



This is a repository copy of *Detection of sub-surface damage in wind turbine bearings using acoustic emissions and probabilistic modelling*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/153963/>

Version: Accepted Version

Article:

Fuentes, R., Dwyer-Joyce, R.S. orcid.org/0000-0001-8481-2708, Marshall, M.B. orcid.org/0000-0003-3038-4626 et al. (2 more authors) (2020) Detection of sub-surface damage in wind turbine bearings using acoustic emissions and probabilistic modelling. *Renewable Energy*, 147 (1). pp. 776-797. ISSN 0960-1481

<https://doi.org/10.1016/j.renene.2019.08.019>

Article available under the terms of the CC-BY-NC-ND licence (<https://creativecommons.org/licenses/by-nc-nd/4.0/>).

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 **Detection of sub-surface damage in wind turbine bearings**
2 **using acoustic emissions and probabilistic modelling**

3 R. Fuentes¹, R. Dwyer-Joyce², M.B Marshall², J. Wheals³, E. J. Cross¹

4 *Dynamics Research Group, Department of Mechanical Engineering, The University of Sheffield¹*

5 *Leonardo Centre for Tribology, Department of Mechanical Engineering, The University of Sheffield²*

6 *Ricardo Innovations ltd³*

7 **Abstract**

8 Bearings are the culprit of a large quantity of Wind Turbine (WT) gearbox failures
9 and account for a high percentage of the total of global WT downtime. Damage within
10 rolling element bearings have been shown to initiate beneath the surface which defies
11 detection by conventional vibration monitoring as the geometry of the rolling surface is
12 unaltered. However, once bearing damage reaches the surface, it generates spalling and
13 quickly drives the deterioration of the entire gearbox through the introduction of debris
14 into the oil system. There is a pressing need for performing damage detection before
15 damage reaches the bearing surface. This paper presents a methodology for detecting
16 *sub-surface* damage using Acoustic Emission (AE) measurements. AE measurements
17 are well known for their sensitivity to incipient damage. However, the background
18 noise and operational variations within a bearing necessitate the use of a principled
19 statistical procedure for damage detection. This is addressed here through the use of
20 probabilistic modelling, more specifically Gaussian mixture models. The methodol-
21 ogy is validated using a full-scale rig of a WT bearing. The bearings are seeded with
22 sub-surface and early-stage surface defects in order to provide a comparison of the
23 detectability at each level of a fault progression.

24 *Keywords:* Acoustic Emission, Condition Monitoring, Bearings, Damage Detection,
25 Probabilistic Modelling, Wind Turbines

26 **1. Introduction**

27 Bearing failures are the leading source of downtime in Wind Turbine (WT) gear-
28 boxes and the root cause for this is attributed to Rolling Contact Fatigue (RCF) [1–3].
29 In the majority of loading conditions, fatigue damage begins its life at the *surface*
30 of materials, where high stresses and imperfections due to manufacturing and surface
31 wear coalesce and lead to crack initiation. The case in bearings is unlike typical fa-
32 tigue damage. Hertzian contact mechanics dictates that, under the compressive load at
33 the contact between a rolling element and a bearing, the location of maximum stress
34 will lie a small distance under the surface at the point of contact between a roller and
35 the bearing surface. This has some important consequences regarding the damage pro-
36 gression of a bearing. A growing crack will spend most of its time under the surface,
37 where it has minimal impact on the operation of the rest of the system. However, once
38 a crack emerges on the surface, the progression of failure is accelerated through contact
39 with the rolling elements and this will generate spalling. At the point of initiation of
40 spalling, the progression of damage is quick as debris is introduced into the rest of the
41 mechanical system, thus accelerating the overall failure of the gearbox.

42 Currently, WTs are designed with an overall target lifetime of 20 years [4], a design
43 requirement which extends to all of their subsystems. However, the average service life
44 of wind turbine gearboxes often falls much below the 20 year target [5]. This is a prob-
45 lem; even though gearboxes are not the most unreliable subsystem, they do cause the
46 most downtime [5]. Minimising gearbox failures is thus a key element in increasing
47 overall wind turbine productivity [2]. Because bearing surface damage releases debris
48 into closed-loop oil systems, sampling the oil quality and checking for debris within
49 the oil system is, to date, still used as a reliable technique for diagnosing the overall
50 condition of WT gearboxes [6, 7]. It is also at this point that vibration-based monitor-
51 ing systems are able to detect the presence of defects. The fact that bearing damage
52 has reached the surface and introduced debris into the oil system motivates the need
53 for detecting fatigue cracks in bearings before they reach this stage, so that preventive
54 maintenance can be carried out and impact to the rest of the gearbox can be minimised.
55 Detecting *subsurface* damage at the incipient stage has been identified as a critical as-

56 pect of wind turbine condition monitoring [3]. Rolling contact fatigue is exacerbated
57 in planetary gearboxes, where the bearing raceway is loaded exclusively in the torque
58 direction. This exerts a compressive load on the same point along the circumference
59 of the raceway, as illustrated in Figure 1. In order to avoid the problems associated
60 with the introduction of debris and accelerated failure, it is highly desirable to be able
61 to detect the *incipient* failure of the bearing, at the point where a fatigue crack has just
62 initiated. There are three critical aspects that will determine the outcome of a damage
63 detection system [8]: 1) the physical sensing system, 2) the damage-sensitive features
64 extracted from the data and 3) the damage identification strategy applied to those fea-
65 tures. This paper addresses these problems. As for the physical sensing system, Acous-
66 tic Emissions (AE) are proposed as a measurement strategy. The damage-sensitive fea-
67 tures extracted from AE data play a fundamental role in the ability to identify damage.
68 In this paper, the state of the art of AE features are reviewed and compared and new
69 features are proposed using advanced signal processing tools. Lastly, a rigorous dam-
70 age identification strategy is proposed that addresses the key challenge of discerning
71 operational and environmental effects from the damage-sensitive features. This is car-
72 ried from a probabilistic modelling point of view, using Gaussian mixture models in
73 combination with dimensionality reduction tools.

74 *1.1. Subsurface cracks*

75 The interest in subsurface cracks has grown since the realisation that fatigue cracks
76 in gearbox bearings tend to start around non-metallic inclusions [10], introduced during
77 the manufacturing process. The presence of these inclusions, coupled with high stress
78 concentrations under the surface, leads to the development of fatigue cracks, often re-
79 ferred to as White Etching Cracks (WEC), White Structure Flaking (WSF) [11, 12], or
80 simply “butterfly” cracks due to their butterfly shape (with the “wings” following a path
81 from the inclusion, out towards the surface). These cracks tend to grow in the region
82 around 1mm under the contact surface of typical WT bearings [13] and it has been
83 proposed that their formation is driven both by chemical and mechanical processes.
84 Chemically, it is the diffusion and release of hydrogen into bearing steel [14], through
85 lubrication and water ingress that drives the formation of WECs. Mechanically, over-

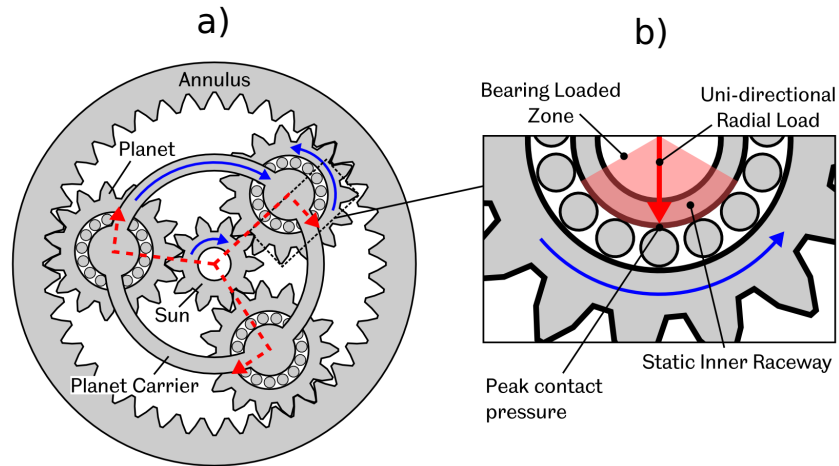


Figure 1: a) Diagram showing gearing setup for a planetary gearbox. Note the planet bearings are constantly loaded in the torque direction, indicated by the red arrows. b) zoom-in to one of the planet bearings, highlighting the loaded zone [9]

86 load events arising from wind gusts, braking and torque reversals drive stress con-
 87 centrations around non-metallic inclusions to yield point and lead to the formation and
 88 growth of WECs. Since the realisation that inclusions in bearing steel directly lead to
 89 subsurface cracks, the quality control of the manufacturing processes has dramatically
 90 improved. However, inclusions will always be present even in today's high standard
 91 of steels. In fact, it has been shown that it is typically the smallest inclusions that lead
 92 to the greatest stress concentrations and therefore the development of WECs [12]. An
 93 example of a WEC at the initial stage of propagation is shown in Figure 2, observed on
 94 a WT bearing section [12].

95 1.2. Damage detection with Acoustic Emissions

96 When considering the dynamic response of a system, it is a generally well-accepted
 97 principle that the physical size of damage is inversely proportional to the frequency at
 98 which its effects will be manifested in its dynamic response [15–17]. Furthermore,
 99 there is a well established relationship between the AE response of a metal and fatigue
 100 crack growth [18]. With this in mind, subsurface damage on bearings represents the
 101 smaller end of the scale, requiring relatively high frequency measurements, when com-

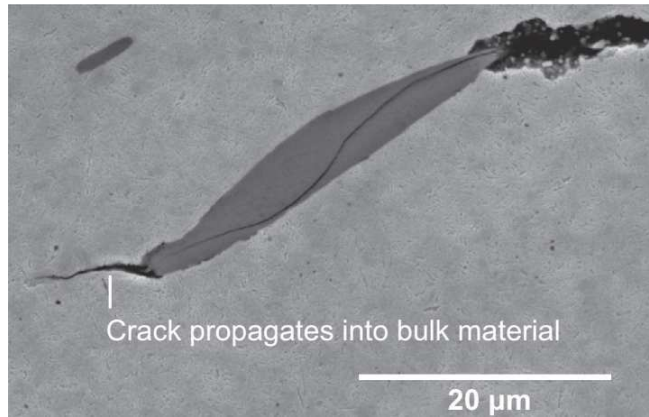


Figure 2: Example of a WEC propagating around a non-metallic inclusion [12]

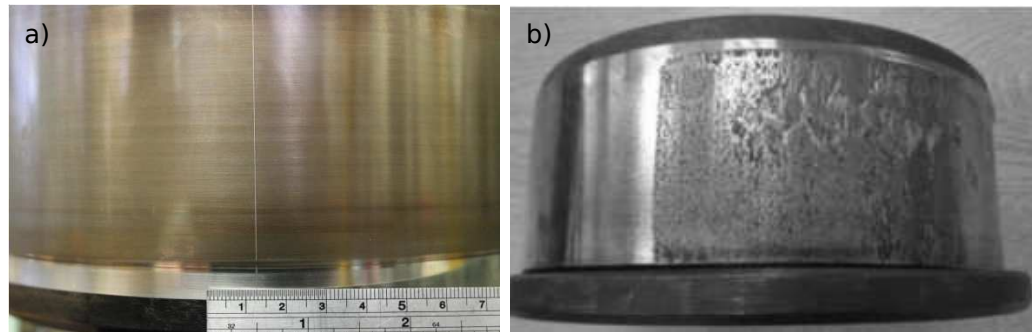


Figure 3: Illustration of damage in representative WT planetary gearbox bearings showing a) a line etch similar, used in [9] which is of similar form to the surface damage introduced in this study (see Figure 6 for the actual profiles) and b) damage arising from real operational conditions, focused on the point at which spalling occurs [12].

102 pared to traditional vibration-based monitoring, in order stand a reasonable chance of
103 being detected. Following this reasoning, this paper investigates the detectability of
104 subsurface cracks using Acoustic Emission (AE) measurements.

105 Acoustic Emissions (AE), when used within a Structural Health Monitoring (SHM)
106 or Non-Destructive Testing (NDT) context, are high frequency stress waves that prop-
107 agate through a material. These waves can be generated by a number of different
108 mechanisms including stress, plastic deformation, friction and corrosion. AE occurs as
109 a result of the release of elastic energy by any one of these mechanisms, which is then
110 propagated through the material as elastic waves. The concept of AE in engineering
111 structures is analogous to the release and propagation of energy that takes place during
112 earthquakes [19], as a result of fractures within fault planes. It is generally accepted
113 that the movement of slip-planes characteristic of micro-cracks during the application
114 of stress and yield [20, 21] leads to the generation of AE. Features extracted from AE
115 measurements have been shown to be successful indicators of the early onset of cracks
116 in various applications [18, 22, 23]. AE testing is a passive method, in the sense that
117 one is listening to the acoustic response of the material when mechanical stress is ap-
118 plied to it. Most materials will have a certain level of AE activity even in an undamaged
119 state when stress is applied to them, the technical term for this is the Kaiser effect [24].
120 However, when defects such as cracks or spalling are present in the material, the AE
121 response when stress is applied will tend to be more frequent, of higher amplitude, and
122 may have different spectral characteristics depending on the material properties of the
123 medium where the waves propagate. Friction processes also tend to generate AE, as the
124 impact between micro-asperities that occurs during contact of two materials releases
125 elastic energy into the system [25]. There is no clear definition of the frequency range
126 that constitutes an AE measurement. This will depend on the physics of the particular
127 defect; a typical AE stress wave generated from the initiation of a crack in steel can
128 range from 50kHz to 2MHz [21, 26].

129 AE being now a popular measurement technique, its application to bearing mon-
130 itoring hasn't gone without attention, but the problem of detecting small subsurface
131 cracks is currently far from solved. One of the barriers to investigating this problem
132 is that subsurface damage is hard to find, validate and measure in an operational bear-

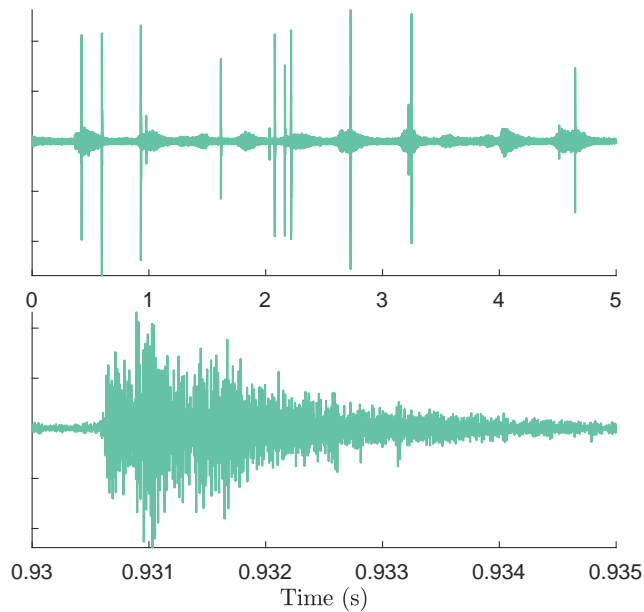


Figure 4: Example of a) an AE signal measured from an undamaged bearing in operation and b) a zoom-in to one of the “hits” characteristic of AE measurements.

133 ing and much more so within a wind turbine, so all investigations resort to laboratory
 134 experiments under controlled environments. The major problem with this lies in the
 135 defect. Most investigations (including some by these authors) use artificially seeded
 136 defects at the surface to show that a given methodology is able to detect incipient dam-
 137 age [26–36]. An example of such defect is illustrated in Figure 3a. The problem with
 138 surface defects (from an investigative point of view) is that they represent the later
 139 stages of damage, quickly lead to spalling and are also relatively easy to detect as con-
 140 tact between rollers and surface damage releases large amounts of acoustic energy into
 141 the system. An example of a late damage state is shown in Figure 3b, at which point
 142 detection through AE or even vibration becomes trivial with the current state of re-
 143 search. A number of studies have attempted to tackle this issue and offer investigations
 144 on bearings run from undamaged conditions all the way to failure, usually under an
 145 environment that accelerates failure [37–39]. While these investigations significantly
 146 advanced the general field of AE-based monitoring, they have not been carried out
 147 with statistical rigour. More specifically, the detectability of incipient damage is only

148 viewed qualitatively and no attempt is made at quantifying it.

149 Detecting damage using data-driven methods can be cast as a *novelty detection*
150 problem. First a statistical model of the system in its undamaged state is characterised
151 and an alarm threshold is established. Further observations are judged based on a
152 *novelty index*; a distance metric that quantifies how far the new observation lies from
153 the baseline undamaged state. The decision of whether an observation belongs to a
154 damaged or undamaged state is then given by whether the novelty index lies above or
155 below the threshold. Under this approach the detectability is judged solely in terms of
156 change from a baseline condition and can be quantified in terms of false positives and
157 false negatives. It also provides ample opportunities for objectively judging different
158 statistical and signal processing methods in terms of their false positive and negative
159 rate performance. A central part of this paper involves performing this comparison
160 using rigorous methods of novelty detection.

161 A critical aspect of any damage detection system is its ability to separate the ef-
162 fects of damage from those of normal gearbox operation. This is particularly true of
163 AE measurements. The recorded AE signals will not only contain the high frequency
164 stress waves produced by crack growth and plastic strain, they will also contain a sig-
165 nificant amount of AE energy that is unrelated with any damage mechanism. A large
166 part of this “benign” energy being released as AE in a gearbox will come from inter-
167 nal stresses and friction. A key challenge is then to separate the AE activity related to
168 normal operation from those that are not. In WT gearboxes, operational variability will
169 naturally arise from varying wind speeds. In variable-speed WTs, changes in speed will
170 result in a wide variety of AE activity, with or without the presence of defects. On a
171 basic level, varying rotational speed will introduce variability into any features derived
172 from a frequency domain or periodicity based analysis of the AE signals, such as those
173 presented in [40, 41]. Varying speed will also introduce changes to the background
174 noise resulting from friction throughout the gearbox and will also lead to changes in
175 oil temperature. Oil viscosity is highly dependent on temperature and this completely
176 changes the lubrication regime of the gearbox, which has a direct effect on the resulting
177 background AE noise signature [42]. AE monitoring of fixed-speed WTs will also suf-
178 fer from this type of variability. Varying wind speed introduces changes in the internal

179 loads of the gearbox, and the release of energy as AE in materials is directly related to
180 applied loads. Furthermore, load variation also causes variation in AE energy through
181 changing temperatures and the effect that this has on the lubrication regime. Higher
182 temperatures generally lead to lower oil viscosities, which in turn drives an increase
183 in friction through higher asperity contact between the surfaces of roller elements and
184 bearing raceways.

