trustworthy the

identification of fault linkage points will be.

After establishing a relationship between error values and critical sample ratios, we

could reach a conclusion on the minimum sample ratios necessary for accurate fault analyses.

667  We found a minimum sample ratio that would be appropriate for various cases, with a 95%

668  success rate (**Table 2**). The success rate measured in these cases is based on the use of a fully

669  automated wavelet method (**Fig. 17** and **Table 2**). With the use of other tools, as well as a

670  human input (fault segment and curvature mapping sensu Kim and Sanderson, 2005) a higher

671  success rate will be likely achieved.

672       The critical values in **Fig. 17** show the minimum sampling ratios calculated for the

673  three downsampling approaches considered in Section 6.1. The Strict approach can be taken

674  as reflecting the minimum sampling length/fault length ratio ($\delta$) for large fault zones

675  comprising fault segments of varied dimensions (see also **Table 2**). These were commonly

676  observed in the datasets gathered in SE Crete where fault-segment length is variable, but with

677  some segments >3.5 km long. Based on these constraints, the point of data loss over a wide

678  range of data sets was calculated in this work and resulted in the estimate of the following $\delta$

679  values:

680  a) A $\delta$ of 1.02% ± 0.02 if one uses a Strict approach for the sampling of throw data. This

681     value is particularly important when in the presence of fault zones that are >3.5 km long,

682  b) A $\delta$ of 4.167% ± 0.18 for a Moderate approach, in which the choice of sampling ratio

683     prioritizes the identification of the longest segment in a fault zone,

684  c) A minimum $\delta$ of 5.882% ± 1.26 is necessary to identify segments in a fault zone using a

685     Lenient sampling approach.

686

687  For a typical fault zone that is longer than 20 km, such as Ierapetra's with its largest segments

688  c. 3.5 km long, the results above indicate that the collection of throw values every 35 m is the

689  minimum sampling ratio one should use. In 3D seismic data, this translates into mapping

30

690    fault throws every 3 lines for a typical volume with a bin spacing of 12.5 m. Moderate and

691    Lenient approaches will respectively translate into the collection of throw data every 140 m

692    and 200 m along the Ierapetra Fault, i.e. every 11 and 16 lines for a similar fault in a 3D

693    seismic volume processed with a bin spacing of 12.5 m. In SW England, sub-seismic faults

694    are 1.65 m to 7.55 m long, and that results in a Strict sampling that varies from 1.68 cm to 7.7

695    cm. A more Lenient sampling would require throws sampled every 9.7 cm and 44.39 cm for

696    such structures.

697         It is worth noting these are not prescriptive sampling distances as, recognizing, the

698    minimum sampling length/fault length ratio ($\delta$) is a function on fault length. Moreover, this

699    same rule also applies to the collection of throw data for T/Z (throw-depth) plots so to

700    prevent the grouping of distinct segments into a single unlinked (coherent) fault.

701

702    *6.3 Implications for T/D and T/Z analyses*

703

704         Ze and Alves (2019) recognized that depositional rates near active normal faults vary

705    significantly on their hanging-wall and footwall blocks, as well as recording variable

706    sediment pathways. This renders the use of expansion indexes and layer-by-layer

707    interpretations of throw troublesome in seismic data imaging relatively old, buried basins.

708    The Strict approach to using a $\delta$ of 1.02% will compensate for any of the issues indicated in

709    Tao and Alves (2019), helping in the identification of early-stage fault segmentation. It will

710    prevent the tendency, in the published literature, of considering the constant-length model as

711    predominant in nature. In order to reduce risk of important data loss in the interpretation of

712    short, minor faults, we recommend the use of a $\delta$ value of 1.0% preventing the loss of

31

713    important fault information. Taking the smallest fault in the area as the reference point for a

714    δ value also gives less room for interpretation error.

715        A limitation concerning the use of T/D and T/Z data in fault analyses is that the scale at

716    which structural geologists acquire and interpret fault throw (or displacement) data is

717    variable. It depends on the inherent scale of the structures of interest, and the aims of the

718    survey or study in question. The chosen scale of observation is also dependent on data

719    resolution and pre-defined structural criteria (e.g. Walsh and Watterson, 1991, Walsh et al.,

720    2002, Walsh et al., 2003, Kim and Sanderson, 2005, Torabi and Berg, 2011). Therefore, to

721    acquire data at a scale that is several orders of magnitude greater than that in which fault

722    segmentation likely occurred, e.g. interpreting deeply buried faults using seismic data of

723    poorer quality persuades interpreters to readily recognize coherent fault-growth models to the

724    detriment of the isolated growth model. This is particularly the case when faults crossing

725    sedimentary basins, but not rooted into basement units (and, therefore, not developed at a

726    crustal scale), are interpreted in seismic data. At what temporal scale is the 'fast-propagation',

727    coherent fault model applicable is another important caveat in many of these models – the

728    time-dependent growth and ultimate linkage of small faults is not easily resolved in seismic

729    data, nor are stratigraphic (age) constraints often accurate enough. For these reasons, we

730    consider that fault segmentation can be systematically overlooked by interpreters when

731    adopting of broad, one-fits-all, attitude to data sampling, against which the Strict δ values

732    suggested in this work should be used in fault analyses, but rarely are.

733    **7. Conclusions**

734        This work shows that the application of a Wavelet-Transform detection system in fault

735    analysis is useful to automate fault mapping and remove human bias from interpretation

736    workflows. With human oversight and adjustments, this system improves the productivity of

32

737    interpreters analyzing complex fault arrays. As a corollary, this work proves the need to

738    consider a threshold sampling ratio ($\delta$) in T/D and T/Z data as necessary, based on the

739    following results:

740        a) A lower sampling ratio ($\delta$) is required when interpreting long, segmented fault zones

741    composed of faults of multiple lengths and heights. This is important as the linkage points

742    between fault segments often coincide with regions of throw minima that are much smaller

743    than the throw maxima of adjacent faults. The adoption of a low sampling ratio is

744    independent of the style of linkage between discrete fault segment, e.g. hard-linkage, soft-

745    linkage, or relay ramps. It is also independent of the type of fault one considers (normal,

746    reverse or strike-slip).

747        b) This work suggests a minimum sampling ratio ($\delta$) of 4.167% for faults that are

748    relatively short, and clearly isolated. This is, however, a rough guideline, as faults in nature

749    can have some unpredictable geometries and a Strict approach ($\delta$ of 1.02% $\pm$ 0.02) may still

750    be the appropriate, in many instances, if recognizing fault segmentation is the main aim of a

751    study.

752        c) For the fault zones we analyzed in the field and at outcrop, a Strict sampling ratio of

753    1.02% will translate into throw data collected every 35 m if a fault zone contains segments

754    greater than 3.5 km. In 3D seismic data, this translates into mapping fault throws every 3

755    lines for a typical volumes with a bin spacing of 12.5 m. Moderate and Lenient approaches to

756    fault measurements will respectively translate into the collection of throw data every 140 m

757    and 200 m for such a fault zone geometry. The smaller sub-seismic faults of SE England

758    require a sampling every 1.65 cm (Strict approach) to 44.39 cm (Lenient approach).

759        d) The final decision regarding the use of Strict sampling ratios of 1.02% $\pm$ 0.02 should

760    be based on all geological information available on the fault zone, or region, being analyzed.