185 In general, studies investigating bearing monitoring, using either vibration or AE
186 measurements, have not taken into account this operational and environmental vari-
187 ability. Recently steps have been taken in [26] to collect AE measurements on a failing
188 bearing at varying loads and rotational speeds, consequently, the notion that load and
189 speed affect AE both directly, and indirectly, through variations in temperature was
190 confirmed. However, again no attempt was made at quantifying the detectability in
191 terms of comparing a statistical model of the undamaged bearing against the rest of the
192 observations.

193 A further complicating factor when performing bearing monitoring in practice is
194 detecting incipient faults from practical measurement locations. Since the early days
195 of bearing monitoring it was recognised that detecting even a large seeded fault was
196 much easier if AE is measured in the inner raceway of the bearing compared to the
197 outer raceway [27]. However, in practice it is clearly less invasive to place a sensor
198 on the outside of the gearbox. Placing a sensor inside a gearbox would often require
199 machining and this can lead to overall reduced gearbox reliability.

200 In summary, the current challenges facing bearing monitoring are clear:

- 201 • To detect incipient failure at the subsurface stage using a rigorous statistical
202 methodology.
- 203 • To perform such detection under the effects of operational and environmental
204 variability.
- 205 • To use practical and non-invasive measurement locations to carry out the moni-
206 toring

207 This paper presents a general methodology for detecting damage in rotational com-

208 ponents and validates this using a rig that is representative of the operational environ-
209 ment of a WT epicyclic gearbox. More specifically, this environment involves realistic
210 varying loads and speeds as well as the changing temperature and oil properties that
211 these create. A focus is placed on detecting subsurface as well as small surface cracks.
212 The approach to inducing subsurface damage used here is different from previous stud-
213 ies; instead of taking an undamaged bearing and running it until failure, subsurface
214 damage was carefully seeded in a bearing in order to be able to carry out a robust and
215 conclusive comparison of damaged and undamaged states. Furthermore, three more
216 surface faults were seeded in a different bearing in increasing sizes ranging from $5\mu\text{m}$
217 to $50\mu\text{m}$ width. The reasoning for using seeded defects here is that it is the most robust
218 and conclusive way of validating a damage detection methodology because one knows
219 the exact state of the bearing at every stage. However, a strong focus was placed on us-
220 ing defects, both subsurface and surface-level, that are representative of the very early
221 stages of fatigue damage in bearings. A description of the experimental rig and the
222 defect-seeding procedures is given in Section 2. Particular attention is devoted to the
223 damage-sensitive feature extraction process of AE signals. In broad terms, there are
224 two major and different types of features extracted from AE signals 1) those based on
225 the characteristics of discrete bursts of energy, often termed AE *hits* and 2) those based
226 on analysis of the periodicity of the global AE signal. While in the literature, studies
227 tend to focus on either one approach or the other, here it is of interest to investigate
228 the performance of each type of damage-sensitive feature. A description of the signal
229 processing methods used to derive the different damage-sensitive features is given in
230 Section 3. The third element of the damage detection methodology is the statistical pat-
231 tern recognition. The methodology used in this paper is to treat the problem as one of
232 novelty detection, where the probability distribution of the damage-sensitive features
233 is modelled for data belonging to an undamaged state. This allows for the computation
234 of a *novelty index* on subsequent observations, and when this exceeds a given alarm
235 threshold this can be indicative of damage. The particular procedure, based on Gaus-
236 sian mixture modelling is described in detail in Section 5.3. The validation of the entire
237 methodology, including the various damage sensitive features and the novelty detection
238 is presented in Section 5, where results are presented for the data-set collected on the

239 experimental rig with the four different stages of damage.

240 **2. Experimental set-up**

241 The experimental rig used in this study was devised to investigate planetary gear-
242 box bearings, due to their propensity to fail prematurely owing to the fact that the load
243 transferred from the rotating outer raceway to the static inner raceway peaks at a fixed
244 position along the circumference of the inner raceway. This loading condition is il-
245 lustrated in Figure 1 and leads to short fatigue lifetimes around the loaded section of
246 the bearing, leaving an “un-used” fatigue life outside of this region. This rig has been
247 used in previous investigations into bearing monitoring using various techniques rang-
248 ing from vibration [43] to AE [35, 36] and ultrasound monitoring [9]. However, this
249 is the first study considering subsurface defects, as well as surface defects with widths
250 under $100\mu\text{m}$.

251 The objective of the rig is to generate a compressive radial load on the inner race-
252 way inside a planetary sun bearing sub-assembly. In order to achieve this, the rig com-
253 prises of two bearings: an inner “test bearing”, which is housed inside an outer “main
254 bearing”. The inner test bearing then houses a stationary shaft, which is connected via
255 two steel lugs to a hydraulic ram, capable of delivering a total compressive load up to
256 1600kN. In order to apply rotation, the inner raceway of the outer bearing is coupled to
257 the outer raceway of the inner bearing. A tensioned pulley is then used to drive these
258 two raceways together, with power delivered from an electric motor. A cross-sectional
259 diagram of the assembly is shown in Figure 7a, while Figure 7b shows a photograph
260 of the entire rig, highlighting the main components. Figure 8 shows a more schematic
261 view of the main components and applied loads. The inner bearing is coloured in red
262 and the outer bearing is coloured in blue. The rolling elements are shown in light grey.
263 The main interest in this investigation is the inner raceway of the inner bearing. To-
264 gether with the shaft, this inner raceway remains stationary, with the rolling elements
265 revolving around it, a compressive load being applied at the bottom (via the hydraulic
266 ram). Due to the constant compressive load, it is here where damage normally initiates
267 in planetary sun bearings and so all seeded defects in this study are located so that their

268 position lies exactly at the bottom of the circumference of this inner raceway.

269 The inner bearing used in this rig is an NU2244, which is typically used in WT
270 gearboxes. Its inner raceway has a bore diameter of 220mm, and the outer raceway has
271 an outer diameter of 400mm. Its maximum dynamic load rating is 1600kN, while its
272 fatigue load limit is of 250kN. Further specifications of this bearing can be found in
273 [44].

274 2.1. Seeded Defects

275 This section details the defects seeded into the inner raceway. For the purposes
276 of this study, two types of defects were seeded in order to emulate increasing damage
277 levels: subsurface and surface defects. Overall, a total of six bearing conditions were
278 examined, summarised in Table 2.1. Note that *two* undamaged bearings were used in
279 the experiment, in order to generate robust training and validation datasets for the data-
280 driven damage detection models. Further discussion on the importance of having two
281 undamaged bearings, for validation purposes is given in Section 5

Label	Condition	Severity	Bearing
UD1	Undamaged		A
UD2	Undamaged		B
D1	Subsurface Damage	800kN	C
D2	Subsurface Damage	1000kN	C
D3	Surface Damage	5 μ m	D
D4	Surface Damage	20 μ m	D
D5	Surface Damage	50 μ m	D

Table 1: Summary of bearing conditions

282 2.1.1. Surface defects

283 Surface defects represent a damage condition in a relatively advanced stage. The
284 objective here was to generate defects as small as possible in order to emulate the early
285 stages of surface damage. In previous work at the University of Sheffield [36], a spark
286 erosion technique was used to etch the surface of a raceway to emulate a surface crack,

287 which generated surface defects of approximately $200\mu m$ width (similar to the damage
288 shown in Figure 3a). In order to achieve a smaller defect, more representative of the
289 early stages of a surface crack, a Cubic Boron Nitrite (CBN) grit was used to scratch
290 the surface using a six-axis Computer Numerical Control (CNC) machine. This was
291 performed at three different pressures, with each one at a different angular position
292 on the raceway. The aim of using three different angular locations along the raceway
293 was to be able to perform a test with three different defect sizes by simply positioning
294 the different defects on the loaded zone of the bearing. The angular positions of the
295 defects were such that only one defect was loaded at any given instant in time, these
296 are shown in Figure 5. The target sizes for the seeded defects were $5\mu m$, $20\mu m$ and
297 $50\mu m$, although the actual profiles obtained are shown in Figure 6. These profiles
298 were taken by first filling the scratches with silicon in order to extract a negative of
299 the profile, and then measuring this with a three-dimensional optical profiler. Note that
300 damage conditions D4 and D5 seem very similar to each other. In fact, most of the
301 profile shown in Figure 6a and b is the pattern of material removed from the inside of
302 etch. The profile of the inner parts of the etches were in fact too small to measure with
303 the available profilers. Note that the curvature of background of each profile shown
304 in Figure 6 is due to the bending of the silicon samples and not the curvature of the
305 bearing.

306 2.1.2. *Subsurface defects*

307 A key element of this paper is the study of the detectability of defects in a bear-
308 ing before they propagate to the surface. To achieve this, a subsurface defect was
309 seeded to a second raceway by means of compressing its outer surface with a rolling
310 element. The compression was applied using a hydraulic press capable of applying up
311 to 2000kN. Subsurface yield was estimated to occur at 1000kN for this bearing, using
312 Hertzian contact mechanics relationships. To ensure that subsurface yield occurred,
313 while also preventing the damage propagating to the surface, the yield process was
314 monitored using AE. Some of the observations on AE from this damage seeding are
315 further discussed in [21].

316 During the seeding of subsurface damage, a large increase in AE energy was ob-

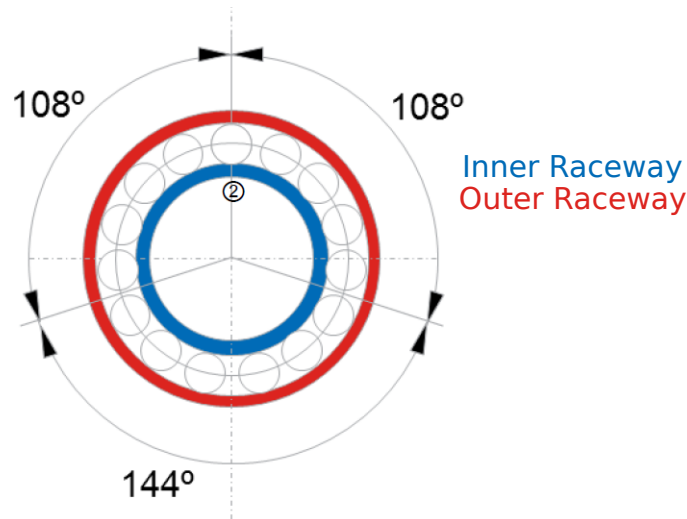


Figure 5: Angular positions of defects along raceway. Note that this angular separation ensures that when one roller passes over one of the defects, none of the rest of the rollers will pass over the other two defects at the same time.

317 served in the 800kN to 1200kN range. Visible surface damage was only found when a
 318 bearing was loaded beyond 1700kN. Although several tests were carried out on numer-
 319 ous raceways, on the final raceway, faults were seeded using two compression levels,
 320 at 800kN and 1000kN. These were applied on the same circumferential indices as for
 321 the surface damage, so that only one damage site is located within the loaded zone of
 322 the raceway.

323 To summarise, two damaged bearing raceways were used for testing. One raceway
 324 contained three surface etches with increasing sizes, to emulate increasing levels of
 325 damage. The second raceway contained two seeded subsurface cracks, with increasing
 326 levels of maximum compressive load.

327 2.2. Test conditions

328 The objective of this study is to perform damage detection of realistic defects in a
 329 realistic operational WT gearbox environment. In order to achieve this, a test schedule
 330 was designed to capture, for each of the raceway conditions (outlined in Table 2.1),
 331 the effects of varying load, speed and temperature. Preliminary tests were conducted,

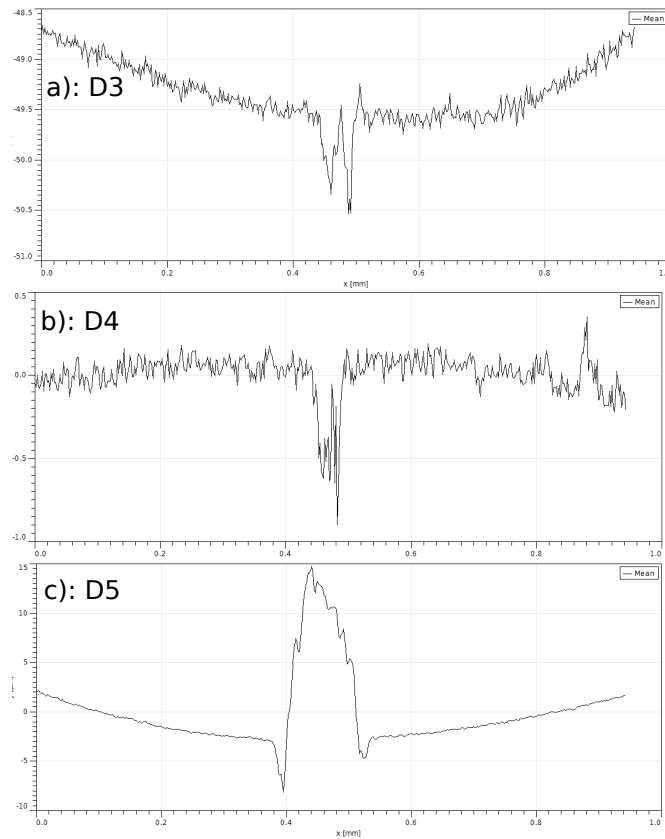


Figure 6: Average (negative) profiles of the three surface damage sites, taken using a 3D optical interferometer profiler, according to Table 2.1. Note that the profiles of the smallest defects only capture to average width, together with excess material on the sides.

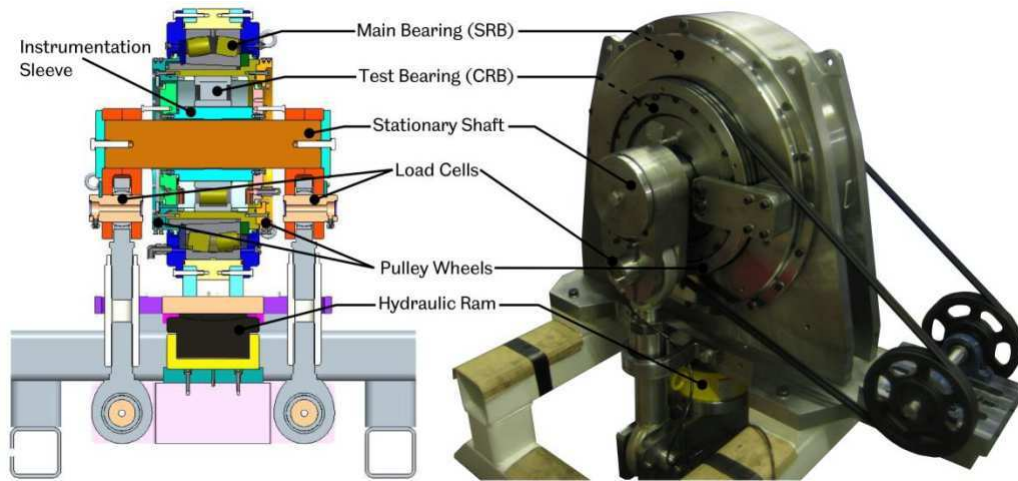


Figure 7: a) Diagram showing fixture, rig bearing and test bearing, outlining the location of AE measurement channels. b) Photograph of rig in the lab.

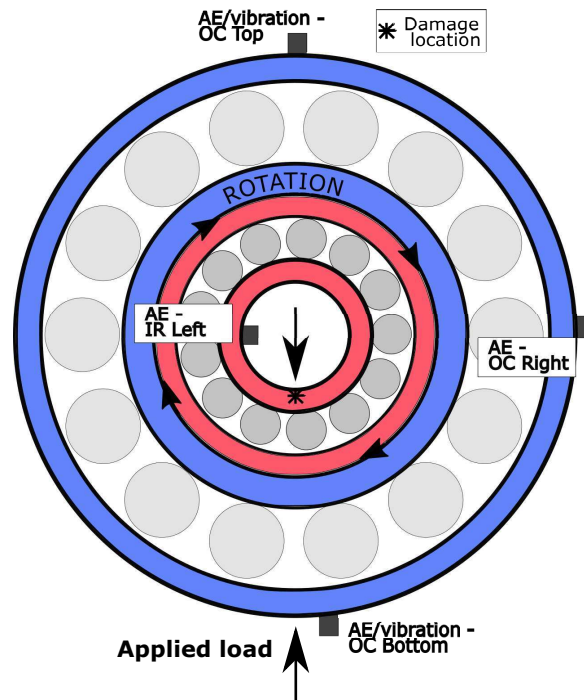


Figure 8: a) Diagram showing fixture, rig bearing and test bearing, outlining the location of AE measurement channels. b) Photograph of rig in the lab.

332 stepping the compressive load at 100kN steps from 0 to 1200kN. This pointed to three
333 major regimes of AE activity, around the low load (0-400kN), medium load (400kN-
334 800kN) and high loads (800kN-1200kN). For this reason three loads were selected for
335 a test schedule: 200kN, 600kN and 1000kN. The temperature of the rig proved difficult
336 to control precisely. The factors that affect the rig and oil temperatures are the operating
337 load, bearing condition (a failed bearing introduces debris into the system and drives
338 the temperature up through friction), ambient temperature, accrued usage time, whether
339 a heat exchanger is present in the oil system and whether this is aided by a cooling
340 device (such as a fan). In this rig the only means of controlling the temperature directly
341 are via the heat exchanger, the operating load and sequencing of the applied loads. In
342 general it is easier to warm up the rig, than to cool it down, as once it runs and a load
343 is applied, it will quickly warm up and reach a stable temperature. Therefore, it was
344 decided to split the tests into low and high temperatures. This split is reasonable, given
345 that the main effect that temperature introduces (to the AE activity) is an increase in
346 friction at higher temperatures from a reduction of viscosity [45, 46]. In order to keep
347 the temperature down, the low temperature runs were performed early in the morning,
348 testing low loads first and keeping cooling fan on.

349 Table 2.2 shows the complete schedule of tests carried out. This is also a realistic
350 scenario given that the temperature in a WT gearbox will vary in a similar fashion.
351 When not operational, or at low wind conditions, modern gearboxes will keep circu-
352 lating the oil through a heat exchanger, to keep it at a stable temperature (and thus the
353 optimal viscosity).

354 *2.3. Instrumentation*

355 Four AE sensors and one accelerometer were used to measure the overall dynamic
356 response of the bearing throughout the tests. The positions are illustrated in Figure 8.
357 Three AE sensors were placed outside the bearing, mounted on the Outer Casing (OC)
358 of the outer raceway. One AE sensor was placed inside the Inner Raceway (IR), by
359 machining a small part of the stationary shaft in order to fit the sensor. Furthermore,
360 the AE sensor located in the IR was sprung-loaded into position. The IR AE sensor
361 was located along the circumference of the IR, at angular orientation of 60deg from

Speed (RPM)	Load (kN)	Temperature
100	200	Low (Fan ON)
100	600	
100	1000	
100	1000	High (Fan OFF)
100	200	
100	600	

Table 2: Bearing test schedule. Each row was performed sequentially, and this schedule was used for every bearing condition.