33

761  If there is any major uncertainty around fault size, one should follow a Strict approach and

762  consider a $\delta$ of 1.02% $\pm$ 0.02.

763      e) Mathematically speaking, a combination of Continuous Wavelet Transforms and

764  Polynomial Regressions allows for an accurate mapping of fault segmentation from T/D and

765  T/Z data. The Continuous Wavelet Transform is used to define fault ranges. A cubic

766  (polynomial) regression model is later applied on these ranges to obtain fault shape in a

767  separate stage. The high reliability of this technique allows for its systematic application

768  using Machine Learning tools.

769      The results in this work are based on mathematical methods tested on a large dataset

770  comprising 415 faults. The method we propose are applied with minimal human intervention,

771  meaning results can be directly linked to the mathematical equations. The results also

772  demonstrate the significant impact data sampling techniques can have on the resulting

773  interpretation of fault location, and growth modes, particularly whenever small faults are

774  quickly lost due to sub-scale imaging or incorrect measuring approaches. Significant changes

775  to the perception of the entire fault zone can be seen when a single fault becomes

776  indistinguishable. For these reasons, we recognize that fault segmentation is systematically

777  overlooked in the published literature when adopting a broad, one-fits-all, attitude to data

778  sampling, against which the $\delta$ values suggested in this work should be used.

779

780  **Acknowledgements**

785  (Univ. Canterbury, NZ) for the provision of an editable version of Fig. 2 and M. Glen for

786  providing the map of Kilve and Watchet in Fig. 4d. The authors also thank editors I. Alsop

787  and S. Mukherjee, together with two anonymous reviewers, for their constructive comments

788  to an earlier draft of this paper.

789

790  **Data statement**

791  The data that support the findings of this study are available on request from the

792  corresponding author, Tiago M. Alves.

793

794  **References**

795  Alves, T.M., 2024. Networks of geometrically coherent faults accommodate Alpine tectonic

796      inversion offshore SW Iberia. Solid Earth, 15, 39–62. https://doi.org/10.5194/se-15-39-

797      2024

798  Alves, T.M., and Cupkovic. T., 2018. Footwall Degradation Styles and Associated

799      Sedimentary Facies Distribution in SE Crete: Insights into Tilt-Block Extensional

800      Basins on Continental Margins. Sedimentary Geology, 367, 1–19.

801      https://doi.org/10.1016/j.sedgeo.2018.02.001

802  Alves, T.M., Mattos, N.H., Newnes, S., and Goodall, S., 2022. Analysis of a Basement Fault

803      Zone with Geothermal Potential in the Southern North Sea. Geothermics 102, 102398.

804      https://doi.org/10.1016/j.geothermics.2022.102398

805  Aydin, A., Ozbek, A., and Cobanoglu, I., 2004. Tunnelling in Difficult Ground: A Case

806      Study from Dranaz Tunnel, Sinop, Turkey. Engineering Geology, 74, 293–301.

807      https://doi.org/10.1016/j.enggeo.2004.04.003

808    Baudon, C., and Cartwright, J., 2008. The kinematics of reactivation of normal faults using

809        high resolution throw mapping. Journal of Structural Geology, 30, 1072-1084.

810        https://doi.org/10.1016/j.jsg.2008.04.008

811    Caputo, R., Catalano, S., Monaco, C., Romagnoli, G., Tortorici, G., and Tortorici, L., 2010.

812        Active Faulting on the Island of Crete (Greece). Geophysical Journal International, 183,

813        111–26. https://doi.org/10.1111/j.1365-246X.2010.04749.x

814    Cartwright, J., Bouroullec, R., James, D., and Johnson, H., 1998. Polycyclic motion history of

815        some Gulf Coast growth faults from high-resolution displacement analysis. Geology,

816        26, 819-822. https://doi.org/10.1130/0091-7613(1998)026<0819:PMHOSG>2.3.CO;2

817    Cloetingh, S., Van Wees, J.D., Ziegler, P.A., Lenkey, L., Beekman, F., Tesauro, M., Förster,

818        A., Norden, B., Kaban, M., and Hardebol, N., 2010. Lithosphere Tectonics and

819        Thermo-Mechanical Properties: An Integrated Modelling Approach for Enhanced

820        Geothermal Systems Exploration in Europe. Earth-Science Reviews, 102, 159–206.

821        https://doi.org/10.1016/j.earscirev.2010.05.003

822    Dart, C.J., McClay, K., and Hollings, P.N., 1995. 3D analysis of inverted extensional fault

823        systems, southern Bristol Channel basin, UK. In: Buchanan, J.G. and Buchanan, P.G.

824        (Eds). *Basin Inversion*, Geological Society of London, Special Publications, 88, pp.

825        393-413. https://doi.org/10.1144/GSL.SP.1995.088.01.21

826    Fernández-Blanco, D., de Gelder, G., Lacassin, R., and Armijo, R., 2019. Geometry of

827        flexural uplift by continental rifting in Corinth, Greece. Tectonics, 39, e2019TC005685.

828        https://doi.org/10.1029/2019TC005685

829    Ferrill, D.A., Morris, A.P., McGinnis, R.N., Smart, K.J., Wigginton, S.S., and Hill, N.J.,

830        2017. Mechanical Stratigraphy and Normal Faulting. Journal of Structural Geology, 94,

831        275–302. https://doi.org/10.1016/j.jsg.2016.11.010

832     Ferrill, D.A., Samrt, K.J., and Morris, A.P., 2020. Fault failure modes, deformation

833          mechanisms, dilation tendency, slip tendency, and conduits v. Seals. *In:* Ogilvie, S.R.,

834          Dee, S.J., Wilson, R.W., Baileu, W.R. (Eds), Integrated Fault Seal Analysis. Geological

835          Society, London, Special Publications, 496, pp. 75-98. https://doi.org/10.1144/SP496-

836          2019-7

837     Gaki-Papanastassiou, K., Karymbalis, E., Papanastassiou, D., and Maroukian, H. 2009.

838          Quaternary Marine Terraces as Indicators of Neotectonic Activity of the Ierapetra

839          Normal Fault SE Crete (Greece). Geomorphology, 104, 38–46.

840          https://doi.org/10.1016/j.geomorph.2008.05.037

841     Glen, R.A., Hancock, P.L., and Whittaker, A., 2005. Basin inversion by distributed

842          deformation: the southern margin of the Bristol Channel Basin, England. Journal of

843          Structural Geology, 27, 2113-2134. https://doi.org/10.1016/j.jsg.2005.08.006

844     Gloaguen, R., Marpu, P.R., and Niemeyer, I., 2007. Automatic extraction of faults and fractal

845          analysis from remote sensing data. Nonlinear Processes in Geophysics, 14, 131–138,

846          https://doi.org/10.5194/npg-14-131-2007

847     Goff, J.A., 1991. A Global and Regional Stochastic Analysis of Near-Ridge Abyssal Hill

848          Morphology. Journal of Geophysical Research: Solid Earth, 96, 21713–21737.

849          https://doi.org/10.1029/91JB02275

850     Grisel, O., Mueller, A., Lars, Gramfort, A., Louppe, G., Fan, T.J., Prettenhofer, P., et al.,

851          2023. Scikit-Learn/Scikit-Learn: Scikit-learn 1.3.0. Zenodo.

852          https://doi.org/10.5281/ZENODO.8098905

853     Gudmundsson, A., 2012. *Rock Fractures in Geological Processes*. Cambridge University

854          Press, Cambridge, 578 pp.

855       https://doi.org/10.1017/CBO9780511975684

856   Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D.,

857       Wieser, E., et al., 2020. Array Programming with NumPy. Nature, 585, 357–62.

858       https://doi.org/10.1038/s41586-020-2649-2

859   He, D., Lu, R., Huang, H., Wang, X., Jiang, H., and Zhang, W., 2019. Tectonic and

860       Geological Setting of the Earthquake Hazards in the Changning Shale Gas

861       Development Zone, Sichuan Basin, SW China. Petroleum Exploration and

862       Development, 46, 1051–64. https://doi.org/10.1016/S1876-3804(19)60262-4

863   Heiberger, R.M., and Neuwirth, E., 2009. Chapter 11 – Polynomial Regression. In:

864       Heiberger, R.M. and Neuwirth, E. (Eds), R Thorugh Excel – A spreadsheet interface for

865       statistics, data analysis, and graphics.  Springer Verlag New York, pp. 269-284.

866       https://doi.org/10.1007/978-1-4419-0052-4

867   Huenges, E., Kohl, T., Kolditz, O., Bremer, J., Scheck-Wenderoth, M., and Vienken, T.,

868       2013. Geothermal Energy Systems: Research Perspective for Domestic Energy

869       Provision. Environmental Earth Sciences, 70, 3927–33. https://doi.org/10.1007/s12665-

870       013-2881-2

871   James, G., Witten, D., Hastie, T., and Tibshirani, R. 2013. Chapter 7 - Moving Beyond

872       Linearity. In: James, G., Witten, D., Hastie, T., Tibshirani, R. (Eds.), An Introduction to

873       Statistical Learning: With Applications in R, Springer Texts in Statistics, 103, New

874       York: Springer, pp. 265–301.

875   Jentsch, A., Jolie, E., Jones, D.G., Taylor-Curran, H., Peiffer, L., Zimmer, M., and Lister, B.,

876       2020. Magmatic Volatiles to Assess Permeable Volcano-Tectonic Structures in the Los

877       Humeros Geothermal Field, Mexico. Journal of Volcanology and Geothermal

878       Research, 394, 106820, https://doi.org/10.1016/j.jvolgeores.2020.106820

879   Kalbermatten, M., van de Ville, D., Turberg, P., Tuia, D. and Joost. S. 2012. Multiscale

880       analysis of geomorphological and geological features in high resolution digital

881    elevation models using the wavelet transform. Geomorphology, 138, 352-363.

882    https://doi.org/10.1016/j.geomorph.2011.09.023

883    Kim, Y.-S, and Sanderson, D.J., 2005. The Relationship Between Displacement and Length

884    of Faults: A Review. Earth-Science Reviews, 68, 317–334.

885    https://doi.org/10.1016/j.earscirev.2004.06.003

886    King, J.J., and Cartwright, J.A., 2020. Ultra-slow throw rates of polygonal fault systems.

887    Geology, 48, 473-477. https://doi.org/10.1130/G47221.1

888    Kozłowska, M., Brundzinski, M.R., Friberg, P., and Currie, B.S., 2017. Maturity of nearby

889    faults influences seismic hazard from hydraulic fracturing. Proceedings of the National

890    Academy of Sciences, 115, E1720-E1729. https://doi.org/10.1073/pnas.171528411

891    Laubach, S.E., Lamarche, J., Gauthier, B.D.M., Dunne, W.M., Sanderson, D.J., 2018. Spatial

892    arrangement of faults and opening-mode fractures. Journal of Structural Geology, 108,

893    2-15. https://doi.org/10.1016/j.jsg.2017.08.008

894    Lee, G., Gommers, R., Wohlfahrt, K., Wasilewski, F., O'Leary, A., Nahrstaedt, H., Sauvé, A.

895    et al. 2022. PyWavelets/Pywt: V1.4.1. Zenodo.

896    https://doi.org/10.5281/ZENODO.1407171

897    Lee, J., Stockli, D.F., and Blythe, A.E., 2023. Cenozoic slip along the southern Sierra Nevada

898    normal fault, California (USA): A long-lived stable western boundary of the Basin and

899    Range. Geosphere, 19, 878-899. https://doi.org/10.1130/GES02574.1

900    Mallat, S. G. 2009. *A Wavelet Tour of Signal Processing*. 2nd ed. San Diego: Academic

901    Press, 895 pp. https://doi.org/10.1016/B978-0-12-374370-1.X0001-8

902    Mechernich, S., Reicherter, K., Deligiannakis, G., and Papanikolaou, I., 2023. Tectonic

903    geomorphology of active faults in Eastern Crete (Greece) with slip rates and earthquake

904    history from cosmogenic $^{36}$Cl dating of the Lastros and Orno faults. Quaternary

905    International, 652, 77-91. https://doi.org/10.1016/j.quaint.2022.04.007

906    Misra AA, and Mukherjee S. 2018. Atlas of Structural Geological Interpretation from

907        Seismic Images. Wiley Blackwell. ISBN: 978-1-119-15832-5.

908    Moska, R., Labus, K., and Kasza, P., 2021. Hydraulic Fracturing in Enhanced Geothermal

909        Systems—Field, Tectonic and Rock Mechanics Conditions—A Review. Energies, 14,

910        5725. https://doi.org/10.3390/en14185725

911    Mukherjee S., 2019. Particle tracking in ideal faulted blocks using 3D co-ordinate geometry.

912        Marine and Petroleum Geology 107, 508-514.

913        https://doi.org/10.1016/j.marpetgeo.2019.05.037

914    Nicol, A., Walsh, J., Childs, C, and Manzocchi, T., 2020. Chapter 6 - The growth of faults.

915        *In:* Understanding Faults – Detecting, Dating and Modeling (Eds: Tanner, D, Brandes),

916        Elsevier, pp. 221-255. https://doi.org/10.1016/B978-0-12-815985-9.00006-0

917    Nixon, C.W., McNeill, L.C., Gawthorpe, R.L., Shillington, D.J., Michas, G., Bell, R.E.,

918        Moyle A., Ford, M., Zakharova, N.V., Bull, J.M., and de Gelder, G., 2024. Increasing

919        fault slip rates within the Corinth Rift, Greece: A rapidly localising active rift fault

920        network. Earth and Planetary Science Letters, 636, 118716.

921        https://doi.org/10.1016/j.epsl.2024.118716

922    Ostertagová, E., 2012. Modelling Using Polynomial Regression. Procedia Engineering, 48,

923        500–506. https://doi.org/10.1016/j.proeng.2012.09.545

924    Peacock, D.C.P., Nixon, C.W., Rotevatn, A., Sanderson, D.J., and Zuluaga, L.F., 2017.

925        Interacting faults. Journal of Structural Geology, 97, 1-22.

926        https://doi.org/10.1016/j.jsg.2017.02.008

927    Plawiak, R.A.B., Carvalho, M.J., Sombra, C.L., Brandão, D.R., Mepen, M., Ferrari, A.L., and

928        Gambôa, L.A.P., 2024. Structural controls of the migration of mantle-derived $CO_2$

929        offshore in the Santos Basin (Southeastern Brazil). Frontier in Earth Sciences, 11,

930        https://doi.org/10.3389/feart.2023.1284151

931 Pollard, D.D. and Segall, P., 1987. Theoretical displacements and stresses near fractures in

932     rock: with applications to faults, joints, veins, dikes, and solution surfaces. In:

933     Atkinson, B.K. (Ed*.), Fracture Mechanics of Rock*, Academic Press, London, pp. 277-

934     349.