362 the damaged zones (the bottom). In these tests, the IR sensor is the closest to the
363 damage, which is the main AE source of interest. Because of the compressive load
364 applied to the inner raceway, the bottom is in direct contact with the rollers, while the
365 top develops a slight clearance. This defines the acoustic path of the AE stress waves
366 to propagate from the source (the damage), down through the rollers, outer bearing
367 and casing and around the circumference of the rig. Therefore, of the three sensors on
368 the Outer Casing, the OC bottom location is closest to the AE source, followed by the
369 OC-right and OC-top (as illustrated in Figure 8). Note that due to symmetry, it was
370 deemed reasonable to not include a measurement position on the OC left side.

371 The three OC-AE sensors used in this study were Mistras 3MICRO-30D, fitted with
372 a differential cable for noise reduction. The IR sensor was a Mistras NANO-30, which
373 is a non-differential sensor. The Micro30D has a marked resonance at approximately
374 350kHz, while the Nano30 has a flatter response in the range of 200kHz-500kHz. This
375 is relevant as the sensor frequency response shapes the acquired signals significantly.
376 Compared with vibration sensors, the frequencies of interest are much broader, in the
377 range of 50kHz - 2MHz. It is therefore much harder to achieve a flat frequency response
378 across all frequencies of interest, and so one must accept the significant mechanical
379 filtering that the sensor applies to the “true” underlying AE stress wave. It must also be
380 noted that sensor-to-sensor variability is usually much more significant in AE sensors,
381 compared with vibration instrumentation. One tri-axial accelerometer was mounted

382 next to the AE sensor at the OC-top location.

383 All data acquisition was recorded with a National Instruments (NI) C-DAQ system.

384 This comprised of several modules, with data recorded at the following sample rates:

- 385 • NI-9223, four-channel analogue module sampling at 1MHz with 16-bit Ana-
386 logue to Digital (ADC) conversion.
- 387 • NI-9234, three-channel IEPE module sampling at 51.2kHz.
- 388 • NI-9213, thermocouple, for temperature measurements (one sample per test).

389 Several operational parameters were also acquired in order to assess the influence
390 of each one on the AE response. These parameters were:

- 391 • Test Bearing Speed (RPM).
- 392 • Left-side and Right-side Load (kN).
- 393 • Oil Temperature.
- 394 • Casing Temperature.

395 Data was acquired in “trials” of ten second duration. For each bearing condition
396 in Table 2.1 and each operational condition in Table 2.2, ten different different trials
397 were gathered. The resulting data-set thus comprises approximately 1260 time series
398 records, each taking approximately 322MB of memory, totalling 491GB of data.

399 **3. Damage-sensitive features from AE measurements**

400 The goal of this section is to discuss the signal processing required to derive dam-
401 age sensitive features from AE data. In broad terms, there are two approaches to this
402 problem and both will be covered here. The fundamental idea behind AE monitoring
403 is to detect the energy release characteristic of the interaction between stress and the
404 plastic deformation around a crack in the form of high frequency stress waves. These
405 short bursts of high frequency data, often referred to as “hits” have been illustrated
406 in Figure 4. The first approach is to identify these hits and to characterise them. On

407 a non-rotating system with low levels of background noise, the simple presence of a
408 hit could be indicative of damage. In this setting, simply characterising the rate of
409 generation of AE hits, or “hit count”, would be a sufficient damage-sensitive quan-
410 tity. However, in rotating systems a certain amount of background AE activity will be
411 expected, as already discussed in Section 1.2. In this setting, it is more desirable to
412 work with damage-sensitive features of individual AE hits. These could be as simple
413 quantities such as the energy, duration and amplitude of hits, but could also be based
414 on more complex models, such as Fourier or autoregressive coefficients. Detecting an
415 AE hit and computing features from these short bursts of energy represents one of the
416 two major different strategies for signal processing. This strategy will be discussed in
417 more detail in Section 3.1. One key advantage of dynamic response data originating
418 from rotating machinery is its tendency to be periodic, and this can be taken advantage
419 of for efficiently deriving damage sensitive features. In terms of AE, one would expect
420 a burst of AE energy to occur every time one roller passes over a damaged area. This
421 information can be encoded efficiently with Fourier coefficients of rectified signal en-
422 velopes. This type of periodicity-based analysis represents the second major approach
423 to deriving damage-sensitive features, and is discussed in more detail in Section 3.4.

424 *3.1. Hit-based features*

425 One of the key points of AE data, from a signal processing point of view, is that
426 the bursts captured by the AE acquisition system are very short in comparison to the
427 large amount of time that needs to be spent monitoring. Because of the high sample
428 rates required to capture these high frequency waves, this means that a lot of noise is
429 recorded, in comparison to the amount of useful AE bursts. To put this in context,
430 the bursts recorded from a yielding steel specimen may last in the order of 2000 μs .
431 If one were to monitor at 1MHz for 1 second and expect 15 bursts (which is roughly
432 how many bursts are expected in the rig in this paper at 100 RPM), this would mean
433 that approximately 3% of the data points are informative and the rest is noise. Given
434 the high sample rate, data storage and handling becomes an issue if one wishes to
435 monitor for long periods of time. This has led the AE monitoring community to develop
436 hit-identification strategies, where an AE hit is defined as a burst large enough for

437 it to be likely to be caused by material fracture. Non-rotating systems are naturally
438 “quiet” in their undamaged states and it is normally straightforward to identify hits by
439 setting a threshold on the overall AE signal; the value of the threshold would mostly
440 be determined by the background noise level of the environment and the electrical
441 noise. Rotating systems on the other hand are noisy; there are more sources of acoustic
442 noise, and bursts that are not related to damage but generated by friction, roller impact,
443 arising from minor misalignments and transient loads. These lead to significant and
444 often varied levels of background noise. Having a constantly changing noise level
445 introduced by periodic friction complicates the basic process of identifying an AE hit.
446 Figure 4a illustrates an AE signal with varying background noise, where some AE hits
447 are evident well above the background noise. The problem is that the lower energy hits
448 that may lie close to, or even be buried under, the noise. In order to identify these AE
449 hits, a special adaptive threshold methodology had to be devised.

450 AE data streams comprise millions of points and the feature extraction process
451 being described here is applied to hundreds of data files. Efficient computation of
452 features is therefore required. In order to achieve efficient compression of AE sig-
453 nals, while preserving the information contained within them, a multi-level Discrete
454 Wavelet Transform (DWT) was applied to all AE signals in this study. Each level of a
455 DWT first splits the signal, using a half-band quadrature mirror filter into its low and
456 high frequency components and decimates the signal by half. Each level of decom-
457 position comprises wavelet coefficients, each representing half of the frequency band
458 of the level above with half the number of points. Multilevel DWT is a popular data
459 compression tool in the general context of signal and image processing [47]. Its ap-
460 plication to AE data makes sense given that the information in the signals is contained
461 in a short bandwidth, dictated by the resonance of the sensor. In this study, the sen-
462 sors used all had resonances in the range of 100-300kHz. For this reason, a two-level
463 DWT was applied that split the signal into two frequency bands, of 0-250kHz and 250-
464 500kHz. Only the lower frequency wavelet coefficients were used effectively reducing
465 the number of data points by half, while keeping all the information of the sensor reso-
466 nant frequencies. If one were to perform monitoring using a broad-band measurement
467 technique, such as a Laser Doppler Vibrometry (LDV) or fibre brag-grating, this step

468 should be applied carefully so as not to throw away important information. However,
469 piezoelectric transducers will always be resonant around a narrow band and so per-
470 forming a DWT that retains only that band is bound to be an efficient pre-processing
471 step.

472 3.1.1. AE hit identification

473 The objective of hit identification is to establish the existence of an AE hit and its
474 location in time. As discussed above, simple threshold strategies, which work well in
475 non-rotating machinery, fail to correctly identify all of the AE hits in a gearbox setting.
476 The problem is that when non-stationary background noise is present in the system,
477 the appropriate threshold that separates a high energy event from background noise
478 will change with time. If a threshold were to be applied directly to AE data in this
479 setting, it will either capture all high energy hits and leave out the lower energy ones,
480 or be set low enough to capture low energy events but also be triggered constantly by
481 background noise.

In order to correctly identify AE hits, a thresholding strategy is required that iden-
tifies the presence of a hit, within a constantly changing noise floor. The methodology
developed here makes use of a hit *identification function*, which computes the differ-
ence between the local signal energy E_t and a lagged version of itself at E_{t-a} . The
difference is then normalised against the local noise level at $t - a$. The resulting hit
identification function $H(t)$ captures rapid changes in energy against the local back-
ground noise. The local energy can be defined as a moving Root Mean Square (RMS)
statistic within a given short time period. The identification function is formally defined
as:

$$H(t) = \frac{E(t) - E(t - a)}{E(t - a)} \quad (1)$$

482 where a represents the lag of the local energy difference. Its value is critical to the
483 success of the identification function, it should represent the expected time over which
484 an AE event will reach its maximum energy. In this study, the lag was tuned empirically
485 and a value of $a = 500\mu s$ was used. The presence of a hit is established when the
486 identification function exceeds a prescribed threshold, T_H . A value of $T_H = 2$ was
487 used in all hit identification procedures presented in this study. This can be interpreted

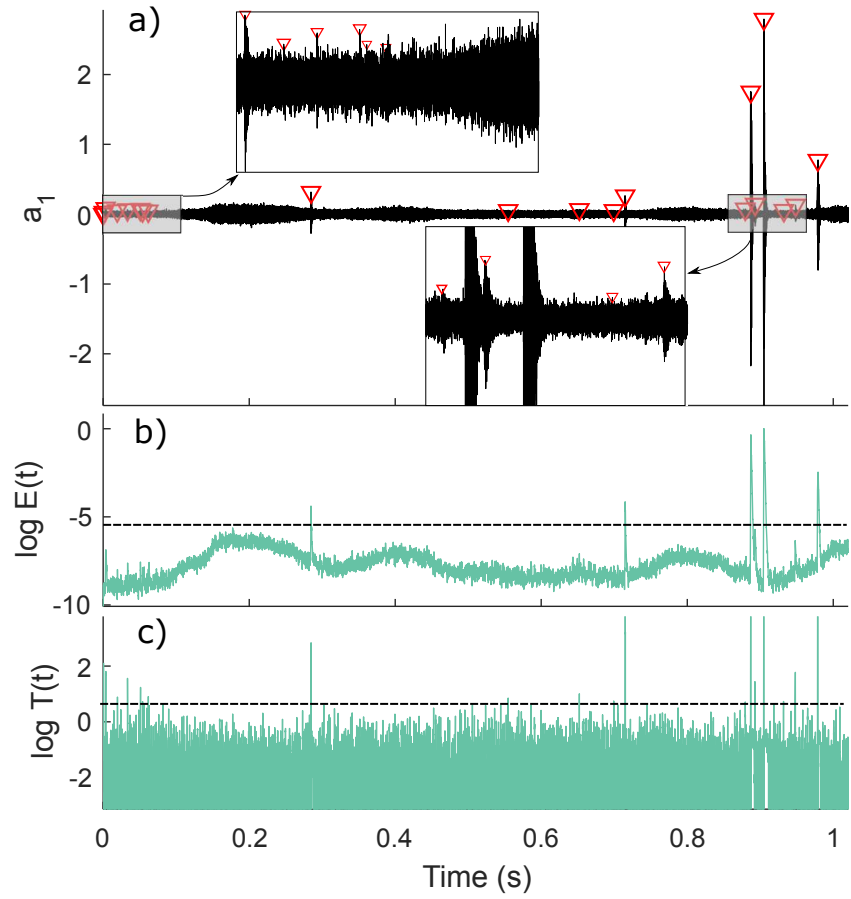


Figure 9: Illustration of effects of setting a threshold either too low or too high, showing a sample of raw AE data with periodically varying noise level

488 as a hit being defined when the value of the local energy quickly rises to two times its
489 baseline level.

490 In summary, the steps taken for detecting the presence of a hit, for every AE channel
491 are:

- 492 • Decimate signal with a wavelet decomposition
- 493 • Take an envelope $E(t)$, of the wavelet coefficients, to capture the amplitude mod-
494 ulation of the process. The envelope could consist of a Hilbert transform or a
495 moving RMS to compute the local signal energy.
- 496 • Compute the hit identification function given by equation (1).
- 497 • Set a threshold over $H(t)$ and record all instances where this threshold is ex-
498 ceeded.

499 These steps are illustrated in Figure 9, using a one-second AE recording of an un-
500 damaged bearing, taken from the OC-top location (see Figure 8). Figure 9a shows the
501 wavelet coefficients of the 0-250kHz band for this data. Figure 9b shows the local en-
502 ergy function, $E(t)$, while Figure 9c shows the hit identification function derived using
503 Equation (1), including the threshold of 2, the exceedance of which defined an AE hit.
504 The AE hits identified from exceedances of $H(t)$ are shown with triangular markers.
505 Note that using this adaptive thresholding methodology, it is possible to identify AE
506 hits across the entire scale of energies. Figure 9a illustrates this by zooming-in into two
507 regions where low-energy AE hits are present which would have clearly not been iden-
508 tifiable with a simple threshold over the raw data, the wavelet coefficients or the local
509 energy. This procedure is important as it enables the characterisation of individual hits
510 even across the entire range of energy levels.

511 Once the presence of a hit has been identified by the adaptive thresholding strategy,
512 further steps are required to identify the precise start and end times of each AE hit.
513 For a given hit, a rough start time is already given by the time of exceedance of the
514 threshold over the hit identification function. The end time is defined as either a) the
515 point at which the local energy decays back to within 10% of the local baseline level

516 (before the hit started) or b) the start time of another AE hit occurring before the energy
517 decays back to the local baseline.

518 Up to this point in the processing, only a rough start time for the hit has been
519 identified, based on the local energy function. However, this will only be accurate to
520 within the given time window used to evaluate $E(t)$. Furthermore, stress waves in
521 a material can propagate in various different modes. The fastest mode will tend to
522 that of longitudinal waves, but this will also carry the least amount of energy. Shear,
523 surface and possibly also Lamb waves (depending on the thickness of the material) may
524 arrive after the first arrival of longitudinal waves, all carrying much more energy. This
525 time delay carries information regarding the total distance a stress wave has travelled,
526 so it is important to capture the precise time of the arrival of the longitudinal mode.
527 Employing a threshold for onset identification, the longitudinal wave will invariably
528 be missed and the arrival of the shear or surface modes is more likely to be captured
529 instead. To overcome this, the methodology proposed by Kurz [48], based on Akaike's
530 Information Criteria (AIC), is used here in order to identify the precise moment of the
531 onset of AE waves. This method computes a cumulative variance of a hit, forwards and
532 backwards and creates an AIC function as the superposition of these two. The point
533 at which this function reaches a minimum indicates the highest change of information
534 (or variance) in the signal and thus the onset of the AE wave can be established by
535 looking for a minimum of this function. An illustration of the AIC function indicating
536 the minimum, where the onset is defined is shown in Figure 10.

537 3.2. AE hit summary features

538 Once a table of start and stop positions has been extracted from the AE data for
539 every channel, it is relatively straightforward to go back to the signal and save only
540 the waveforms at those time instances. This is the strategy that has been adopted; it
541 significantly reduces the amount of data stored, and focuses all the post-processing on
542 the data points corresponding to AE hits only, which as discussed before, comprise
543 only a small percentage of the data points in the signal.

544 There are numerous features that can be extracted once the waveform has been
545 captured. Because an AE waveform is a transient event, there are some key features

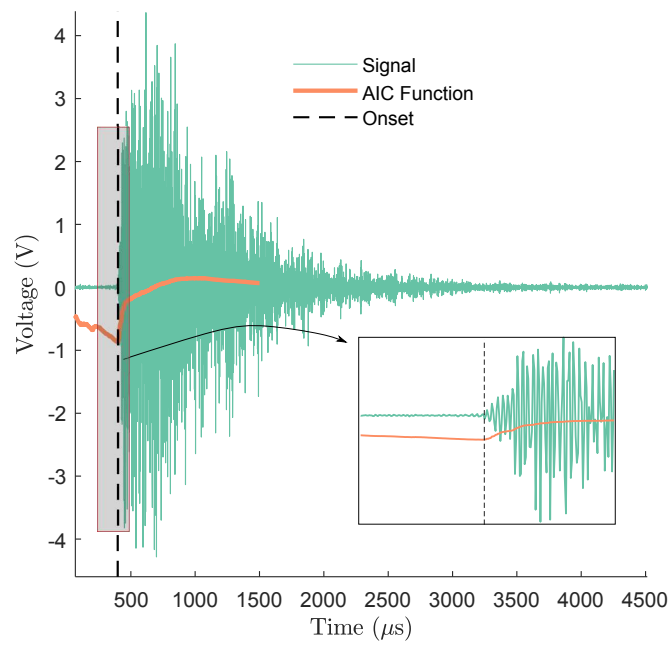


Figure 10: Illustration of AIC onset function on a sample AE hit. The onset is defined where the function reaches its minimum, indicating the greatest change in variance within the signal window.

546 that can characterise it in general, but simple, terms.

547 Possibly the most informative feature is the energy contained in the waveform.
548 Different sources of AE will release stress waves at widely different energy levels.
549 The energy is easily computed as the sum of squares of the data points. The power
550 normalises the energy by the duration of the signal. In the case of a transient waveform,
551 such as that of an AE hit, energy, power and RMS will all be related to each other
552 since the duration is a function of the total energy, because of the exponential decay in
553 amplitude. This means that there is some correlation between these variables.

554 The rise time is defined as the time difference between the waveform onset and
555 its maximum amplitude. The information this carries is valuable because, due to the
556 difference in speeds of different wave modes, some will arrive first and some later, thus
557 giving a rough indication of how far the source is from the sensor. In practice, in a
558 steel structure, waves will propagate as longitudinal, transversal, surface, and possibly
559 Lamb waves. These waves will all travel at different speeds and will carry different
560 proportions of the total energy of the wave-front. The Lamb wave modes may or may
561 not be excited, as their existence requires that the wavelength of the AE be of the
562 same order of magnitude to the thickness of the material it is travelling through. An
563 investigation of Lamb waves is outside the scope of this paper, but their use should not
564 be discarded and is marked as future work. In steel, longitudinal, shear and surface
565 waves arrive in that order. The amount of energy they carry is also given in that order.
566 Therefore the first arrival will always be from a longitudinal wave, and the maximum
567 amplitude will tend to be recorded at the arrival of a surface wave. The usefulness of
568 this is that the rise-time of an AE hit is a useful feature as it gives an indication of
569 how far the wave has travelled. Waves that come from far away will have high rise-
570 time (separation between longitudinal and surface waves) while the opposite is true for
571 short rise times.