935 Proakis, J.G., 1992. Multirate Digital Signal. In: Proakis, J.G. (Ed.), Advanced Digital Signal

936     Processing, Macmillan, 141–204.

937     https://books.google.co.uk/books?id=4_RSAAAACAAJ

938 Purba, D.P, Adityatama, D.W., Umam, M.F., and Muhammad, F., 2019. Key Considerations

939     in Developing Strategy for Geothermal Exploration Drilling Project in Indonesia.

940     Proceedings, 44th Geothermal Reservoir Engineering, Stanford University, Stanford,

941     California, SGP-TR-214.

942 Raschka, S., and Mirjalili, V., 2018. Predicting Continuous Target Variables with Regression

943     Analysis." In: Raschka, S., Mirjalili, V. (Eds.), Python Machine Learning: Machine

944     Learning and Deep Learning with Python, Scikit-Learn, and TensorFlow, Packt

945     Publishing, Second edition, fourth release, pp. 309–346.

946 Rawnsley, K.D., Peacock, D.C.P., Rives, T., and Petit, J.P., 1998. Joints in the Mesozoic

947     sediments around the Bristol Channel Basin. Journal of Structural Geology, 20, 1641-

948     1661. https://doi.org/10.1016/S0191-8141(98)00070-4

949 Robertson, J., Roberts, G.P., Iezzi, F., Meschis, M., Gheorghiu, D.M., Sahy, D., Bristow, C.,

950     and Sgambato, C., 2020. Distributed normal faulting in the tip zone of the South

951     Alkyonides Fault System, Gulf of Corinth, constrained using $^{36}$Cl exposure dating of

952     late-Quaternary wave-cut platforms. Journal of Structural Geology, 136, 104063.

953     https://doi.org/10.1016/j.jsg.2020.104063

954      Rotevatn, A., Jackson, C.A.-L., Tvedt, A.B.M., Bell, R.E., and Blækkan, I., 2019. How do

955          normal faults grow? Journal of Structural Geology, 125, 174-184.

956          https://doi.org/10.1016/j.jsg.2018.08.005

957      Saeidi, O., Rasouli, V., Vaneghi, R.G., Gholami, R., and Torabi, S.R., 2014. A modified

958          failure criterion for transversely isotropic rocks. Geoscience Frontiers, 5, 215-225.

959          https://doi.org/10.1016/j.gsf.2013.05.005

960      Shen, S., Li, H., Chen, W., Wang X., and Huang, B., 2022. Seismic Fault Interpretation

961          Using 3-D Scattering Wavelet Transform CNN, IEEE Geoscience and Remote Sensing

962          Letters, 19, 1-5, Art no. 8028505. doi: 10.1109/LGRS.2022.3183495

963      Scholz, C.H., and Aviles, C.A., 2013. The Fractal Geometry of Faults and Faulting. In: Das,

964          S., Boatwright, J., Scholz, C.H. (Eds), Geophysical Monograph Series, American

965          Geophysical Union, Washington, D.C., 147–155. https://doi.org/10.1029/GM037p0147

966      Sifuzzaman, M., Islam, M.R., and Ali, M.Z., 2009. Application of Wavelet Transform and its

967          advantages compared to Fourier Transform. Journal of Physical Sciences, 13, 121-134.

968          ISSN: 0972-8791

969      Tao, Z., and Alves, T.M., 2019. Impacts of Data Sampling on the Interpretation of Normal

970          Fault Propagation and Segment Linkage. Tectonophysics, 762, 79–96.

971          https://doi.org/10.1016/j.tecto.2019.03.013

972      Torabi, A., Alaei, B., and Libak. A., 2019. Normal Fault 3D Geometry and Displacement

973          Revisited: Insights from Faults in the Norwegian Barents Sea. Marine and Petroleum

974          Geology, 99, 135–55. https://doi.org/10.1016/j.marpetgeo.2018.09.032

975      Torabi, A., and Berg, S.S., 2011. Scaling of fault attributes: A review. Marine and Petroleum

976          Geology, 28, 1444-1460. https://doi.org/10.1016/j.marpetgeo.2011.04.003

977    Torabi, A., Rudnicki, J., Alaei, B., and Buscarnera, G., 2023. Envisioning faults beyond the

978         framework of fracture mechanics. Earth-Science Reviews, 238, 104358.

979         https://doi.org/10.1016/j.earscirev.2023.104358

980    Trippetta, F., Petricca, P., Billi, A., Collettini, C., Cuffaro, M., Lombardi, A. M., Scrocca, D.,

981         Ventura, G., Morgante, A., and Doglioni, C., 2019. From mapped faults to fault-length

982         earthquake magnitude (FLEM): a test on Italy with methodological implications. Solid

983         Earth, 10, 1555–1579. https://doi.org/10.5194/se-10-1555-2019

984    Tvedt, A.B.M., Rotevatn, A., Jackson, C.A.-L., Fossen, H., and Gawthorpe, R.L., 2013.

985         Growth of normal faults in multilayer sequences: A 3D seismic case study from the

986         Egersund Basin, Norwegian North Sea. Journal of Structural Geology, 55, 1-20.

987         https://doi.org/10.1016/j.jsg.2013.08.002

988    Underhill, J.R. and Paterson, S. 1998. Genesis of tectonic inversion structures: seismic

989         evidence for the development of key structures along the Purbeck–Isle of Wight

990         Disturbance. Journal of Geological Society, London, 155, 975-992.

991         https://doi.org/10.1144/gsjgs.155.6.0975

992    Varela. C.L., and Mohriak, W.U., 2013. Halokinetic rotating faults, salt intrusions, and

993         seismic pitfalls in the petroleum exploration of divergent margins. American

994         Association of Petroleum Geologists Bulletin, 97, 1421-1446.

995         https://doi.org/10.1306/02261312164

996    Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D.,

997         Burovski, E., et al., 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing

998         in Python. Nature Methods, 17, 261–72. https://doi.org/10.1038/s41592-019-0686-2

999    Walsh, J., Nicol, A., Childs, C., 2002. An alternative model for the growth of faults, Journal

1000        of Structural Geology, 24, 1669-1675. https://doi.org/10.1016/S0191-8141(01)00165-1

1001    Walsh, J.J., Bailey, W.R., Childs, C., Nicol, A., and Bonson, C.G., 2003. Formation of

1002        Segmented Normal Faults: A 3-D Perspective. Journal of Structural Geology, 25,

1003        1251–1262. https://doi.org/10.1016/S0191-8141(02)00161-X

1004    Walsh, J.J., and Watterson, J., 1991. Geometric and Kinematic Coherence and Scale Effects

1005        in Normal Fault Systems. Geological Society, London, Special Publications, 56, 193–

1006        203. https://doi.org/10.1144/GSL.SP.1991.056.01.13

1007    Wang, Y., 2015a. Frequencies of the Ricker wavelet. Geophysics, 80, 1MA-Z50.

1008        https://doi.org/10.1190/geo2014-0441.1

1009    Wang, Y., 2015b. Generalized seismic wavelets. Geophysical Journal International, 203,

1010        1172-1178. https://doi.org/10.1093/gji/ggv346

1011    Wang, Y., Li, J., Wang, Z.-F., and Chang, H., 2022. Structural Failures and Geohazards

1012        Caused by Mountain Tunnel Construction in Fault Zone and Its Treatment Measures: A

1013        Case Study in Shaanxi. Engineering Failure Analysis, 138, 106386.

1014        https://doi.org/10.1016/j.engfailanal.2022.106386

1015    Yielding, G., 2015. Trapping of buoyant fluids in fault-bound structures. *In:* Industrial

1016        Structural Geology: Principles, Techniques and Integration (Eds: Richards, F.L.,

1017        Richardson, N.J., Rippington, S.J., Wilson, R.W., Bond, C.E.). Geological Society,

1018        London, Special Publications, 421, 29-39. https://doi.org/10.1144/SP421

1019    Ze, T., and Alves, T.M., 2016. The Role of Gravitational Collapse in Controlling the

1020        Evolution of Crestal Fault Systems (Espírito Santo Basin, SE Brazil). Journal of

1021        Structural Geology, 92, 79–98. https://doi.org/10.1016/j.jsg.2016.09.011

1022    Zhang, Q., Alves, T.M., and Martins-Ferreira, M.A.C., 2022. Fault Analysis of a Salt

1023        Minibasin Offshore Espírito Santo, SE Brazil: Implications for Fluid Flow, Carbon and

1024        Energy Storage in Regions Dominated by Salt Tectonics. Marine and Petroleum

1025        Geology, 143, 105805. https://doi.org/10.1016/j.marpetgeo.2022.105805

**Figure and table captions**

1026

1027

1028 Figure 1 – Schematic representation of how tectonic faults interact and link in nature. Faults

1029 evolve from isolated to interacting faults by linking vertically and laterally. The ratio of

1030 $d_{max}/L$ (maximum displacement vs. length) increases as lateral propagation occurs in a fault.

1031 Stage 1 corresponds to the formation of isolated, non-interacting fault segments. Stage 2

1032 relates to the start of fault interaction, overlap and joint growth. Stage 3 represents a fully

1033 linked pair of faults that grow together from that moment onwards. Figure is modified from

1034 Kim and Sanderson (2005).

1035

1036 Figure 2 – Schematic representation of normal-fault evolution. Isolated propagating faults

1037 (left) consist of isolated segments that coalesce to form long, interlinked fault strands. The

1038 coherent constant-length growth model (right) assumes that lateral fault propagation is rapid

1039 but vertical propagation is limited. Figure modified from Nicol et al. (2020).

1040

1041 Figure 3 – Diagram summarizing the way fault-throw data are measured at outcrop, or using

1042 stratigraphic markers in seismic data. The diagram is modified from Ze and Alves (2019) and

1043 based on the Ierapetra Fault Zone, SE Crete, one of the faults analyzed in this work. Throw

1044 measurements are usually taken relative to a correlative surface that is present on the footwall

1045 and hanging-wall blocks of faults. However, this can be made difficult by fault scarp erosion,

1046 and by the covering of the immediate hanging-wall depocentre to the fault by strata. Heave

1047 corresponds to the lateral displacement accommodated by a fault during its movement. Fault

1048 displacement is the resultant vector of throw and heave.

1049

1050   Figure 4 – a) World map indicating the location of the regions where T/D and T/Z data were

1051   acquired for this study. b) Location of the seismic surveys interpreted in SE Brazil from

1052   which fault-throw data were acquired. c) Location of the Ierapetra Fault relative to other fault

1053   families, local sedimentary basins and regional basement terrains. d) Map of SW England's

1054   coast highlighting the locations where fault-throw data were acquired at a sub-seismic scale

1055   (see black squares on the map). Figure 4b is modified from Alves and Cupkovic (2018).

1056   Figure 4d is modified from Glen et al. (2005).

1057

1058   Figure 5 – Examples of faults analyzed in this work, from where throw measurements were

1059   acquired. a) Some of the salt-related faults at the scale of industry seismic data acquired from

1060   a high-resolution seismic survey shot in SE Brazil. b) Panoramic view of the central part of

1061   the Ierapetra Fault Zone and its constituting fault segments. In the parentheses are shown the

1062   height of footwall blocks associated with what is a > 25 km long normal fault zone. c) and d)

1063   Faults in the SW England (Bristol Channel) at the sub-seismic scale.

1064

1065   Figure 6 – Normalized Ricker wavelet, a symmetrical wavelet used to represent signal

1066   changes in the time domain Wang, 2015a, 2015b). In this work, the time domain was

1067   replaced with by a spatial component (length or height) in order to apply the *Ricker* wavelet

1068   theory to the identification of fault segments.

1069

1070   Figure 7 – Graphical example of the Continuous Wavelet Transform technique used to

1071   identify discrete fault segments (Step 1 in this work, Section 4.1) at the lower polynomial

1072   degree 3. Note the obvious correlation between frequency band strength and the throw

1073   maxima recorded for each fault segment. Note that fault segmentation using this technique

46

1074 results in the smallest segments being ignored by the algorithm. This figure thus stress the

1075 fact that a Continuous Wavelet Transform cannot identify throw maxima in the smaller fault

1076 segments – it is focused on picking the greatest throw maxima in a given T/D and T/Z

1077 dataset.

1078

1079 Figure 8 – Workflow suggested in this paper for the identification of fault segments using a

1080 Machine Learning approach.

1081

1082 Figure 9 – Example of the improved fault recognition resulting from applying gradient

1083 measurements from the point of threshold minima (Step 2 in this work, Section 4.2). Step 2

1084 focused on finding the nearest throw minimum representing the linkage point between two

1085 fault segments. Threshold values can be changed in the algorithm, with a stricter threshold

1086 resulting in the identification of only the larger fault segments, and a looser threshold

1087 resulting in multiple fault segments being found. If no frequency minimum is found before

1088 reaching the end of the dataset, the last value picked by the algorithm is taken as the end of

1089 the segment. In Step 2, some of the smallest fault segments were still overlooked by the

1090 algorithm but not on such a scale as revealed in Step 1 (see **Figs. 7 and 8**).

1091

1092 Figure 10 – Graph used to estimate noise floor in the data used in this work. The rapid

1093 descent recorded with increasing wavelength sizes represents the reduction in noise occurring

1094 as a result, as small changes in throw are filtered out by the algorithm. Once this noise is

1095 filtered out, and the curve approaches a flat, we can be confident that the remaining data is

1096 accurate. Peak rate values, when plotted against the frequency of data, show that adopting a

1097 threshold peak rate of 0.04 is a valid approach.

1098

1099    Figure 11 – Example of the improved fault recognition after applying a peak rate threshold to

1100    a Continuous Wavelet Transform (Step 3 in this work, Section 4.3). a) Fault R2_H3

1101    interpreted in high-resolution seismic data from SE Brazil. b) Fault L2 H4-1 from onshore SE

1102    Brazil. Note the improved results in Step 3 when compared with Step 2, but with some

1103    smaller peaks being still overlooked in parts of the fault segments analyzed. The adoption of

1104    a 0.04 peak rate (see **Fig. 10**) returned positive results in Step 3 - all Peaks that are clearly not

1105    part of discrete segments were ignored, without overlooking any possible faults.