572 Other features that are collected are the peak amplitude of the signal, the total du-
573 ration and the decay time. The duration is defined, during the hit extraction process,
574 as a decay after the peak amplitude to a level within a specified tolerance of the base-
575 line noise, immediately before the hit. The duration will tend to be a function of the
576 energy in the waveform, but also of the physical mechanism exciting the wave, and it

577 is therefore a useful feature. Once all of these features for each hit are computed, they
578 are assembled so that inference with a Bayesian network can be performed with them.

579 In summary, six summary statistics are taken from every AE hit:

- 580 1. Maximum Amplitude
- 581 2. Power
- 582 3. Energy
- 583 4. Rise-time
- 584 5. Decay time
- 585 6. Duration

586 These features are commonly used in the detection of damage from AE measure-
587 ments [22, 36, 41]

588 3.3. *AE hit Auto Regressive coefficients*

589 Whilst the hit summary statistics may provide sufficient information for detection,
590 their general drawback is that they provide a simplistic representation of the signal,
591 they also require a significant amount of pre-processing (such as the identification of
592 an accurate onset), which can be prone to error. An alternative is to represent the signal
593 in terms of a model that automatically captures the main characteristics of the signal.
594 In this paper, Auto Regressive (AR) models are used as a damage-sensitive feature that
595 provides a more detailed representation of the individual AE hits. In this paper, AR
596 model weights, w , are fit via linear regression to every single AE hit extracted using
597 the procedure outlined above. As with the summary statistics, this AR model is fit
598 to the single-level Discrete-Wavelet-Transformed signal, thus halving the number of
599 points to compute and focusing on the frequency range of interest (250kHz to 500kHz
600 in this case).

601 3.4. *Modulated signal envelope features*

602 If one signal processing paradigm were to be singled out as having enabled large-
603 scale fault identification in rotation machinery, taking frequency decompositions of
604 rectified signal envelopes would easily win. The idea is simple; every time a periodic

605 load is applied to a gearbox component (roller-bearing, gear-tooth contact for example)
606 which contains a sizeable defect, this will generate a high frequency burst of energy,
607 which can often be sensed at other points in the gearbox. These high-frequency burst
608 will not necessarily be evident behind a frequency spectrum of the dynamic response,
609 because the actual frequency content of these burst will be dictated by various reso-
610 nant frequencies at which these bursts are transmitted. These resonant frequencies are
611 characteristic of the gearbox assembly, not of the bursts and so will be generally ex-
612 cited during normal gearbox operation. What *is* characteristic of these bursts is that
613 they happen at periodic intervals and this period can single out the particular compo-
614 nent that is generating them. The result of this is that the amplitude modulation of
615 the dynamic response signal contains more information about damage processes than
616 does the signal itself. One simple and well-established way of extracting this amplitude
617 modulation is via the use the Hilbert-Huang transform. A simple frequency analysis,
618 via a Discrete Fourier Transform (DFT), of the signal envelope has been shown to high-
619 light well defects in many different types of rotating machinery [49]. Whilst this idea
620 was originally applied to vibration signals which measure dynamic response in a much
621 lower frequency range, this technique has been applied to AE measurements with a
622 good degree of success [28, 33, 40].

623 In this paper, the DFT coefficients of AE signal envelopes are used as a damage-
624 sensitive feature. This provides a useful point of comparison, given that there so far,
625 this type of feature have been widely used in the majority of papers investigating dam-
626 age detection using vibration and AE in rotating systems [28, 33, 40, 41], and more
627 specifically, detection of sub-surface damage.

628 As discussed in the previous sub-sections, the AE signals are originally sampled
629 at relatively high sample rates (1MHz in this case), in order to capture the high fre-
630 quencies at which the stress waves characteristic of AE travel (250kHz-500kHz in this
631 case). By contrast, the frequencies at which one would expect to find evidence of dam-
632 age in the amplitude modulation of these signals is much lower, belonging in the range
633 between zero and the low hundreds of Hertz. In the specific case of this bearing rig, the
634 ball pass frequency between the rolling elements and the bearing is estimated at 15Hz,
635 hence, there is a large disparity between the bandwidth of the original envelopes and

636 the frequency content of interest. A DFT at this original sample rate would yield very
637 poor frequency resolution and would also be computationally expensive. In order to
638 remedy this, in this paper, the Hilbert transform is applied to the one-level Discrete-
639 Wavelet-Transformed signal. This halves the number of points used for computation
640 of the envelope and also only takes the envelope over the frequency bandwidth of inter-
641 est, thus eliminating the potential for noise to be introduced here. After this envelope
642 is derived from the DWT, these wavelet coefficients are further low-pass filtered and
643 down-sampled down to an effective sample rate of 100kHz, one-tenth of the origi-
644 nal sample rate. At this point, a Short Time Fourier Transform (STFT) is applied to
645 down-sampled envelopes, with a window length of 250000 points (2.5 seconds). This
646 is enough to capture a potential of 45 cycles of damage-related AE bursts in every win-
647 dow, with frequency resolution of 0.2 Hz. This frequency resolution is appropriate to
648 capture the damage process at the expected ball pass frequency of 15Hz. Note that
649 for further processing (the probabilistic modelling detailed in Section 4), this damage-
650 sensitive feature was truncated to 2502 spectral lines, which yields an effective analysis
651 bandwidth of 12.5kHz.

652 **4. Probabilistic Modelling**

653 The process of detecting damage from the observed AE damage-sensitive features
654 is a problem of searching for outliers in statistical data. An outlier can be defined as
655 an observation that is different enough from the rest of the observations that it is likely
656 to have been generated by a different mechanism [50]. There are two fundamental
657 elements of outlier analysis in data. The first is a statistical model of the reference
658 (undamaged) condition data. This model is often assembled as a *probability density*.
659 The second element is a statistical distance that measures how far any given observation
660 is to the centre of the data mass, relative to the reference probability density. For the
661 purposes of this paper, an “observation” shall be defined as a multivariate vector of
662 damage-sensitive features, evaluated at one instance in time. In the context of AE
663 data (as discussed in Section 3) this could comprise, for example, a vector of AE hit
664 statistics, autoregressive model coefficients, or rectified signal envelope spectra.

665 In a large number of application domains [50], including SHM and condition mon-
666 itoring [16, 51], it is common to assume that the underlying probability density of the
667 reference data can be safely modelled as a Gaussian distribution. Under the Gaussian
668 assumption, the *Mahalanobis Squared Distance* (MSD) is a good measure of the rela-
669 tive closeness between observations and the reference set. The approach of modelling
670 damage-sensitive features with single-Gaussian distributions and using an MSD as a
671 novelty index is now in wide-spread use in the field of statistical outlier analysis [50]
672 as well as in the field of Structural Health Monitoring (SHM) [16, 51, 52]. However, a
673 major drawback of the single-Gaussian distribution approach is that it is unsuitable for
674 modelling the probability density of data that has been generated by multiple regimes.
675 In monitoring contexts, multiple regimes often arise from changing environments and
676 operation. In the case of bearings, varying loads, speeds and temperatures will generate
677 differing characteristic responses in the AE features. These will manifest themselves as
678 multiple modes in the probability density of AE features. One way of modelling these
679 complex probability densities is through the use of mixture distributions, discussed in
680 the following section.

681 4.1. Dimensionality reduction

682 One characteristic of some of the damage-sensitive features being used in this study
683 is that they are high-dimensional. The AR coefficients comprise vectors of 150 dimen-
684 sions while AE envelope spectra contain 2500 dimensions. Generally speaking, most
685 novelty detection schemes rely on computing distance metrics between feature spaces
686 of new observations against feature spaces representing normal conditions. In this set-
687 ting, it is a well recognised result and an effect of the “curse of dimensionality”, that
688 the contributions of individual dimensions to a distance metric tend to get masked in
689 high dimensional feature spaces [53]. This presents a problem - if damage will only
690 introduce a change to a handful of dimensions within a high dimensional feature, this
691 may not show up when computing a novelty index. Furthermore, high dimensional fea-
692 tures also present a problem when computing covariance matrices of statistical models.
693 It is a generally well-accepted rule that one needs at least twice as many observations
694 as there are dimensions in a feature space in order to begin to accurately capture the

695 covariance structure of the data.

In order to remedy these problems, the authors turn to the use of standard dimensionality reduction techniques. Here Principal Component Analysis (PCA) is used as a dimensionality reduction strategy. Its use is wide-spread within the field of statistical learning [54, 55] as a data visualisation and pre-processing tool. PCA can be viewed as a class of linear Gaussian models [56] and is represented by a linear transformation from a d -dimensional data/feature space \mathbf{y} into a lower-dimensional space \mathbf{x} , called the principal component scores,

$$\mathbf{y} = \mathbf{C}\mathbf{x} + \varepsilon \quad (2)$$

696 where \mathbf{C} is a $d \times d$ PCA rotation matrix and ε is an additive Gaussian noise term. It is an
697 orthogonal transform designed so that \mathbf{x} contains a rotated version of \mathbf{y} aligned in the
698 directions of greatest variance. The resulting PCA scores, \mathbf{x} , therefore contain most
699 of the information (in terms of variance of the original data-set) within the first few
700 dimensions. It is therefore common practice to use only a handful of the dimensions of
701 \mathbf{x} , this could be determined by observing how much variance is lost as one throws away
702 dimensions. In this paper, however, the dimensionality of the PCA scores is fixed to
703 five, so as to enable a comparison between different damage-sensitive features without
704 introducing the effects of differing dimensionality of the feature space into the outlier
705 analysis process.

706 The PCA rotation matrix, \mathbf{C} must be estimated. By definition, \mathbf{C} is the eigen-
707 decomposition of the covariance of the data/feature set \mathbf{y} , so it can be estimated through
708 an eigenvalue analysis. However, this may not scale well to very high dimensions. An
709 alternative approach is an iterative Expectation Maximisation (EM) algorithm, which
710 leads to the notion of probabilistic PCA [57, 58]. This paper uses the EM approach
711 described in [57] to learn matrix \mathbf{C}

712 In the scenario being investigated, one does not have access to damaged condition
713 data sets. The PCA rotation matrix, \mathbf{C} is learned using the same data used for the
714 novelty detector (described below, and so it is only representative of the undamaged
715 class. Note that detection would be *easier* if \mathbf{C} could be learned using both damaged
716 and undamaged data sets, as the main differences in both of these sets would likely

717 represent the greatest variance, leading to two different columns in \mathbf{C} describing each
718 condition. This would lead to a significant separation of the two different conditions in
719 the PCA scores, \mathbf{x} . However, what is of interest here is the performance of a novelty
720 detector on a projection of the undamaged class only. The discussions that follow will
721 discuss the problem of novelty detection. This novelty detection is carried out in the
722 lower-dimensional domain of PCA scores, which is denoted throughout the paper as \mathbf{x} .

723 4.2. Gaussian Mixture Models

724 A natural way of dealing with probability densities that arise from multi-regime
725 processes is to partition the space of features into the different regimes and to fit a
726 density model to each region of the feature space. Ideally, one would have a label of
727 the regime type associated with each observation. However, in practice this is difficult
728 to attain and there may be more naturally occurring clusters than one has labels for.
729 It is therefore more desirable to work with algorithms that automatically partition the
730 space into different regions. This task is generally referred to as *clustering* and there
731 is a wide choice of algorithms available for carrying this out [55]. Here, the focus
732 is on novelty detection, so whatever clustering scheme is used should also define a
733 probability density over the feature space. The Gaussian Mixture Model (GMM) is
734 used here because 1) it is a flexible density estimator for multi-modal data, 2) there are
735 efficient algorithms for clustering, or data partitioning with GMMs and 3) it is possible
736 and straight-forward to derive a novelty index in order to perform damage detection
737 with this model. This subsection discusses these three important points.

738 The concept of using a GMM for novelty detection within an SHM and condi-
739 tion monitoring context has been investigated in some recent studies[59–61], including
740 some by these authors[62]. The approach to using a GMM to achieve novelty detection
741 taken here follows that of [62], where a novelty index is derived using the probability
742 density of the GMM. The parameters of a GMM model with K components comprise
743 of the set of K means, covariances and mixing proportions. For notational simplic-
744 ity, these can be encoded in the vector $\boldsymbol{\theta} = \{(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k), (\mathbf{S}_1, \dots, \mathbf{S}_k), (\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_k)\}$.
745 Learning the appropriate parameters of a GMM involves choosing a parameter set, $\boldsymbol{\theta}$
746 that maximises an objective function. Damage detection then involves evaluation of

747 the likelihood of new observations, given the reference parameter set θ .

748 4.2.1. The GMM likelihood function

The likelihood of the model plays a central role in the damage detection approach of this paper: it is the objective function used for parameter estimation, it defines the probability density and is also therefore a useful novelty index. The Gaussian mixture is defined as a weighted sum of Gaussians, so its probability density can be written as,

$$p(\mathbf{y}|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}|\mu_k, \mathbf{S}_k) \quad (3)$$

where, $\mathcal{N}(\mathbf{y}|\mu_k, \mathbf{S}_k)$ represents the normal Gaussian density of the k^{th} component, with mean vector μ_k and covariance matrix \mathbf{S}_k . The term π_k defines the relative contribution of component k to the total density, also known as responsibilities. Equation (3) defines the *likelihood* of the model, given observed data \mathbf{y} . This quantity plays an important role in determining the model parameters as well as deriving appropriate novelty indices. In practice, it is easier and more computationally stable to work with the *log-likelihood* of the model. This is true both for parameter learning, as well as for evaluation of novelty indices for damage detection. Written explicitly as a sum over all observations, the log-likelihood for the GMM is,

$$\ln \mathcal{L}(\theta) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \ln\{\pi_k p(\mathbf{y}_i|\mathbf{z}_k)\} \quad (4)$$

749 where z_{nk} denotes the posterior responsibility of component k for generating the i^{th}
750 observation. This is a convenient form for the log-likelihood function as it is given
751 as a sum of logarithms This results from a formulation in terms of hidden variables
752 (see [56, 58]). The model parameters are estimated using maximum-likelihood learn-
753 ing, which involves maximising the log-likelihood of Equation (4). Unfortunately,
754 evaluating this quantity involves knowledge of the partitioning of the data-set into the
755 different clusters. This cluster-assignment is not known a-priori which makes this an
756 optimisation task with missing data, or hidden variables. The formulation of the GMM
757 log-likelihood function in terms of these hidden variables has already been discussed
758 in Section 4.2.1. The Expectation Maximisation (EM) algorithm was derived to deal

759 with exactly this type of problem; it is a general framework for performing maximum-
760 likelihood parameter optimisation of models with hidden data [63]. Furthermore, it has
761 been shown that EM is a suitable learning algorithm for a wide class of linear Gaussian
762 models with latent variables [56], to which GMMs belong to. EM is used in this paper
763 as a learning strategy for all GMM models of damage-sensitive features.

764 4.2.2. *GMM model selection through cross-validation*

765 A common issue is that of choosing an appropriate model order; the number of
766 components, K , of the GMM. Too many components will result in the GMM over-
767 fitting the density estimate (assign very high density in regions of data where observa-
768 tions exists, and no density elsewhere). On the other hand, using too few components
769 would fail to correctly capture the complexities of the true underlying distribution. Ide-
770 ally, the number of components should be determined directly from the training data.
771 In this paper, the authors turn to the use of cross-validation in order to perform this task.
772 Cross validation procedures divide the available training data into different training and
773 testing subsets, so that the model error is evaluated on previously un-seen data and so
774 the generalisation performance of the algorithm can be evaluated. For this paper, k-fold
775 cross-validation is used, as this is a well-known robust and efficient way of performing
776 cross-validation in a wide variety of statistical models [54]. K-fold cross-validation di-
777 vides the available training data into K different partition, or folds. These are selected
778 at random. Then, the GMM model is trained using EM, on $K - 1$ folds, thus leaving
779 one fold out of the training set. The model error is then computed for the left-out fold,
780 and the process until every fold is held-out at least once as a test set. The benefit of this
781 is that the model error can be quantified entirely in terms of its performance on previ-
782 ously unseen data, and so this gives a good insight into the generalisation performance
783 of the model.

784 As for the model error, two different quantities will be used in order to assess
785 the quality of the model fit. The first is the Bayesian Information Criterion (BIC)
786 [64], a quantity that is derived from the model likelihood, but penalises models of
787 higher complexity (more components). The BIC considers the intrinsic function of the
788 GMM as a density estimator. However, it does not consider its function as a classifier,

789 more specifically in the novelty detection case, a one-class classifier. Therefore, the
790 second error metric considered in the cross-validation procedure in this paper is the
791 misclassification rate. However, in order to treat the GMM model as a classifier, one
792 must also consider that the overall model comprises both the density estimation *and* the
793 novelty threshold together. How this is determined plays a central role in the success
794 or otherwise of the damage detection scheme. The threshold estimation scheme used
795 in this paper is discussed below. The error metric considered is the exceedance rate of
796 the estimated threshold over the GMM density model observed on the held-out fold of
797 the cross validation procedure.

798 *4.3. Detection thresholds*

799 An important aspect of any novelty detection scheme is the definition of the detec-
800 tion threshold. The detection threshold defines the decision boundary between which
801 observations are classed as normal or abnormal. In this case, because the negative
802 log-likelihood of a Gaussian mixture is being used as a novelty index, the detection
803 threshold will be given in terms of this quantity. There is no ultimate gold standard
804 method for estimating an appropriate novelty threshold, and there is a wide variety of
805 approaches to this problem in the literature [65]. In the problem at hand in this paper,
806 the primary requirement is that the threshold minimise the number of false positive and
807 false negatives. The other consideration of practical importance in this work is whether
808 the training data, from the undamaged condition, contains outliers. Outliers in a train-
809 ing set will manifest themselves as observations with high novelty indices. If these are
810 taken into consideration for the threshold determination for example, through the use
811 of maxima or percentiles of the training novelty indices, the resulting threshold will
812 be biased by these outliers. This would in turn mean that observations from potential
813 damaged classes may fall under the threshold. On the other hand, placing the thresh-
814 old too close to the majority of training novelty indices will result in an over-sensitive
815 classifier that produces a high number of false positives. In this work, a robust ap-
816 proach to threshold estimation was used, following ideas from Monte Carlo sampling
817 and Extreme Value Statistics (EVS), which have been used in the past as robust means
818 of threshold estimation [51, 66, 67]. When deciding on an appropriate novelty thresh-

819 old, one is interested in the distribution of maximum values of the novelty indices. The
 820 distributions of the exponential family (to which the Gaussian belongs) are generally
 821 poor at predicting extreme values, due to the fact that the tails of the distribution decay
 822 exponentially to infinity. This is not a particularly good representation of the max-
 823 ima (and minima) of physical events. Extreme value theory offers an alternative for
 824 modelling the tails of statistical distributions. It dictates that the extremes will be gov-
 825 erned by either one of three distributions: Gumbel, Frechet or Weibull. The Gumbel
 826 distribution is of particular interest since it is the limiting distribution of the maxima
 827 of Gaussian random variables, so it would be a suitable distribution for modelling the
 828 tails of a Mahalanobis (un-squared) distance. However, when faced with more complex
 829 probability densities of the undamaged class data (as is the case in this application of
 830 bearing monitoring), other novelty indices might be used which may not conform to
 831 the Gaussian assumption. Such is the case of the negative log-likelihood of the GMM,
 832 being used here as a novelty index. In this case, the Generalised Extreme Value (GEV)
 833 distribution offers a solution to modelling the tails of arbitrary probability distributions
 834 [68]. Its use has previously been investigated in the context of SHM [69].