1106

1107    Figure 12 - Examples of regression curves modelling fault shape in T/D and T/Z data using a

1108    cubic model (Approach 4). Overall, this was the method that returned a better correlation

1109    between the fault segments identified in our dataset and the segments identified by the

1110    algorithm used. a) Fault R2 H2 analyzed from high-resolution 3D seismic data from SE

1111    Brazil. b) Segmented fault zone R2 H3 interpreted in SE Brazil using high-resolution seismic

1112    data. c)  Fault L2 H4-1 from offshore SE Brazil.

1113

1114    Figure 13 – a) Visualization of T/D plots before and after a critical sampling ratio is applied.

1115    b) Example of the changes in fault shape when sampling ratio is reduced to an Integral Error

1116    of 11.6%.

1117

1118    Figure 14 – Change in error rate observed while the number of samples is reduced. a) In

1119    keystone fault 6-11, Modulus Error increases at a constant rate, whereas integral and intersect

1120    errors vary erratically due to loss of fault intersection points. b) Fault C24 records a rapid

1121    oscillation of error values is recorded. In most cases the error gradually increases when

1122    decreasing the sampling ratio, but there are some examples of minima in Integral and

1123    Intersect errors occurring due to a sample coinciding exactly with a fault segment linkage

1124    point (see Section 5.1 in this article).

1125

1126    Figure 15 – Total distribution of minimum sample ratios ($\delta$) for all datasets in this work.

1127    Results are shown separately for three different downsampling approaches: Strict, Moderate,

1128    and Lenient (see Section 6.2 in this article).

1129

1130    Figure 16 – Error distribution after a critical sampling ratio is applied to the throw data in this

1131    work.

1132

1133    Figure 17 – Graph showing the minimum sampling ratio calculated for Strict, Moderate and

1134    Lenient approaches to T/D and T/Z sampling. The sampling ratio ($\delta$) values corresponding to

1135    a 95% success rate in fault-segment recognition are highlighted, with each data point

1136    represented by a vertical line.

1137

1138    Table 1 - Key statistics concerning the box plot in Fig. 16.

1139

1140    Table 2 – Minimum sampling ratios ($\delta$) calculated based on a 95% success rate in fault-

1141    segment recognition for each downsampling approach: Strict, Moderate, and Lenient. See

1142    **Fig. 17** for a graphical representation of these values.

Table 1 - Key statistics concerning the box plot in Fig. 16.

| Error Type | Min. | Q1 | Median | Mean | Q3 | Upper | Max. |
|---|---|---|---|---|---|---|---|
| Integral | 2.48% | 55.5% | 89.5% | 71.8% | 96.8% | 99.97% | 99.97% |
| Modulus | 0.401% | 4.50% | 6.72% | 11111% | 10.1% | 17.7% | 49.1% |
| Intersection | 0% | 21.0% | 40.4% | 36.6% | 59.7% | 86.0% | 86.0% |
| Reduction (%) | 10% | 19.5% | 27.3% | 31.7% | 52.9% | 90.4% | 90.4% |

Table 2 – Minimum sampling ratios ($\delta$) calculated based on a 95% success rate in fault-segment recognition for each downsampling approach: Strict, Moderate, and Lenient. See Fig. 17 for a graphical representation of these values.

| Method | Critical Value | Uncertainty |
|--------|----------------|-------------|
| Lenient | 5.882% | $\pm$0.37% |
| Moderate | 4.167% | $\pm$0.18% |
| Strict | 1.020% | $\pm$0.02% |

Figure 1

Figure 2

**Cross-section of the Ierapetra Fault Zone (SE Crete)**

W                                                                  E

Eroded
footwall scarp

Holocene surface

**Footwall
block**

Displacement

Throw

Alluvial fan

Holocene surface      **Hanging-wall
block**

Heave

Present-day sea level

Figure 3

Figure 4

a)

| H1/Seafloor | H2 | H3 | H4 | H5 | H6 (Top salt) |

Bf: border fault  Rf: radial fault  Sf: Small crestal fault  Lf: Listric fault
(Vertical exaggeration = 4x)

b)

North ... South

Kavousi fault segment
(1320 m)
Ha fault segment A
(992 m)
(1121 m)
Ha fault segment B
(465 m)
(815 m)
Episkopi fault segment
(186 m)
Pachia Ammos fault
300 m

c)

Southwest ... Northeast
Fault
50 cm

d)

West ... East
Fault in b
Fault
Fault
30 cm

Figure 5

Figure 6

Step 1
Application of a discrete wavelet transform (fault R2 H3)

Figure 7

**Step 1**

Convert fault data to a standard format and required sample size

↓

Perform CWT (Continuous Wavelet Transform) on entire dataset

↓

Find peaks and troughs in each wavelet band, excluding wavelet peaks that are not in current scan range

↓

Select smallest wavelength that is above the threshold of peaks per sample

↓

Take highest peak of this wavelet band as the peak of a new fault segment

**Step 2**

Lowest wavelet trough value is considered the lateral tip of the new fault segment

↓

Lowest wavelet trough on the opposite side of the peak is considered as the other lateral tip of a new fault segment

→

**Step 3**

Gradient descent is applied to fault troughs/minima on a throw-distance plot. A smoothing filter is applied at this stage

↓

Fault Range is considered in the next scan

↓

Gaps between faults are joined up by the algorithm in order to find secondary throw minima

Repeat if new scan range is large enough

**Step 4**

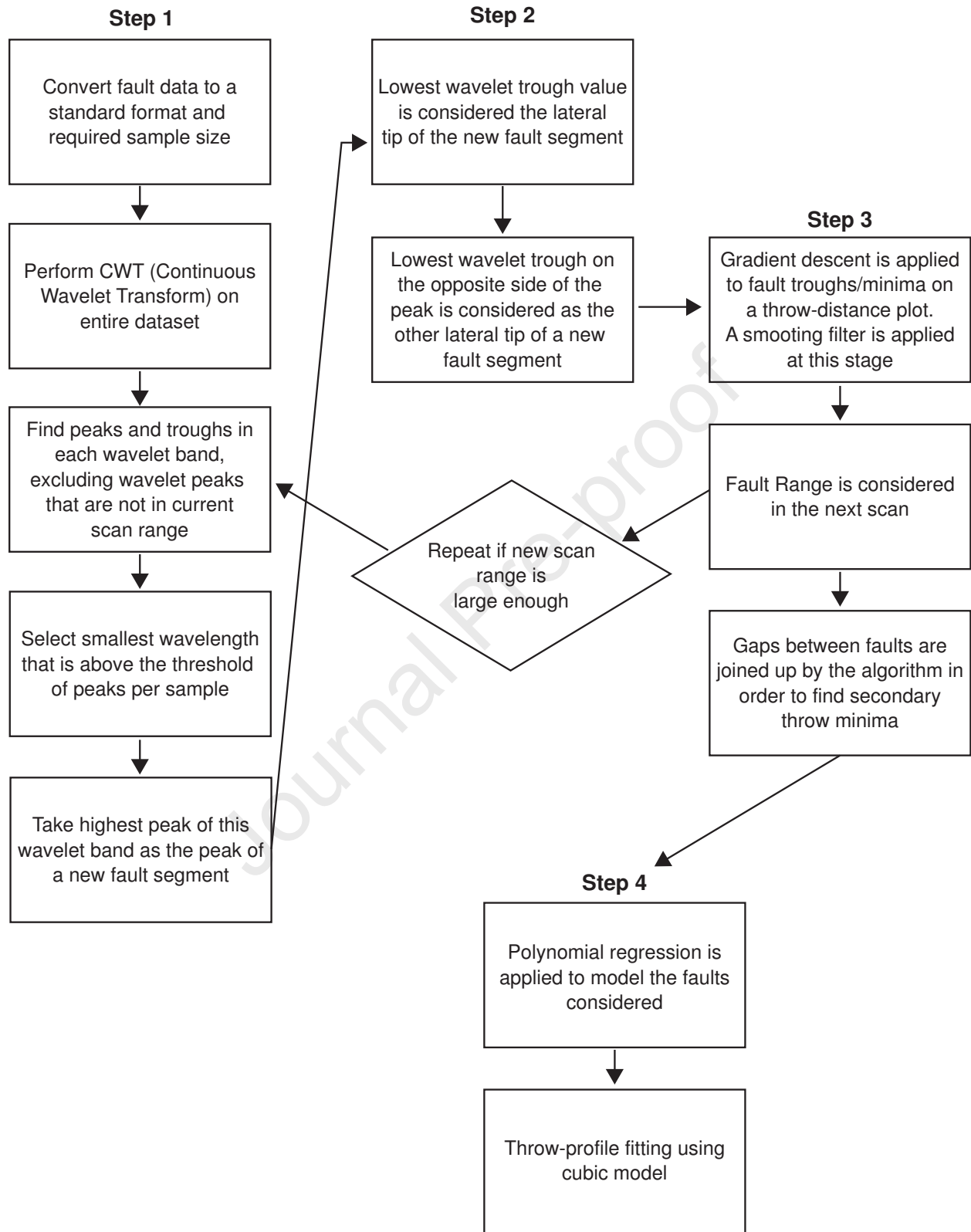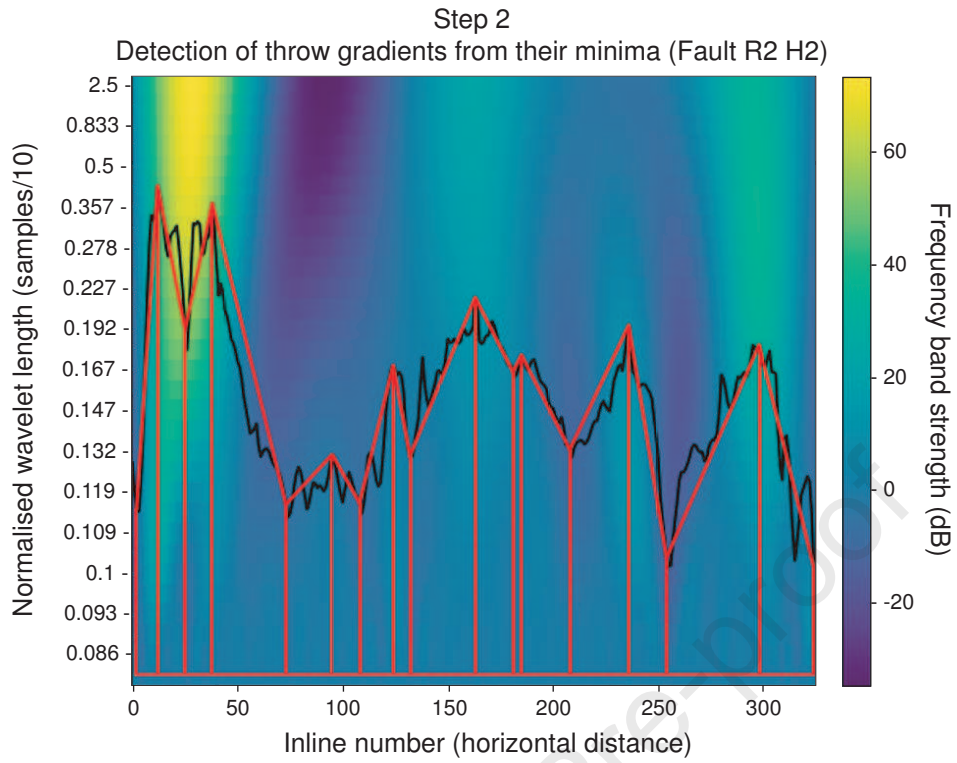Polynomial regression is applied to model the faults considered

↓

Throw-profile fitting using cubic model
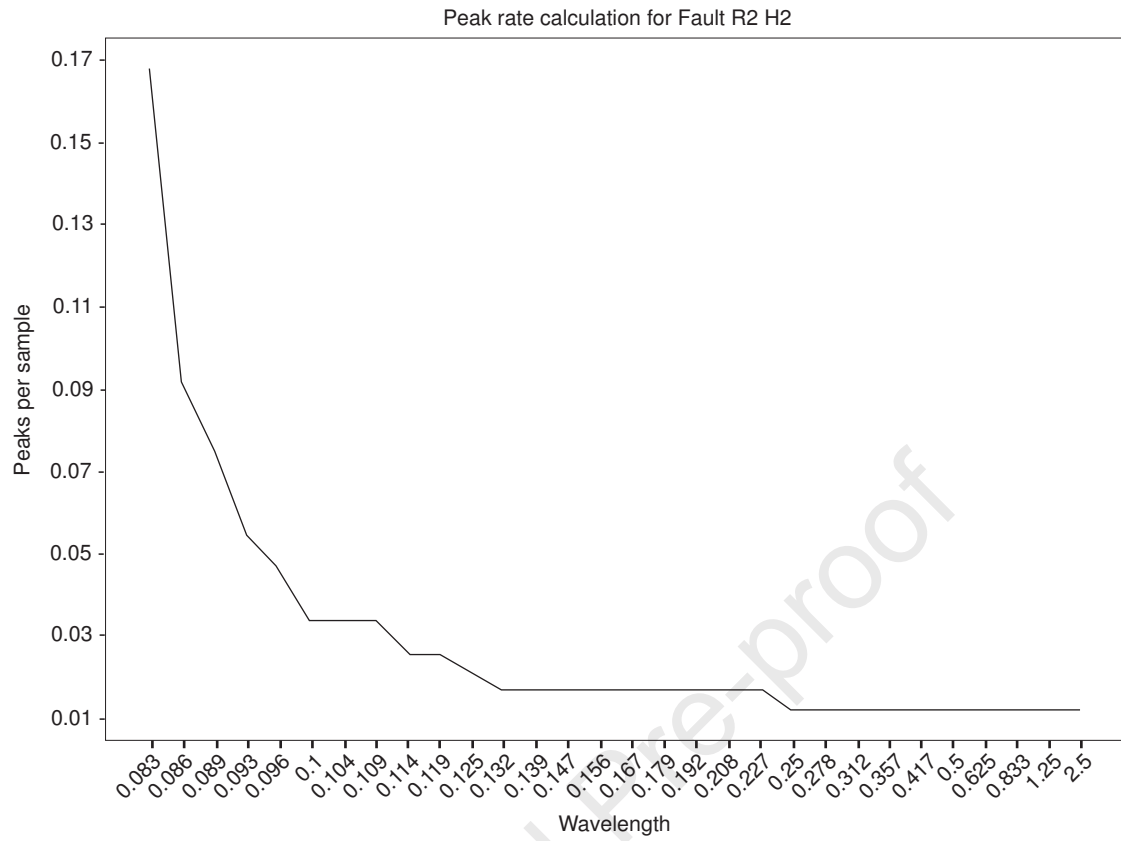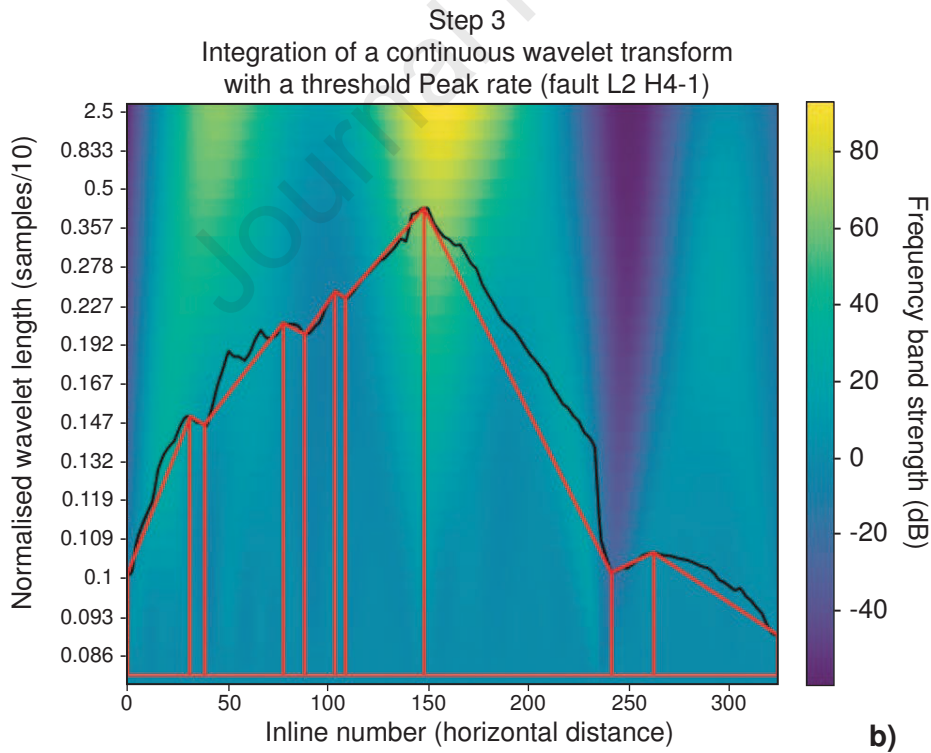
Figure 8