Here, the novelty threshold is defined using a GEV distribution fitted to random
 samples of negative log-likelihoods drawn from the estimated GMM density. The
 threshold estimation procedure follows the Monte Carlo approach outlined in [51],
 except that here, a GEV is used instead of an empirical cumulative distribution func-
 tion when assessing confidence levels. The threshold estimation procedure is outlined
 in Algorithm 1. The idea of it is to draw N_s random samples from a GMM with param-
 eter set θ representing the normal condition, and to evaluate their novelty index (the
 negative log-likelihood). This process is repeated for N_t different trials, in which the
 maxima of each trial is recorded and stored in a vector \mathbf{z} . A GEV distribution is then
 fitted to this vector of maxima, from which a threshold, T , can be estimated using the
 GEV Cumulative Distribution Function (CDF),

$$GEV_{cdf}(z) = \exp \left\{ - \left(1 + \xi \frac{z - \mu}{\psi} \right)^{-1/\xi} \right\}, 1 + \xi \frac{z - \mu}{\psi} > 0 \quad (5)$$

835 where μ and ψ are location and scale parameters respectively and ξ is an additional
 836 parameter which determines the type of distribution the GEV fit belongs to (from the

837 family of Gumbel, Frechet or Weibull).

Algorithm 1 Threshold Estimation

procedure MC-GEV THRESHOLD($\theta, N_s, N_t, C_{onf}$)

for $i = 1 : N_t$ **do**

$\mathbf{y} \leftarrow$ draw N_s samples from GMM with parameter θ

$\mathbf{L} = -\log p(\mathbf{y}|\theta)$

$z_i \leftarrow \max(\mathbf{L})$

end for

 Fit a GEV distribution to vector of maxima, \mathbf{z}

$T \leftarrow \operatorname{argmin} \|GEV_{cdf}(\mathbf{z}) - C_{onf}\|_2^2$

return T

end procedure

838 *4.4. Procedure Summary*

839 To summarise, the procedure for using a Gaussian mixture for damage detection
840 suggested here is as follows:

- 841 1. Select a (damage sensitive) feature vector \mathbf{y} to represent the data from a healthy
842 condition. Any operational and environmental changes should be captured in
843 \mathbf{y} . In this paper, three features are considered from the AE data: hit summary
844 statistics, hit AR coefficients and signal envelopes.
- 845 2. Select a subset of \mathbf{y} to use for training the model, \mathbf{y}_{train} . Select another subset
846 to test the model predictions on: \mathbf{y}_{test} .
- 847 3. Project high dimensional features, \mathbf{y} onto a suitably lower-dimensional domain
848 \mathbf{x} .
- 849 4. Use cross-validation to decide on an appropriate number of clusters, K , for a
850 GMM model, the training data set.
- 851 5. Train a GMM, using the EM algorithm, on the entire \mathbf{x}_{train} set, using the number
852 of clusters K determined from the cross-validation step.

- 853 6. Evaluate negative log-likelihood function on \mathbf{y}_{train} , point-by-point, using equa-
854 tion (3), and set the detection threshold \mathcal{T} using the procedure described in Al-
855 gorithm 1.
- 856 7. Evaluate $-\log \mathcal{L}$ on any new observations, and check whether this falls above or
857 below the detection threshold.
- 858 8. If the model correctly captures the variability of the healthy condition, exceedances
859 of \mathcal{T} indicate damage, or other previously unseen (and possibly benign) changes.

860 5. Experimental results

861 This section details the results of applying the damage detection framework de-
862 scribed in Section 4 to the AE damage sensitive features computed using the methods
863 outlined in Section 3 to the wind turbine bearing experimental set-up, described in
864 Section 2.

865 Recall that damage detection is performed by fitting a Gaussian mixture model to
866 the PCA projection of the damage-sensitive AE features. The results of interest are the
867 quantification of detection performance of the identified GMM model on observations
868 from bearings with increasing levels of damage. The damage levels are outlined in
869 Table 2.1.

870 All GMM models were trained using data from condition UD1, from the first un-
871 damaged bearing, tested on the other half of the UD1 set and then validated on the data
872 set from UD2. Each AE channel is being considered separately, therefore each will
873 generate a different number of hits at different time indices. The envelope spectrum-
874 based features also generate a much less dense quantity of features per trial. In order
875 to keep a consistent indexing and to enable a quantification of false positives and neg-
876 atives across the three different features and four different AE channels, the decision
877 threshold was applied to the maximum novelty index of each 10 second recording. If
878 this falls above the detection threshold, then the trial is classified as normal, otherwise
879 it is classed as outlying. This makes it possible to quantify false positives and neg-
880 atives, because the ground truth of the state of each trial is available. On the other
881 hand, it would be impossible (at this stage) to establish the ground truth as to whether

882 an individual AE feature was generated by a damage or a benign mechanism, within
883 the stream of hundreds of thousands of AE hits. Bearing this in mind, the random
884 reshuffling to divide the UD1 data set into training and testing sets was carried out by
885 reshuffling the trial indices, rather than individual feature indices. This ensures that the
886 same training and testing data is used for all AE channels and features. The splitting
887 of available undamaged condition data into training and testing helps to understand the
888 generalisation performance of the GMM as a density estimator. However, validating
889 the detection performance on a second undamaged bearing measures the robustness of
890 the entire damage detector - including the choice of damage-sensitive features.

891 Whilst the aim was to keep the environmental variations consistent across the dif-
892 ferent bearing conditions in order to carry out a fair comparison, the reality was that
893 the temperatures inevitably fluctuated slightly between the different conditions. Most
894 critical is the oil temperature, as this relates directly to the viscosity and therefore the
895 lubrication regime. All tests were carried out at a low and a high temperature regime,
896 but the precise temperature of each regime varied slightly. In general, operating the
897 gearbox at higher temperatures leads to increased levels of AE activity due to the higher
898 level of asperity contact. Of the two data sets collected on undamaged bearings, UD1
899 was collected at a slightly lower temperature than UD2. Training a novelty detector
900 on the lower temperature data set and validating its predictions on a data set from a
901 higher temperature provides a more robust validation than would be by testing against
902 data from similar temperatures. For this reason, the data set from UD1 was used for
903 training, and UD2 for validation. This presents a harder problem than would be by
904 training on UD2 and validating on UD1.

905 The temperature distributions for UD1 and UD2 are illustrated in Figure 11a, using
906 a kernel density estimate (with a band-width of 4). This shows the density of trials at
907 the different temperatures, split by training, testing and validation sets. Note that this
908 shows the training and validation distributions after applying the random shuffling to
909 select training and testing sets from UD1. Two things are clear. The first is that there
910 are two clear regimes of temperature, shown by the two different modes of the density
911 curves. The second is that, for the validation set, the temperature distribution is shifted
912 towards higher temperatures. This means that UD1 contains a major regime around

913 the 25°C range which is not completely captured in UD2. Conversely, UD2 contains a
914 number of trials above 50°C that are not present in either the training and testing sets
915 of UD1. Figure 11b shows the temperature distributions of the damage condition data
916 sets. Note that the data for the lower damage levels (D1, D2 and D3) were collected
917 at slightly higher temperatures than the training and testing sets, but overall lower than
918 the validation set. The modes of the D1 and D2 - the two subsurface damage sets lie
919 at 50°C, which is still within the reach of the density of the training set. The data for
920 the highest level of damage (D5) was collected at relatively low temperatures, with the
921 mode for the lower temperature regime having a significant level of density in the range
922 between 20°C and 25°C.

923 *5.1. Damage sensitive features*

924 Three damage-sensitive features were computed from the raw AE data in order
925 to carry out a comparison of damage detection performance, as described in Section 3.
926 Two of these were based on individual AE hits: summary statistics and Auto Regressive
927 (AR) model coefficients. One of the features was designed to capture the periodic
928 nature of the amplitude modulation of the AE signals, so the frequency spectra of AE
929 envelopes was used. The hit summary statistics have a relatively low dimensionality
930 of $d = 5$. The model order for the AR coefficients was chosen to be 150, so in this
931 case $d = 150$. The envelope spectra originally had a much higher dimensionality of
932 $d = 2500$, although this was truncated to 150 as it was found that the higher dimensions
933 (belonging to higher frequencies) only contributed in terms of added noise and did not
934 add a significant amount of information.

935 Recall that all features were computed from a discrete-wavelet-transformed domain
936 not directly from the raw data. Only the high frequency component of the single level
937 DWT is used, in order to leave out the low frequency bands where the AE sensors are
938 not resonant. The features are, therefore, only representative of the 250kHz-500kHz
939 frequency band. These features are illustrated in Figure 12 (for the OC top sensor),
940 which shows the median and inter-quantile ranges of the three features (along each col-
941 umn), grouped by different levels of damage. The shaded area represents the percentile
942 regions of $\pm 25\%$ around the median, while the dashed lines represent the extrema; the

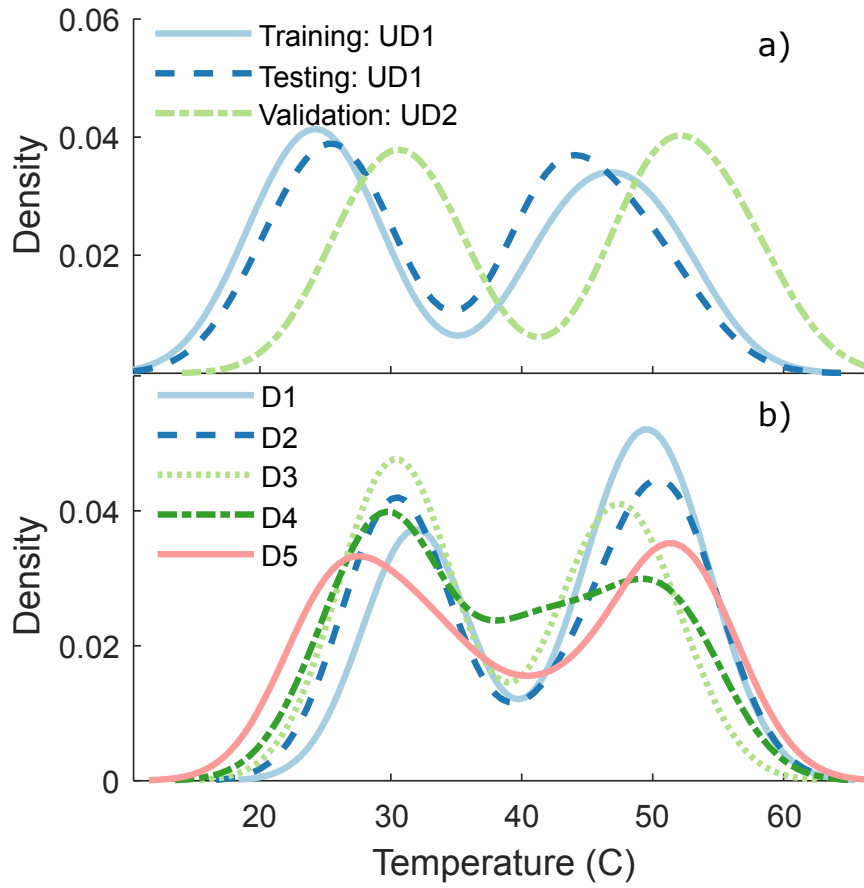


Figure 11: Kernel density of temperature ranges for a) the training, testing and validation trials and b) the damage condition trials.

943 1% and 99% percentiles. The AE features observed in this gearbox are characterised
944 by a large number “extreme” events, which completely mask non-robust measures of
945 location and scatter (a regular mean and standard deviation). The visualisation in terms
946 of inter-quantile ranges of Figure 12 allows a qualitative appreciation of the shape,
947 scatter and extremes of these features. Note that the scale has been adjusted so that the
948 upper 99th percentile of the undamaged feature vectors is visible on the plots. This
949 makes it harder to examine the details of the features, but enables an appreciation for
950 the large difference between the average process and the extreme AE events. The de-
951 velopment from non-damaged to damaged is visually clear across all features in Figure
952 12. A common factor between all three features is an increase in variance of the feature
953 vector as damage progresses, and this is most evident towards state D5.

954 These robust measures of location and spread of the features make it possible to
955 visualise how in this case, damage does not manifest itself as a sharp change to the
956 average feature vector. Instead, the baseline characteristic AE activity remains largely
957 the same, with the addition of extra AE activity that is characteristic to the damage pro-
958 cess. This is true for both hit-based features, but not so for the envelope spectra, which
959 is capturing information across relatively large time-scales. It is evident from Figure
960 12c that the envelope spectra completely shifts its median as damage is introduced and
961 progresses, along with an increase in variability. While this may be a desirable property
962 for this damage-sensitive feature, it should also be noted that its variability in the un-
963 damaged condition is much larger, and this will result in greater variance in any density
964 model fitted to it. This adversely affects detectability, especially at the lower damage
965 levels. Of the three features compared in Figure 12, AE hit statistics, which have the
966 lowest dimension (6), also have the lowest variability in the undamaged state.

967 In order to reduce the dimensionality, PCA was applied to all three damage sen-
968 sitive features, using only the training data (from the undamaged condition) to derive
969 the PCA rotation matrix \mathbf{C} (see Equation (2)). The dimension of the PCA scores was
970 chosen to be $m = 5$ as this captures most of the variance contained in the original full-
971 dimensional features. Figure 14 provides an illustration of the first 5 PCA scores for
972 each damage-sensitive feature. Because of the large quantity of individual feature vec-
973 tors, which greatly vary in scale, a visualisation of the PCA scores of individual feature

974 vectors is not effective. Instead, Figure 14 shows the standard deviation of each PCA
975 score for individual trials (10-second recordings). In order to see how environmental
976 and operational changes have an effect on the PCA scores, the bottom row of Figure 14
977 shows the applied load to the bearing and casing temperature for each trial. Focusing
978 on the two undamaged conditions, it is clear that temperature drives the variance of
979 the first principal component on all three features. In particular, note that the higher
980 temperature range of UD2 has the greatest variance of all UD1 and UD2 trials. The
981 PCA scores of AE hit statistics and envelope spectra show a clear increase in overall
982 variance as damage is introduced, which increases as damage progresses. This effect
983 is marked by large increases in the first score (this represents the direction of greatest
984 variance). The situation is different for the hit-based AR coefficients, where the effect
985 of damage is more markedly seen from the second score onwards. The reason for this
986 is that an AR model is insensitive to scale, so a simple change in the overall energy
987 of the waveform will not yield a different coefficient set. A change in the shape of the
988 waveform, on the other hand will lead to a change in the coefficient set. The fact that
989 damage is more evident from the third score onwards indicates that the AE waveforms
990 generated by a damage mechanism “look” largely the same in this case, but have subtle
991 differences. This is consistent with the visualisation of the median and inter-quantile
992 ranges of AR coefficient vectors provided in Figure 12b. Damage does not completely
993 change the shape of the AR coefficient set, it subtly changes and generates more vari-
994 ability in some of the dimensions.

995 Another way to provide a qualitative view of how the damage process affects the
996 three different damage sensitive features is to visualise how the *probability density* of
997 these features changes with different classes of damage. Figure 13 provides a visu-
998 alisation of a two-dimensional kernel density estimate evaluated on all three damage
999 sensitive features for conditions UD1, D1,D2,D3,D4 and D5¹. The kernel density was
1000 evaluated on the first and second PCA components of each feature for the OC Top sen-
1001 sor location. Note that this is helpful even for the lower-dimensional features such as
1002 AE hit summary statistics, as it allows for visualisations of a decomposition of the data

¹UD2 was omitted since it is visually similar to UD1

1003 in the two directions of greatest variance. The contours of Figure 13 indicate regions of
1004 equal probability density, and only four contours have been drawn, so as to divide the
1005 probability density into four quantile regions. Two major regimes are clear across all
1006 three features for the undamaged state. These represent the high and low temperature
1007 regimes. The first damage level is not strongly visually evident, but the progression of
1008 damage is clear across all three features, albeit in different ways. As damage appears,
1009 there is a clear change in the *shape* of the density. In D2, the second mode that was ev-
1010 ident in UD1 is now masked by a much higher density closer to the core density. This
1011 is the case for all three features. While this is indicative of a change, this will be harder
1012 to detect given the relative closeness of the change to the majority of the data mass.
1013 As surface damage appears, all three features develop regions with high density away
1014 from the core of the data mass. Given this difference in the shape of the probability
1015 density of the features as damage progresses, these observations are bound to generate
1016 high negative log-likelihoods with respect to a model trained on UD1.

1017 5.2. Cross validation analysis of GMM

1018 In order to gain an insight into the output variance and generalisation performance
1019 of the novelty detectors, cross validation was used as described in Section 4.2.2. Two
1020 metrics were considered: the exceedance of the detection threshold as an error metric
1021 and the Bayesian Information Criterion. It is important that this is carried out using
1022 only the data available for training, as this is the scenario one is faced with when
1023 performing a realistic monitoring task. A 10-fold cross validation procedure was used,
1024 considering cluster sizes in the range of $K = 1, \dots, 15$. Using any further than 10 folds
1025 on this data set tends to result in ill-conditioned covariance matrices for some of the
1026 GMM components. The 10-fold cross validation results using misclassification rate are
1027 shown in Figure 15, for all three damage sensitive features. The curves in Figure 15
1028 show the median of the cross-validated output as a solid line, while the greyed-out area
1029 encloses the regions between the 5th and 95th percentiles.

1030 In general, the variance of the misclassification rate decreases as the number of
1031 components of the GMM is increased. The AR model coefficients yield the best per-
1032 formance in terms of variance. The envelope spectra yield the best performance in

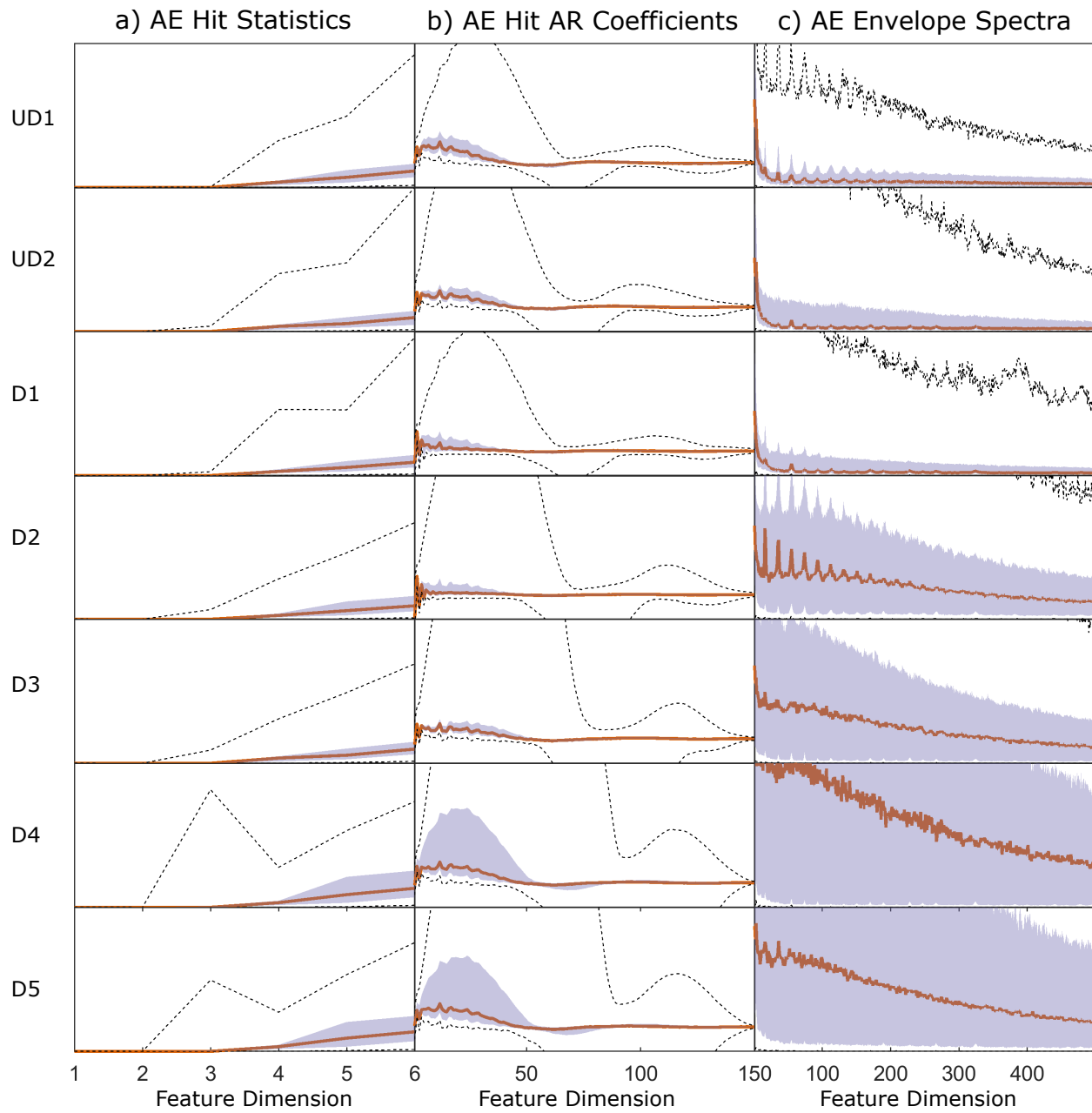


Figure 12: Median and inter-quantile ranges of the three features: a) AE hit statistics, b) AE AR coefficients and c) AE envelope spectra, grouped by different levels of damage. The shaded area represents the percentile regions of $\pm 25\%$ around the median, while the dashed lines represent the extrema; the 1% and 99% percentiles. Note that the vertical axes are not labelled as they correspond to normalised features, but they show the same scale for each feature type.

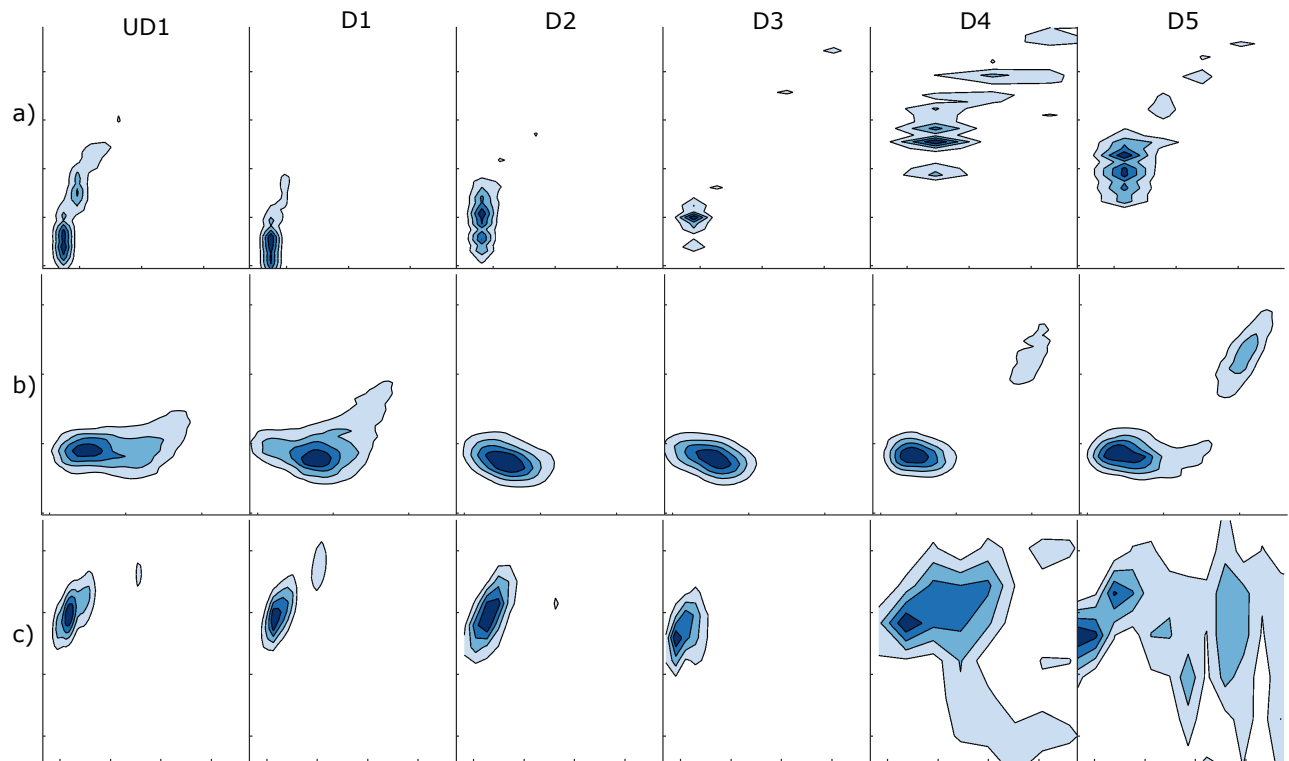


Figure 13: Two-dimensional empirical density estimates of the first and second PCA scores for the three damage sensitive features used: a) hit-based features, b) AR coefficients and c) envelope spectra. These are shown for increasing levels of damage as outlined in Table 2.1

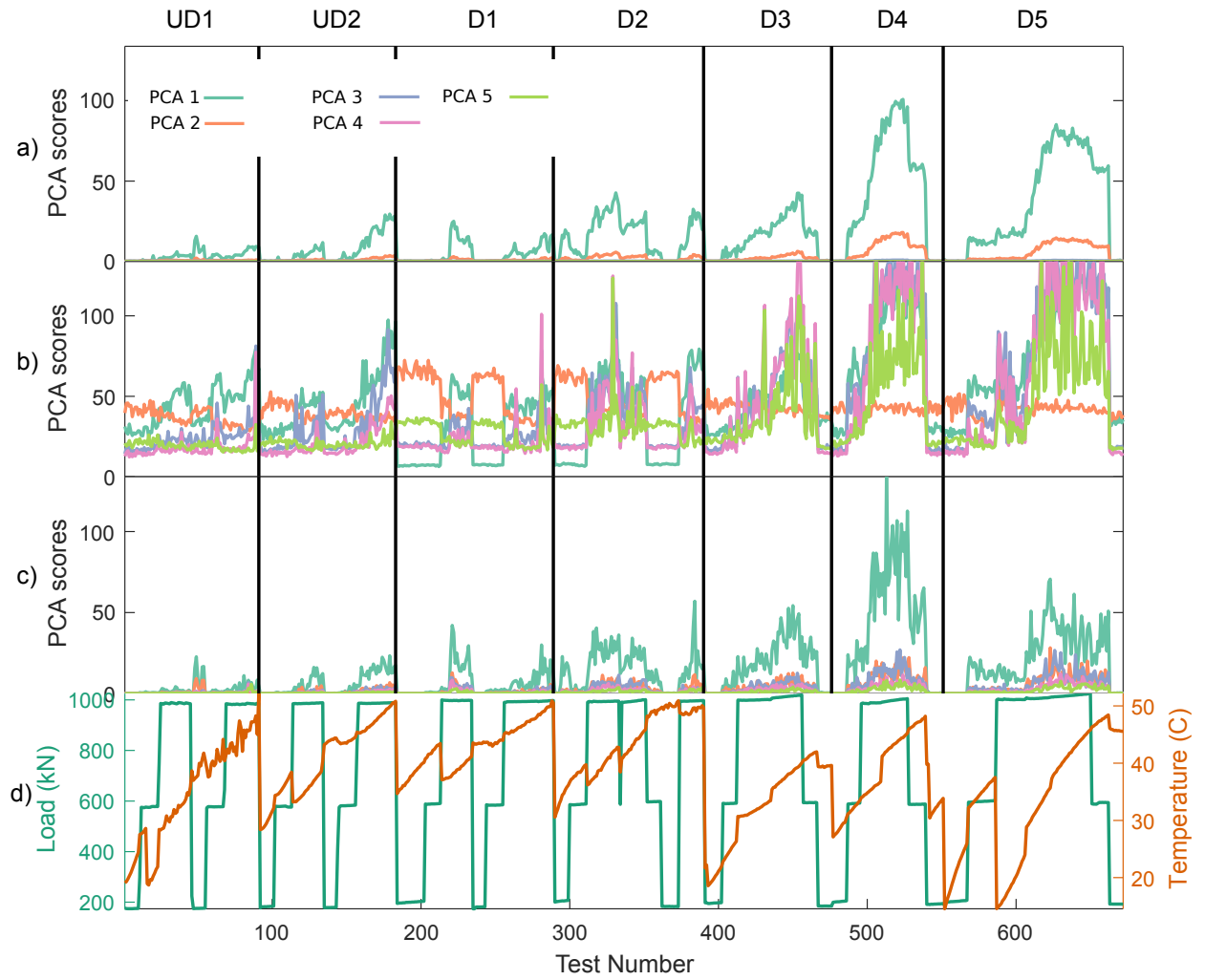


Figure 14: Standard deviation of the first five PCA scores computed per 10-second test for all bearing states. The three damage-sensitive features are shown: a) AE hit statistics, b) AE AR coefficients and c) AE envelope spectra. d) shows the variation in load and temperature throughout these tests.

1033 terms of median misclassification rate, reaching 100% correct classification but with
1034 very high variance due to a few stray folds.

1035 Figure 16 shows the cross validation results, using the BIC as an error metric. Note
1036 that in this figure, the vertical axes, representing the BIC are not shown as each sub-
1037 figure has a different scale. However, it is the trend and the variance that are important.
1038 In general, as the number of clusters increase, the BIC increases, indicating a better fit,
1039 but so does the variance. A high BIC with high variance is indicative of model over-fit.
1040 The selection of the appropriate model order has to balance high BIC scores, low BIC
1041 variance as well as low misclassification rate and low variance around this.

1042 Several observations can be made from examination of Figures 15 and 16. The
1043 first is that while the average BIC increases with increasing cluster numbers (as would
1044 be expected), the misclassification rate does not have such drastic improvements and
1045 tends to converge quickly. The second observation is that all four AE sensor locations
1046 behave in a similar fashion for each of the damage-sensitive features considered. This
1047 is more evident in the misclassification rate. Furthermore, each of the three different
1048 features has a markedly different optimal number of clusters. The misclassification
1049 rate provides greater insight than the BIC on this point. For each feature and AE
1050 sensor location, the model order for the GMM was selected as the first to generate
1051 a misclassification where 95% of the folds (the upper bounds in Figure 15) have a
1052 misclassification under 0.01%. These are marked with vertical lines in Figure 15. It
1053 is interesting to note that the AE envelope spectra generate consistently low *median*
1054 misclassification rates, but with high 95th percentiles. This behaviour is due to one or
1055 two outlying observations, the source of which is likely to be contaminating ambient
1056 noise, which sometimes overpowers the envelope spectrum if this is at much higher
1057 amplitudes than regular AE activity. It is also worth considering that hit-based features
1058 generate in the range of 500 to 1000 observations per trial. Hence outliers tend to
1059 hide well above the 95th percentile. On the other hand, the envelope spectra generate
1060 only 8 observations per 10-second trial, which means that a small number of outlying
1061 observations make it seem like the spread of this error metric is large. Furthermore,
1062 this also means that there is less resolution on the misclassification rate. Considering
1063 this, the criteria for setting the GMM model order on AE envelope spectra features was

1064 that the median reached a 0% misclassification rate.

1065 5.3. Damage detection results

1066 This section presents the results of the damage detection process. The objective is
1067 to quantify the performance of each detector. A *detector*, in this context is a GMM
1068 of a damage-sensitive feature at a sensor node. Each detector has a different model
1069 order, established during the cross-validation procedure described above, and has its
1070 own threshold, established using the GEV procedure of algorithm 1.

1071 After making a decision on the GMM model orders for each sensor location and
1072 damage-sensitive feature, a GMM was trained using the entire training data-set. The
1073 novelty of subsequent observations is assessed by evaluating the Negative Log-Likelihood
1074 (NLL) of each feature vector against the reference GMM. This is described by Equa-
1075 tion (4). Even though the NLL already represents a logarithmic scale of the original
1076 Euclidean distance between the GMM centres relative to their variance, the resulting
1077 NLLs evaluated over the entire range of bearing conditions still results in orders-of-
1078 magnitude difference in scale. For the purposes of visualisation, log NLL is used here,
1079 noting that this is just a practical transformation for visualising results. Figure 17 shows
1080 a kernel density of the log NLL, for all four sensor nodes and three damage-sensitive
1081 features. Each sub-plot in Figure 17 represents a sensor-feature combination, and den-
1082 sities are shown for each bearing condition. Note that all sub-figures in Figure 17 have
1083 been zoomed-in on the vertical axes, to focus on the low-density regions, where dam-
1084 age is most evident. The thresholds identified using the GEV approach are shown as
1085 vertical dashed lines (note the same log-transformation has been applied to the thresh-
1086 old). In this setting, changes to the baseline AE sound-scape should be evident as
1087 regions of higher density of the (log) NLL above the threshold.

1088 As it has already been illustrated in the previous sections, the effect of damage
1089 is different across all three damage-sensitive features being considered. In the case
1090 of hit-based features, when damage is present, the original density of the features is
1091 preserved and additional bursts of energy related to the damage process are generated.
1092 This is evident from Figure 17. In the case of the envelope spectra-based features,
1093 there is a marked shift in overall mass of the density of the NLL toward the right. The

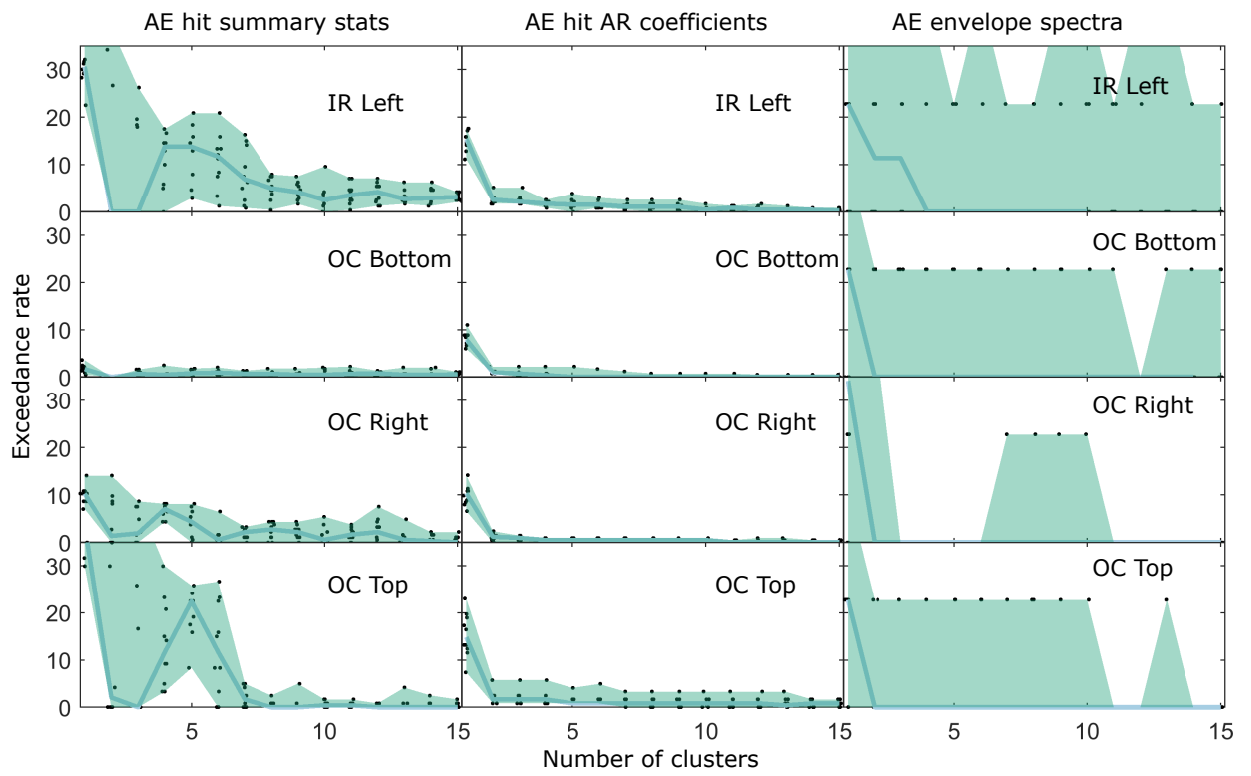


Figure 15: 10-fold cross-validated output of GMM misclassification rate with increasing number of clusters, showing results for the three damage sensitive features at the four sensor locations.

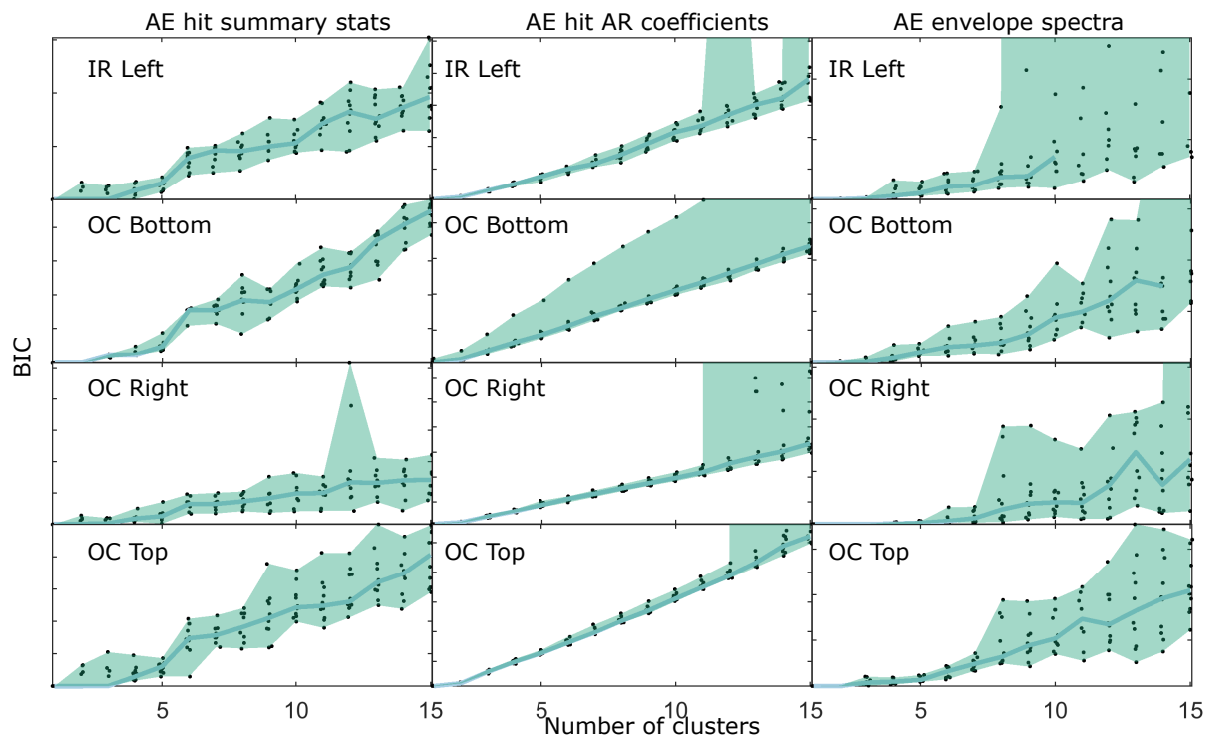


Figure 16: 10-fold cross-validated output of GMM Bayesian Information Criterion (BIC) with increasing number of clusters, showing results for the three damage sensitive features at the four sensor locations.

1094 importance of an appropriate threshold is highlighted in here. Note that there clearly an
1095 appreciable shift in probability mass early on in the development of damage. However,
1096 the detection threshold represents denotes the point beyond which probability mass
1097 of the training features will be negligibly small. If the training features have high
1098 variance, this will lead to larger threshold, hence lower detectability, even when this
1099 may be visible by a visual comparison of the densities of the novelty indices. This is
1100 the case for the AE envelope spectra, where even though there is a considerable shift
1101 in probability mass of the NLL, the original high variance of the features places the
1102 threshold at a relatively high position. This highlights the difference and difficulty of
1103 detecting damage using no prior information of the damage process, as opposed to a
1104 retrospective analysis with knowledge of damaged states.

1105 It is impossible to accurately quantify the detection performance based on individ-
1106 ual hit results, simply because one does not have access to the ground truth of whether
1107 a specific AE burst of energy was generated due to damage or due to a “benign” process
1108 within the gearbox. In this case, the available ground truth is the bearing condition of
1109 each 10-second trial. For this reason, summarising features of the NLL of each individ-
1110 ual trial are used to quantify false positive and negative rates. Two such summarising
1111 features are presented here. Noting that it is the extremes of the NLL that flag novelty,
1112 Figure 18 shows the maxima of the (log) NLL for the OC top sensor location, for each
1113 individual trial and for the three features. The vertical divisions in Figure 18 mark the
1114 different bearing conditions, while the horizontal dashed lines indicate the detection
1115 threshold. As before, the bottom row shows the average applied compressive load and
1116 casing temperature.

1117 For the three features considered, the maxima (log) NLL clearly capture the be-
1118 haviour of the undamaged process, as the majority of the trials of UD1 and UD2
1119 fall under the detection threshold. This is a robust validation of the overall detec-
1120 tion methodology, given there are minimal false positives in the validation set (UD2).
1121 In this case, the only false positives come from the AR coefficient features, and at the
1122 highest temperature observed in UD2. All three detectors fail to identify the presence
1123 of the lowest level of subsurface damage (D1). The second level of subsurface dam-
1124 age, D2, is detectable by the three features, albeit only at high loads. It is reasonable

1125 to conclude that it is the applied load that drives the detection, since tests with high
1126 temperature but low load have low detection rates. Moving upwards in the damage
1127 scale, all three surface damage conditions, D3,D4 and D5 are detectable with the three
1128 features. However, note that they all have different degrees of success at this. In gen-
1129 eral, AE hit statistics and envelope spectra do not detect well under low loads. The AR
1130 coefficients on the other hand begin to detect the damage at the lower loads, although
1131 only for the most severe of the surface damage conditions.

1132 Considering only the maxima (log) NLL of each trial is still prone to an increased
1133 rate of false positives seeing it is likely that even in an undamaged state, rare AE events
1134 will be generated that will drive one of the feature vectors to have a high novelty in-
1135 dex. This has the potential to judge an entire observation set based on one erratic event
1136 while ignoring the information contained in the thousands of other feature vectors con-
1137 tained in that time window. A more robust way of quantifying detection performance
1138 would be through the exceedance rate of feature vectors above the detection threshold.
1139 This, in effect, quantifies the probability mass of the NLL that falls above the detection
1140 threshold, as was illustrated using Figure 17. A positive trial is defined as one where the
1141 exceedance rate above the detection threshold falls above the exceedance rate observed
1142 on the training set. Using this definition, the detection rate for all sensor locations and
1143 feature vectors is given in Figure 19. A positive detection rate on UD1 and UD2 im-
1144 plies a false positive, while the same implies a true positive on the damaged conditions.
1145 Overall, it is possible to conclude that all three feature vectors are capable of detect-
1146 ing from the second level of subsurface damage (D2) onwards, with varying degrees
1147 of success depending on the sensor location. The AE envelope spectra is overall the
1148 worst performing, missing D2 altogether on the OC right location, and with overall low
1149 true positive rates. This is attributed to the high variability of the feature vectors, as
1150 seen in Figure 12. The AR coefficients, while having overall the highest true positive
1151 rate across all locations, also have the highest false positive rate on the validation set
1152 (condition UD2), which is undesirable.

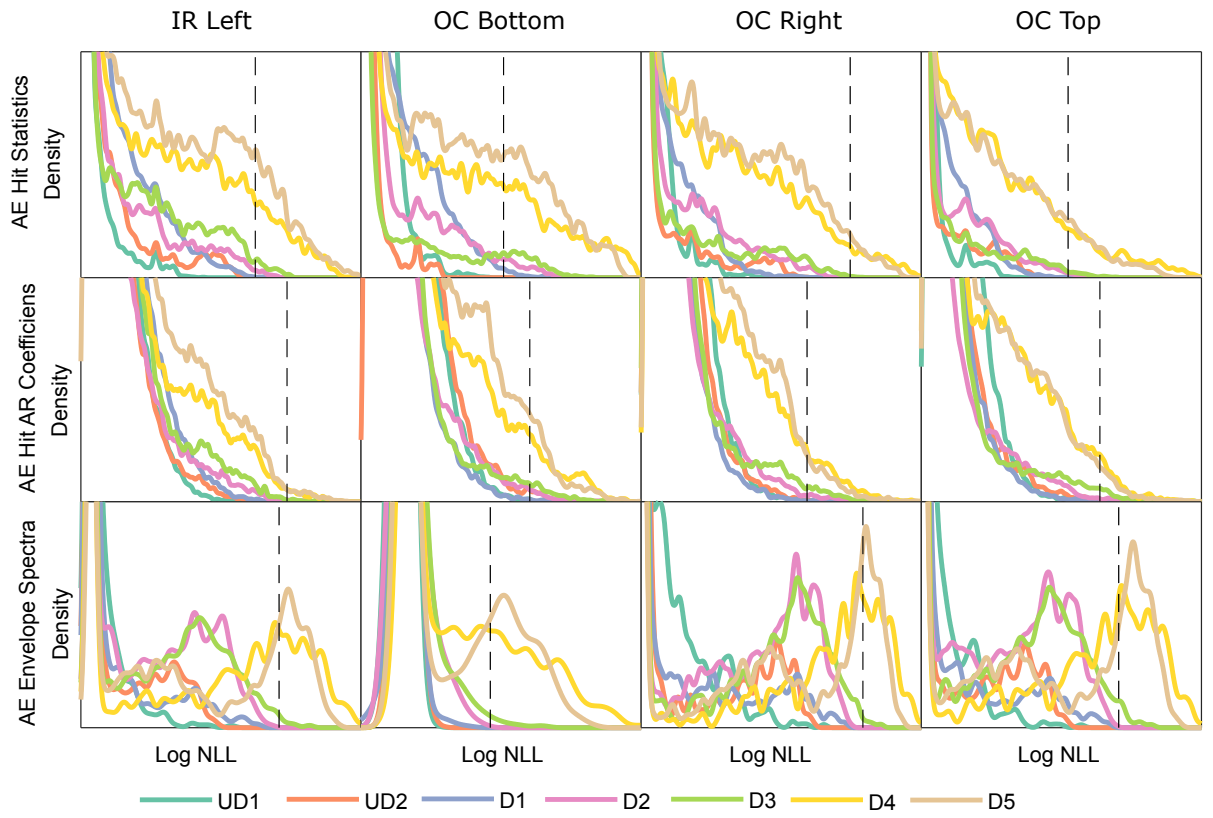


Figure 17: Kernel density of log negative-log-likelihood of GMM model, evaluated on the three different AE features and on the four sensor nodes, for different damage states. The scales have been normalised and adjusted to highlight the tail of the distributions. The vertical dashed lines represent the detection threshold.

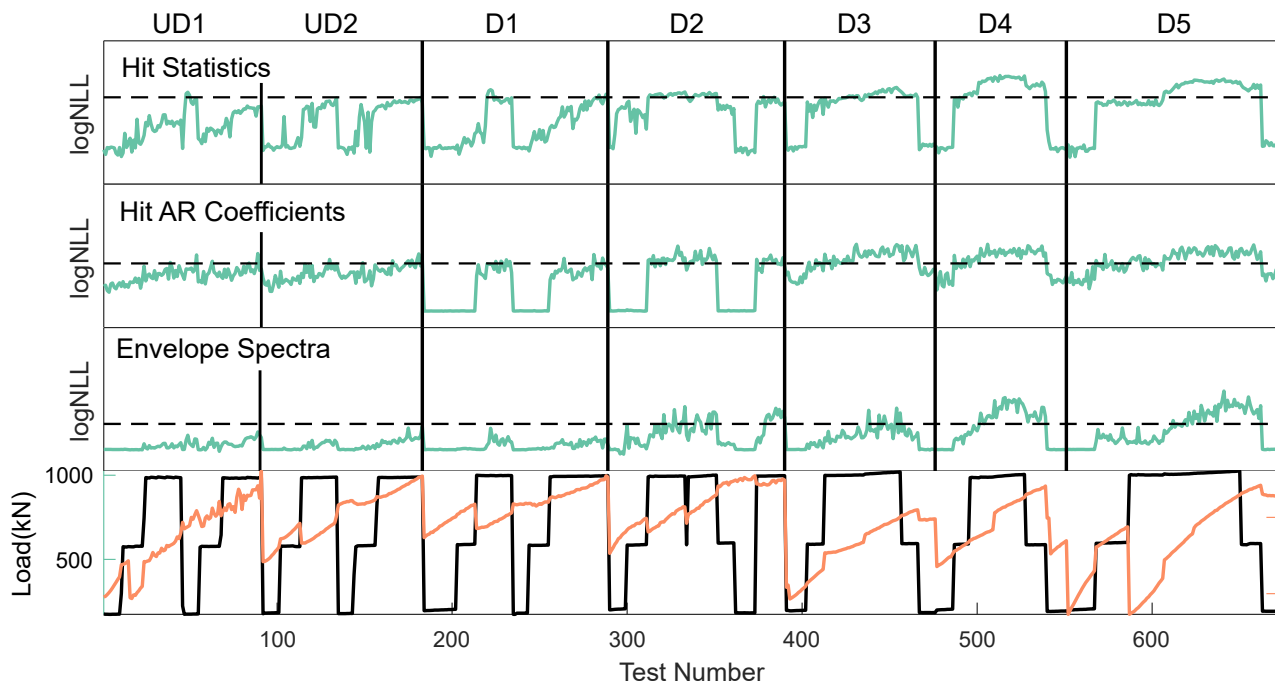


Figure 18: Maxima of the (Log) NLL for each trial, for different bearing conditions and the three different damage-sensitive features (along rows). The horizontal dashed lines represent the detection threshold. Bottom row shows the applied compressive load and casing temperature.

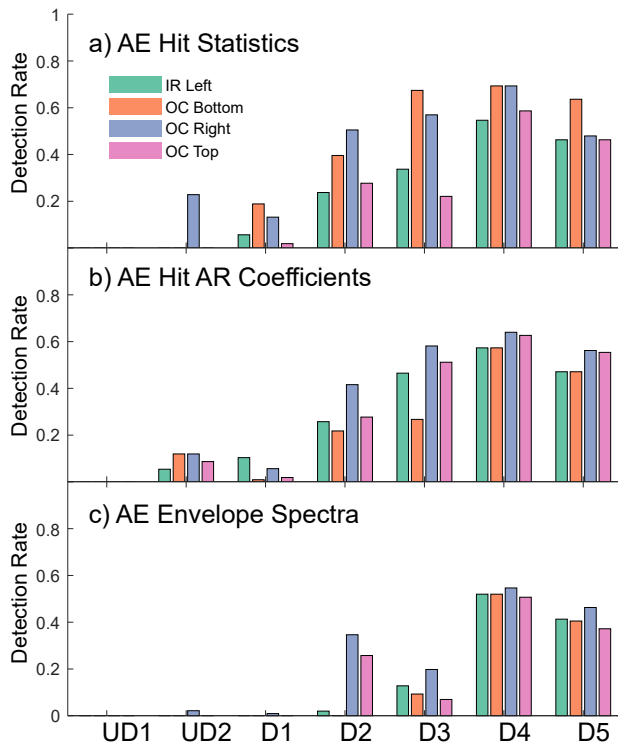


Figure 19: Detection rate for the four different AE sensor locations and the three damage-sensitive feature considered.

1153 **6. Conclusions**

1154 This paper has investigated the problem of detection of sub-surface damage in Wind
1155 Turbine bearings, applying probabilistic modelling to features extracted from Acoustic
1156 Emission measurements. Detecting sub-surface damage on an operational bearing is
1157 not a trivial problem. It is however, an important problem to tackle. Currently, the
1158 majority of gearbox failures can be attributed to bearing failures and this, in turn, ac-
1159 counts for most of the down-time of WTs globally. Damage in bearings starts under the
1160 surface, as a result of Hertzian contact mechanics. While a lot of effort has been placed
1161 in investigating and developing monitoring strategies for surface damage, sub-surface
1162 damage is harder to detect and has received much less attention. This paper focused
1163 on this problem. An emphasis has been placed in using measurements that can be used
1164 in practice, so AE measurements have been taken from practical sensing locations, at
1165 the outer casing of a bearing. An experimental rig was devised in order to replicate the
1166 operational environment of a planetary bearing inside an epicyclic gearbox, as these
1167 suffer the most from early failure before the end of their prescribed fatigue life.

1168 Even though it is sub-surface damage that was of primary interest in this study,
1169 early-stage surface defects were also investigated. A total of five levels of damage were
1170 used, two sub-surface and three early-stage surface defects. This allows to investigate
1171 the detectability of damage throughout its progression. It also helps to build confidence
1172 in the detection scheme by observing that detection rates are higher for larger or more
1173 severe defects.

1174 Measurements were taken at four different locations; one sensing location at the
1175 inner raceway bearing, close to the location of seeded faults and three around the cir-
1176 cumference of the outer casing of the bearing rig. AE data were collected from each
1177 sensor and processed separately in order to evaluate the detectability of the different
1178 levels of defects at each location. Detection of damage was carried out by first extract-
1179 ing three different damage-sensitive features from the raw AE data, and then fitting a
1180 probabilistic model to perform novelty detection. The features chosen in this investi-
1181 gation were hit summary statistics, Auto Regressive (AR) coefficients of the individual
1182 AE hit time histories, and the envelope spectra of the raw AE signals. These features

1183 all capture a different type of information contained in the AE data. The hit-based
1184 summary statistics contain information about the average energy, duration and the dis-
1185 tance each individual stress wave has travelled. The hit AR coefficients on the other
1186 hand provide a greater level of detail as to the spectral characteristics of the waveforms.
1187 Finally, the envelope spectra capture the amplitude modulation of the signal, so any pe-
1188 riodic bursts of energy, which is a characteristic manifestation of damage in dynamic
1189 response data, should be evident in this feature.

1190 The progression of early bearing failure has been illustrated qualitatively, using a
1191 Principal Component Analysis (PCA) projection of the three damage sensitive features
1192 used here. This is shown in Figure 13. In the case of the three different damage-
1193 sensitive features investigated, it is clear that the larger surface-level defects are evi-
1194 dent through a clear change in the probability distribution of the features. This is easy
1195 to spot qualitatively. On the other hand, the change in all damage-sensitive features
1196 arising from sub-surface damage is not necessarily clear from visual examination. Fur-
1197 thermore, in the presence of high levels of noise common in rotating machinery, it
1198 is difficult to establish whether a specific burst of AE energy has been generated by a
1199 damage process or belongs to the background noise. This motivates the need for a prin-
1200 ciple statistical approach to the detection problem. The problem is therefore treated
1201 as one of inference under a probabilistic model. In this particular case, because the
1202 data were gathered under a range of operational conditions, a Gaussian Mixture Model
1203 was used for this task. A GMM was fitted to the damage-sensitive features from an
1204 undamaged bearing and the Negative Log Likelihood (NLL) of the model was used as
1205 a novelty index.

1206 The key result of this investigation is that it is clearly possible to identify sub-
1207 surface damage from a practical measurement location, at the casing of a planetary
1208 gearbox bearing. Using the probabilistic framework presented in this paper, it is pos-
1209 sible to perform such detection under changing environmental and operational condi-
1210 tions. One of the key aspects of AE activity within a bearing environment is that it is
1211 highly dependent on applied load and temperature, as this affects the lubrication prop-
1212 erties. The methods developed here have proved to be robust against these challenges.
1213 This paper validates the approach under an experimental rig where the environment can

1214 be carefully controlled; This is important as temperature, load and lubrication affect the
1215 background AE response and hence the detectability of defects. A clear next step in
1216 this research would be to validate the detectability and the probabilistic approach on
1217 an operational wind turbine, where although there may be less control over the opera-
1218 tional parameters, background noise from gearboxes and other source would be more
1219 realistic.

1220 **7. Acknowledgements**

1221 The authors would like to acknowledge the financial support of Ricardo Innovations
1222 thorough the sponsorship of several PhD programmes as well as the bearing test rig. R.
1223 Dwyer-Joyce would like to acknowledge the EPSRC for funding part of this research
1224 through the fellowship on Tribo-Acoustic Sensors EP/N016483/1. Support for E.J.
1225 Cross through grant number EP/S001565/1 is also acknowledged.

1226 **8. References**

- 1227 [1] F. Oyague, “NREL/TP-500-41160 gearbox modeling and Load simulation of a
1228 baseline 750-kW wind turbine using state-of-the-art simulation codes,” tech. rep.,
1229 National Renewable Energy Laboratory, 2009.
- 1230 [2] W. Yang, P. J. Tavner, C. J. Crabtree, Y. Feng, and Y. Qiu, “Wind turbine con-
1231 dition monitoring: technical and commercial challenges,” *Wind Energy*, vol. 17,
1232 pp. 673–693, may 2014.
- 1233 [3] W. Qiao and D. Lu, “A survey on wind turbine condition monitoring and fault
1234 diagnosis Part I: components and subsystems,” *IEEE Transactions on Industrial*
1235 *Electronics*, vol. 62, pp. 6536–6545, oct 2015.
- 1236 [4] “IEC 61400-1 Wind turbines - Part 1: Design requirements,” 2005.
- 1237 [5] J. Ribrant and L. M. Bertling, “Survey of failures in wind power systems with
1238 focus on Swedish wind power plants during 1997–2005,” *IEEE Transactions on*
1239 *Energy Conversion*, vol. 22, pp. 167–173, mar 2007.
- 1240 [6] D. Coronado and J. Wenske, “Monitoring the oil of wind-turbine gearboxes: main
1241 degradation indicators and detection methods,” *Machines*, vol. 6, no. 2, p. 25,
1242 2018.
- 1243 [7] M. P. Barrett and J. Stover, “Understanding oil analysis: how can it improve
1244 reliability of wind turbine gearboxes,” *Gear Technology*, 2013.
- 1245 [8] C. R. Farrar and K. Worden, *Structural health monitoring: a machine learning*
1246 *perspective*. John Wiley & Sons.
- 1247 [9] T. Howard, *Development of a novel bearing concept for improved wind turbine*
1248 *gearbox reliability*. PhD thesis, The University of Sheffield, 2015.
- 1249 [10] K. Tamada and H. Tanaka, “Occurrence of brittle flaking on bearings used for
1250 automotive electrical instruments and auxiliary devices,” *Wear*, vol. 199, pp. 245–
1251 252, nov 1996.

- 1252 [11] M. H. Evans, “White structure flaking (WSF) in wind turbine gearbox bearings:
1253 effects of ‘butterflies’ and white etching cracks (WECs),” *Materials Science and*
1254 *Technology*, vol. 28, no. 1, pp. 3–22, 2012.
- 1255 [12] T. Bruce, E. Rounding, H. Long, and R. Dwyer-Joyce, “Characterisation of white
1256 etching crack damage in wind turbine gearbox bearings,” *Wear*, vol. 338, pp. 164–
1257 177, 2015.
- 1258 [13] M.-H. Evans, A. Richardson, L. Wang, R. Wood, and W. Anderson, “Confirm-
1259 ing subsurface initiation at non-metallic inclusions as one mechanism for white
1260 etching crack (WEC) formation,” *Tribology International*, vol. 75, pp. 87–97, jul
1261 2014.
- 1262 [14] M.-H. Evans, A. Richardson, L. Wang, and R. Wood, “Effect of hydrogen on
1263 butterfly and white etching crack (WEC) formation under rolling contact fatigue
1264 (RCF),” *Wear*, vol. 306, pp. 226–241, aug 2013.
- 1265 [15] D. Alleyne and P. Cawley, “The interaction of Lamb waves with defects,” *IEEE*
1266 *Transactions on Ultrasonics, Ferroelectrics and Frequency Control*, vol. 39,
1267 pp. 381–397, may 1992.
- 1268 [16] S. Doebling, C. R. Farrar, B. Prime, M, and D. Shevitz, “Damage identification
1269 and health monitoring of structural and mechanical systems from changes in their
1270 vibration characteristics: A literature review,” 1996.
- 1271 [17] K. Worden, C. R. Farrar, G. Manson, and G. Park, “The fundamental axioms of
1272 structural health monitoring,” *Proceedings of the Royal Society A: Mathematical,*
1273 *Physical and Engineering Sciences*, vol. 463, pp. 1639–1664, 2007.
- 1274 [18] J. Baram and M. Rosen, “Fatigue life prediction by distribution analysis of acous-
1275 tic emission signals,” *Materials Science and Engineering*, vol. 41, no. 1, pp. 25–
1276 30, 1979.
- 1277 [19] C. U. Grosse and M. Ohtsu, *Acoustic Emission Testing*. Gardners Books, 2010.

- 1278 [20] C. R. Heiple and R. O. Adams, “Acoustic Emission produced by deformation of
1279 metals and alloys - A review: Part I,” *Journal of Acoustic Emission*, vol. 6, no. 6,
1280 pp. 177–204, 1987.
- 1281 [21] R. Fuentes, T. P. Howard, M. B. Marshall, E. J. Cross, and R. S. Dwyer-Joyce,
1282 “Observations on Acoustic emissions from a line contact compressed into the
1283 plastic region,” *Proceedings of the Institution of Mechanical Engineers, Part J:
1284 Journal of Engineering Tribology*, vol. 230, no. 11, pp. 1371–1376, 2016.
- 1285 [22] K. M. Holford, R. Pullin, S. L. Evans, M. J. Eaton, J. Hensman, and K. Worden,
1286 “Acoustic emission for monitoring aircraft structures,” *Proceedings of the
1287 Institution of Mechanical Engineers, Part G: Journal of Aerospace Engineering*,
1288 vol. 223, pp. 525–532, aug 2009.
- 1289 [23] J. Hensman, K. Worden, M. Eaton, R. Pullin, K. Holford, and S. Evans, “Spatial
1290 scanning for anomaly detection in acoustic emission testing of an aerospace struc-
1291 ture,” *Mechanical Systems and Signal Processing*, vol. 25, no. 7, pp. 2462–2474,
1292 2011.
- 1293 [24] J. Kaiser, *Untersuchungen über das auftreten von geräuschen beim zugversuch*.
1294 PhD thesis, Technical University of Munich (TUM), 1950.
- 1295 [25] V. M. V. M. Baranov, *Acoustic emission in friction*. Elsevier Science, 2007.
- 1296 [26] A. Cockerill, A. Clarke, R. Pullin, T. Bradshaw, P. Cole, and K. M. Holford, “De-
1297 termination of rolling element bearing condition via acoustic emission,” *Proceed-
1298 ings of the Institution of Mechanical Engineers, Part J: Journal of Engineering
1299 Tribology*, vol. 230, no. 11, pp. 1377–1388, 2016.
- 1300 [27] T. Yoshioka, “Detection of rolling contact sub-surface fatigue cracks using acous-
1301 tic emission technique,” *Lubrication Engineering*, vol. 49, no. 4, pp. 303–308,
1302 1993.
- 1303 [28] J. Shiroishi, Y. Li, S. Liang, T. Kurfess, and S. Danyluk, “Bearing condition diag-
1304 nostics via vibration and acoustic emission measurements,” *Mechanical Systems
1305 and Signal Processing*, vol. 11, pp. 693–705, sep 1997.

- 1306 [29] N. Jamaludin, D. Mba, and R. H. Bannister, "Condition monitoring of slow-
1307 speed rolling element bearings using stress waves," *Proceedings of the Institution*
1308 *of Mechanical Engineers, Part E: Journal of Process Mechanical Engineering*,
1309 vol. 215, pp. 245–271, nov 2001.
- 1310 [30] A. Morhain and D. Mba, "Bearing defect diagnosis and acoustic emission," *Pro-*
1311 *ceedings of the Institution of Mechanical Engineers, Part J: Journal of Engineer-*
1312 *ing Tribology*, vol. 217, pp. 257–272, apr 2003.
- 1313 [31] K. R. Al-Balushi, A. Addali, B. Charnley, and D. Mba, "Energy index technique
1314 for detection of Acoustic Emissions associated with incipient bearing failures,"
1315 *Applied Acoustics*, vol. 71, pp. 812–821, sep 2010.
- 1316 [32] B. Eftekharijad, M. Carrasco, B. Charnley, and D. Mba, "The application of
1317 spectral kurtosis on Acoustic Emission and vibrations from a defective bearing,"
1318 *Mechanical Systems and Signal Processing*, vol. 25, pp. 266–284, jan 2011.
- 1319 [33] D. Pandya, S. Upadhyay, and S. Harsha, "Fault diagnosis of rolling element bear-
1320 ing with intrinsic mode function of acoustic emission data using APF-KNN,"
1321 *Expert Systems with Applications*, vol. 40, pp. 4137–4145, aug 2013.
- 1322 [34] P. Kankar, S. C. Sharma, and S. Harsha, "Fault diagnosis of rolling element
1323 bearing using cyclic autocorrelation and wavelet transform," *Neurocomputing*,
1324 vol. 110, pp. 9–17, jun 2013.
- 1325 [35] R. Fuentes, T. Howard, E. J. Cross, R. Harald-Hestmo, T. Huntley, B. Marshall,
1326 Mathew, and R. Dwyer-Joyce, "Detecting damage in wind turbine bearings using
1327 acoustic emissions and Gaussian process latent variable models," in *Proceedings*
1328 *of the 10th International Workshop in Structural Health Monitoring*, (Stanford
1329 University, Palo Alto, CA), 2015.
- 1330 [36] J. R. Naumann, *Acoustic emission monitoring of wind turbine bearings*. PhD
1331 thesis, The University of Sheffield, 2015.

- 1332 [37] M. Elforjani and D. Mba, "Assessment of natural crack initiation and its propa-
1333 gation in slow speed bearings," *Nondestructive Testing and Evaluation*, vol. 24,
1334 pp. 261–275, sep 2009.
- 1335 [38] Z. Rahman, H. Ohba, T. Yoshioka, and T. Yamamoto, "Incipient damage detection
1336 and its propagation monitoring of rolling contact fatigue by acoustic emission,"
1337 *Tribology International*, vol. 42, pp. 807–815, jun 2009.
- 1338 [39] Z. Zhi-qiang, L. Guo-lu, W. Hai-dou, X. Bin-shi, P. Zhong-yu, and Z. Li-na, "In-
1339 vestigation of rolling contact fatigue damage process of the coating by acoustics
1340 emission and vibration signals," *Tribology International*, vol. 47, pp. 25–31, mar
1341 2012.
- 1342 [40] B. Kilundu, X. Chimentin, J. Duez, and D. Mba, "Cyclostationarity of Acoustic
1343 Emissions (AE) for monitoring bearing defects," *Mechanical Systems and Signal
1344 Processing*, vol. 25, pp. 2061–2072, aug 2011.
- 1345 [41] A. Rai and S. Upadhyay, "A review on signal processing techniques utilized in
1346 the fault diagnosis of rolling element bearings," *Tribology International*, vol. 96,
1347 pp. 289–306, apr 2016.
- 1348 [42] C. Jiaa and D. Dornfeld, "Experimental studies of sliding friction and wear via
1349 acoustic emission signal analysis," *Wear*, vol. 139, pp. 403–424, aug 1990.
- 1350 [43] I. Antoniadou, T. Howard, R. Dwyer-Joyce, M. Marshall, J. Naumann,
1351 N. Dervilis, and K. Worden, "Envelope analysis using the Teager-Kaiser energy
1352 operator for condition monitoring of a wind turbine bearing," in *Applied Mechan-
1353 ics and Materials*, vol. 564, pp. 170–175, 2014.
- 1354 [44] SKF, "NU2244 Cylindrical Roller Bearing," 2018.
- 1355 [45] O. Reynolds, "On the theory of lubrication and Its application to Mr. Beauchamp
1356 tower's experiments, including an experimental determination of the viscosity of
1357 olive oil," *Proceedings of the Royal Society, London*, vol. 40, pp. 191–203, 1886.

- 1358 [46] D. P. Hess and A. Soom, "Friction at a lubricated line Contact operating at oscil-
1359 lating sliding velocities," *Journal of Tribology*, vol. 112, no. 1, p. 147, 1990.
- 1360 [47] S. Mallat, *A Wavelet Tour of Signal Processing*. 2009.
- 1361 [48] J. H. Kurz, C. U. Grosse, and H. W. Reinhardt, "Strategies for reliable automatic
1362 onset time picking of acoustic emissions and of ultrasound signals in concrete,"
1363 *Ultrasonics*, vol. 43, no. 7, pp. 538–546, 2005.
- 1364 [49] R. B. Randall, *Vibration-based Condition Monitoring*. Chichester, UK: John Wi-
1365 ley & Sons, Ltd, jan 2011.
- 1366 [50] C. C. Aggarwal, *Outlier Analysis*. New York, NY: Springer New York, 2013.
- 1367 [51] K. Worden, G. Manson, and N. R. J. Fieller, "Damage detection using outlier
1368 analysis," *Journal of Sound and Vibration*, vol. 229, no. 3, pp. 647–667, 2000.
- 1369 [52] N. Dervilis, E. Cross, R. Barthorpe, and K. Worden, "Robust methods of inclusive
1370 outlier analysis for structural health monitoring," *Journal of Sound and Vibration*,
1371 vol. 333, pp. 5181–5195, sep 2014.
- 1372 [53] C. C. Aggarwal and P. S. Yu, "Outlier detection for high dimensional data," *Pro-
1373 ceedings of the 2001 ACM SIGMOD international conference on Management of
1374 data*, p. 46, 2001.
- 1375 [54] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*,
1376 vol. 2nded. Springer series in statistics, 2001.
- 1377 [55] C. M. Bishop, *Pattern recognition and machine learning*. Springer-Verlag New
1378 York, 2006.
- 1379 [56] S. Roweis and Z. Ghahramani, "A unifying review of linear gaussian models.,"
1380 *Neural computation*, vol. 11, no. 2, pp. 305–345, 1999.
- 1381 [57] S. Roweis, "EM Algorithms for PCA and SPCA," *Computing*, vol. 10, no. 13,
1382 pp. 626–632, 1997.

- 1383 [58] C. M. Bishop, "Latent variable models," *Published in Learning in Graphical*
1384 *Models*, pp. 371–403, 1999.
- 1385 [59] E. Figueiredo and E. Cross, "Linear approaches to modeling nonlinearities in
1386 long-term monitoring of bridges," *Journal of Civil Structural Health Monitoring*,
1387 vol. 3, no. 3, pp. 187–194, 2013.
- 1388 [60] E. Figueiredo, L. Radu, K. Worden, and C. R. Farrar, "A Bayesian approach based
1389 on a Markov-chain Monte Carlo method for damage detection under unknown
1390 sources of variability," *Engineering Structures*, vol. 80, pp. 1–10, 2014.
- 1391 [61] J. Kullaa, "Structural health monitoring under nonlinear environmental or opera-
1392 tional influences," *Shock and Vibration*, vol. 2014, 2014.
- 1393 [62] R. Fuentes, *On Bayesian Networks for Structural Health and Condition Monitor-*
1394 *ing*. PhD thesis, 2017.
- 1395 [63] A. Dempster, N. Laird, and D. B. Rubin, "Maximum likelihood from incom-
1396 plete data via the EM algorithm," *Journal of the Royal Statistical Society Series*
1397 *B Methodological*, vol. 39, no. 1, pp. 1–38, 1977.
- 1398 [64] R. E. Kass and L. Wasserman, "A reference bayesian test for nested hypotheses
1399 and its relationship to the Schwarz criterion," *Journal of the American Statistical*
1400 *Association*, vol. 90, p. 928, sep 1995.
- 1401 [65] M. A. Pimentel, D. A. Clifton, L. Clifton, and L. Tarassenko, "A review of novelty
1402 detection," *Signal Processing*, vol. 99, pp. 215–249, jun 2014.
- 1403 [66] H. Sohn, D. W. Allen, K. Worden, and C. R. Farrar, "Structural damage classifica-
1404 tion using extreme value statistics," *Journal of Dynamic Systems, Measurement,*
1405 *and Control*, vol. 127, p. 125, mar 2005.
- 1406 [67] D. A. Clifton, S. Hugueny, and L. Tarassenko, "Novelty detection with multi-
1407 variate extreme value statistics," *Journal of Signal Processing Systems*, vol. 65,
1408 pp. 371–389, dec 2011.

- 1409 [68] M. Aitkin and D. Clayton, “The fitting of exponential, Weibull and Extreme Value
1410 Distributions to complex censored survival data Using GLIM,” *Applied Statistics*,
1411 vol. 29, no. 2, p. 156, 1980.
- 1412 [69] H. W. Park and H. Sohn, “Parameter estimation of the generalized extreme value
1413 distribution for structural health monitoring,” *Probabilistic Engineering Mechan-*
1414 *ics*, vol. 21, no. 4, pp. 366–376, 2006.