Figure 9

Figure 10

Step 3
Integration of a continuous wavelet transform
with a threshold Peak rate (fault R2 H3)



**a)**

Step 3
Integration of a continuous wavelet transform
with a threshold Peak rate (fault L2 H4-1)



**b)**

Figure 11

Fault throw-profile ftting using a cubic regression model (fault R2 H2)

Local misfit of throw data with the model
due to a local reduction in segment length

Local misfit due to
footwal erosion

**a)**

Fault throw-profile ftting using a cubic regression model (fault R2 H3)

Local misfit due to
footwal erosion

Local misfit due to
footwal erosion

**b)**

Fault throw-profile ftting using a cubic regression model (fault L2 H4-1)

Local misfit of throw data with
the model due to footwall erosion

**c)**

Figure 12

a)

— T/D data before the application of critical (minimum) sampling rate
-- T/D data after the application of critical (minimum) sampling rate



b)

— Raw T/D data
— Original curve fitting using a cubic approach
-- Curve fitting when data are reduced to an Integral Error of 11.6%.

Figure 13

Keystone fault 6 - 11


a)

Linked keystome and polygonal faults (C24)


b)

Figure 14

Distribution of minimum sampling rates (δ) needed to accurately detect fault segments



Figure 15

Figure 16

Figure 17

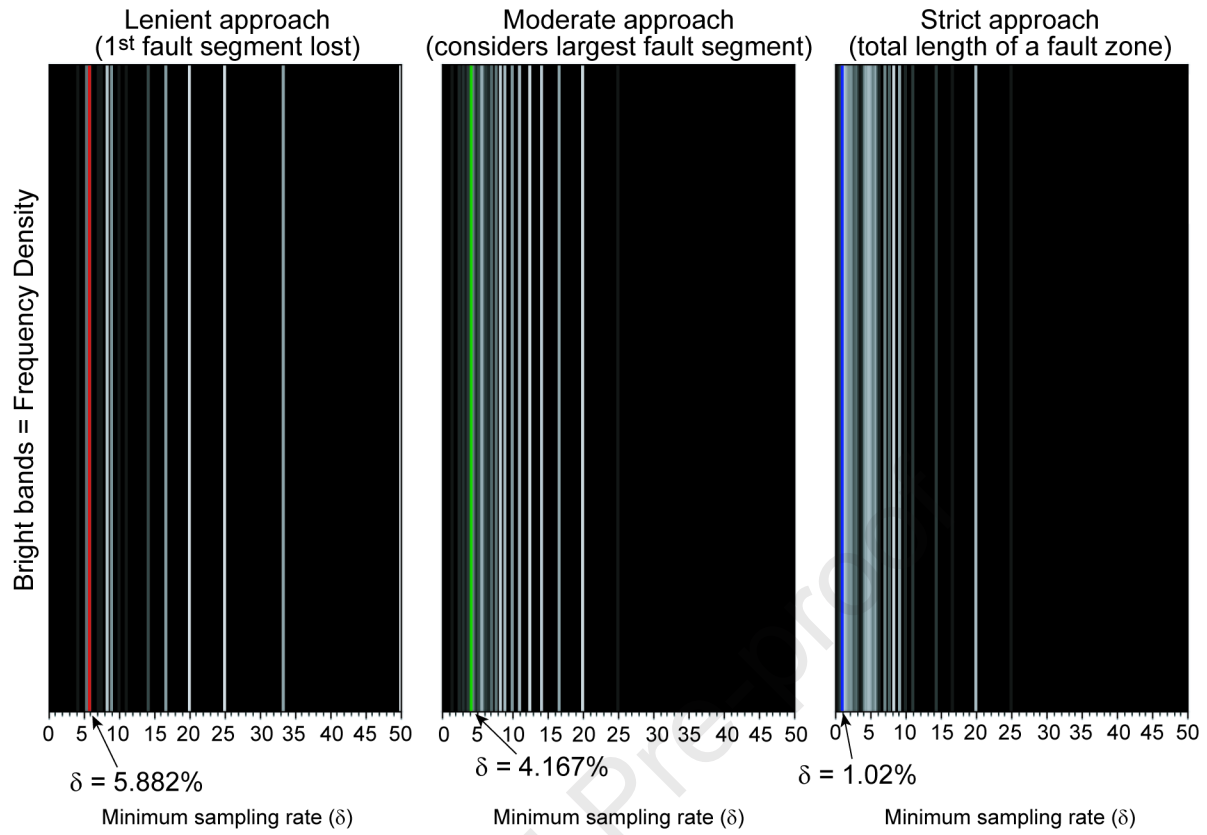**Conflicts of interest statement**

The authors declare that they have no conflicts of interest regarding this work.

**Declaration of interests**

☒ The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

☐ The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: