

1 **Testing for statistically significant differences in predictions obtained from**
2 **competing creep models**

3 **M. Evans**

4

5 *Institute of Structural Materials, Swansea University Bay Campus, SAI 8EN, Swansea, Wales,*
6 *m.evans@swansea.ac.uk.*

7 **ABSTRACT**

8 It is important to be able to predict the creep life and other creep properties of materials
9 used in power plants and aeroengines. Whilst existing studies have compared different
10 creep models using different measures of predictive performance, none have identified
11 whether these differences are real from a statistical significance perspective. This paper
12 proposes several tests based on a review of the literature. These tests were applied to two
13 creep models using failure time data on 2.25Cr-1Mo steel to develop a recommended
14 approach for practitioners to adopt when selecting a creep model. All such tests
15 concluded that the Evans model produced better long-term life predictions using only
16 short-term data, and that this difference was statistically significant at the 5% significance
17 level.

18 **Keywords:** creep, parametric tests, non-parametric tests, mean percentage squared and
19 absolute errors, statistical significance

20

21 Introduction

22 It is important to be able to predict the creep life and other creep properties for materials
23 used in power plants and aeroengines. When this can be done with a high degree of confidence,
24 the results can potentially be used to justify the continued use of aging power plants beyond
25 their original design lives - as a short-term solution to potential energy gaps for example.
26 2.25Cr-1Mo steel is a main stay alloy used for structural components operating at high
27 temperature within such aging power plants - where the usual service conditions for heater
28 tubes is around 840 K and 35 MPa.

29 There are many creep models available in the literature for carrying out such
30 extrapolations and a good review can be found in Evans [1] and Holdsworth [2]. Such models
31 generally fall into two broad categories – those that aim at predicting points on a creep curve
32 (such as the time to failure or the minimum creep rate) and those that aim to predict the position
33 and shape of the whole creep curve as a function of test conditions. In turn, there has been a
34 recent subdivision of the former models into those that are parametric in nature and those that
35 are semi or non-parametric in nature [3]. Many papers also exist in the literature that perform
36 comparisons of these models. For example, Holdsworth [5] compared several models using
37 2.25Cr1Mo, 9CrMoVNb and 18Cr13NiMo steels, whilst Abdallah et. al. [6], Bueno and
38 Sobrinho [7] have concentrated on 2.25Cr-1Mo.

39 In all these studies, differences in the predictive performance of different creep models
40 were observed. However, in none of these studies, (and as far as the author is aware in any
41 other published paper), was there an attempt to assess whether these differences were
42 statistically significant - rather than simply occurring due to chance. Clearly, this is important
43 for determining the selection of the most appropriate creep model for long term life predictions
44 on a particular material. The objectives of this paper are therefore to review the literature to
45 present a range of statistical tests that can be used for this purpose by future researchers. The
46 paper also aims to illustrate how these tests should be applied using 2.25Cr-1Mo steels and two
47 competing creep models as a test bed. It is not to identify a which of these creep models works
48 best for any other material, or different batches of this material. That would require an extensive
49 analysis of data on other materials, To meet these objectives the paper is structured as follows.
50 The method reviews two performance measures and several statistical tests that are present in
51 the literature. This section also describes two parametric creep models and how their unknown
52 parameters can be estimated. In the results section these tests are applied to the predictions
53 made by these creep models using data on 2.25Cr-1Mo. Conclusions are then drawn.

54 Method

55 *Statistical tests*

56 Evans [9,10] has suggested the use of two statistics for comparing the accuracy of
57 competing creep models. Both are based on the approximate percentage prediction error, with
58 one squaring this error and the other using the absolute values for these errors. Both these
59 approaches prevent under and over predictions offsetting each other in the averaging process

$$60 \quad \text{MPSE} \approx \frac{1}{n} \sum_{i=1}^n [A_i - P_i]^2 = \frac{1}{n} \sum_{i=1}^n e_i^2 \quad (1)$$

$$61 \quad \text{MPAE} \approx \frac{1}{n} \sum_{i=1}^n |A_i - P_i| = \frac{1}{n} \sum_{i=1}^n |e_i| \quad (2)$$

62 where MPSE stands for mean percentage squared error and MPAE stands for mean percentage
 63 absolute error. $A_i = \ln(t_F)_i$ is the natural log of the time at which the i th test specimen fails (of
 64 which there are n such failure times in a data set) and P_i is a creep model's prediction of A_i .
 65 Note that in what follows A_i could equally stand for any other creep property such as the
 66 minimum creep rate, time to a specified strain or even strain itself. The reason for working with
 67 natural logs is that e_i is then approximately equal to the percentage difference between the
 68 actual time to failure and a model's prediction of it. This approximation is very good for
 69 percentages of around 30% or less (and obviously the smaller the better). In this paper e_i will
 70 be referred to as the percentage prediction error. Evans [9,10] has also shown how both these
 71 statistics can be decomposed into random and systematic errors.

72 The use of the MPSE has several limitations, however. The first is that an absolute-
 73 value-based measure, such as the mean percentage absolute error (MPAE), is much more
 74 interpretable. Whilst taking the square root of the MPSE (= RMPSE) helps with interpretation
 75 by converting the MPSE into the same units as e , this still can give very misleading assessments
 76 of a creep models adequacy. This is because the percentage prediction errors are squared, and
 77 this makes the MPSE and the RMPSE very sensitive to the presence of any outliers present in
 78 the data, i.e. very poor predictions. Thus, there is the potential for the MPSE to underestimate
 79 the predictive performance of a given creep model. This is not true of absolute percentage
 80 errors which are more robust to the presence of such outliers. But when it comes to developing
 81 tests of statistical significance, the squaring of errors enables several well know distributions
 82 to be used for determining statistical significance (a difference can be statistically significant
 83 no matter how small the difference).

84 *Parametric based tests*

85 When selecting an appropriate creep model, the obvious approach is to select the creep
 86 model that has the smaller MPSE or MPAE. But then there is a need to go one step further and
 87 determine whether this difference is significant for predictive purposes or simply due to the
 88 specific sample of data collected. Let e_{1i} be the percentage prediction errors obtained using
 89 creep model 1 and e_{2i} the percentage prediction errors obtained using creep model 2. The
 90 Diebold-Mariano [11] test (the DM test) then involves creating a new series called the
 91 differential loss d_i which can be defined in different ways. But for this paper, d_i will be defined
 92 in one of the following two ways

$$93 \quad d_i = (e_{1i})^2 - (e_{2i})^2 \quad \text{or} \quad d_i = |e_{1i}| - |e_{2i}| \quad (3)$$

94 The Diebold-Mariano (DM) test statistic is then defined as

$$95 \quad DM = \frac{\bar{d}}{s_{\bar{d}}} \quad (4a)$$

96 where

$$97 \quad \bar{d} = \frac{1}{n} \sum_{i=1}^n d_i \quad \text{and} \quad s_{\bar{d}}^2 = \frac{\sum_{i=1}^n [d_i - \bar{d}]^2}{n(n-1)} \quad (4b)$$

98 Under the null hypothesis that the population mean value for d is zero, (i.e. $\mu = E[d_i] =$
 99 0 , where $E[d_i]$ reads the expected value for d), DM has an asymptotic standard normal
 100 distribution: $DM \sim N(0, 1)$. Thus, there is a statistically significant difference between the
 101 predictions obtained from two creep models if $|DM| > Z_{\text{crit}}$, where Z_{crit} is the two-tailed critical

102 value from the standard normal distribution. Z_{crit} depends of course on the researchers chosen
 103 level of statistical significance (typically either 1%, 5% or 10%).

104 Diebold and Mariano provided some simulation evidence to suggest that their test
 105 performance was versatile and satisfactory relative to any other test in moderately large
 106 samples over a wide range of situations, including heavy-tailed as well as normal prediction
 107 error distributions. However, they did find the test tends to reject the null hypothesis too often
 108 in small to moderately sized samples (i.e., it is oversized or has lower power). Harvey,
 109 Leybourne, and Newbold [12] found that the following modified DM test (or the HLN test)

$$110 \quad \text{HLN} = \sqrt{\frac{n}{n-1}} DM \quad (5)$$

111 performed better in all scenarios that they looked at. They also found that this modified test
 112 was not as over-sized, and that the performance of the test would be acceptable to practitioners.
 113 Under the null hypothesis that the population mean value for d is zero ($\mu = E[d_i] = 0$), HLM
 114 follows a Student t distribution with $n - 1$ degrees of freedom: $\text{HLN} \sim T(n-1)$. Thus, there is a
 115 significant difference between the model's predictions if $|\text{HLN}| > T_{\text{crit}}$ where T_{crit} is the two-
 116 tailed critical value from the Student t distribution. The standard assumptions required for the
 117 above two tests to be valid are that the d_i are normally distributed and that the variance for d is
 118 constant.

119 Tests that were present in the literature prior to the publication of the DM test, include
 120 those proposed by Ashley, Granger and Schmalensee [13] (the AGS test) and that by Morgan-
 121 Granger-Newbold [14,15] (the MGN test). However, these are tests for equality of the
 122 MPSE's between creep models only. The foundation of the AGS test lies in the decomposition
 123 of the mean percentage squared error (MPSE) based on the work of Granger and Newbold [16].
 124 They showed that

$$125 \quad \text{MPSE} = (\bar{A} - \bar{P})^2 + s_e^2 \quad (6a)$$

126 where s_e^2 is the variance in the percentage prediction errors

$$127 \quad s_e^2 = \frac{\sum_{i=1}^n [e_i - \bar{e}]^2}{n} \quad (6b)$$

128 If MPSE_1 is the mean percentage squared error associated with the n predicted creep
 129 properties obtained from creep model 1 and if MPSE_2 is the mean percentage squared error
 130 associated with the n predicted creep properties obtained from creep model 2, then from
 131 Equations (6a)

$$132 \quad \text{MPSE}_1 - \text{MPSE}_2 = (\bar{e}_1^2 - \bar{e}_2^2) + (s_{e1}^2 - s_{e2}^2) \quad (6c)$$

133 where \bar{e}_1 and \bar{e}_2 are the average percentage prediction errors from creep models 1 and 2
 134 respectively, and s_{e1}^2, s_{e2}^2 their variances. These variances are obtained in the same way as in
 135 Equation (6b) – with e_1 and e_2 replacing e . Equation (6c) can be rewritten using $s_i = e_{1i} + e_{2i}$
 136 and $g_i = e_{1i} - e_{2i}$

$$137 \quad \text{MSPE}_1 - \text{MSPE}_2 = (\bar{e}_1^2 - \bar{e}_2^2) + s_{sg} \quad (7a)$$

138 where s_{sg} is the covariance between s and g

139
$$s_{sg} = \frac{\sum_{i=1}^n [s_i - \bar{s}][g_i - \bar{g}]}{n} \quad (7b)$$

140 Equation (7a) comes about because $(s_{e1}^2 - s_{e2}^2) = s_{sg}$. Now, if model predictions P₂,
 141 are more accurate than model predictions P₁, then MPSE₂ will be smaller than MPSE₁ so that
 142 the terms on the right-hand side of Equation (6a) must in combination be positive. Therefore,
 143 a test of whether the MPSE₂ is smaller than MPSE₁, can be based on whether the combination
 144 of the covariance and the difference between the squares of the mean errors is positive. This
 145 test can be carried out by first estimating the parameters of the regression equation using least
 146 squares

147
$$g_i = \lambda_1 + \lambda_2(s_i - \bar{s}) + u_i \quad (8a)$$

148 where u is an error term with mean zero, and \bar{s} is the average of all the s_i values. The least
 149 squares estimators of the parameters λ_1 and λ_2 are

150
$$\hat{\lambda}_1 = \bar{g} \quad \text{and} \quad \hat{\lambda}_2 = \frac{s_{sg}}{(s_{s-\bar{s}})^2} \quad (8b)$$

151 where $s_{s-\bar{s}}$ is the standard deviation in the mean adjusted values for s_i. From Equations (8a,b),
 152 the hypothesis that there is no difference in the MPSE's implies that λ_1 and λ_2 will be zero,
 153 which can be tested via the null hypothesis, H₀: $\lambda_1 = \lambda_2 = 0$. The alternative hypothesis, that
 154 the predictions P₂ have a smaller MPSE, implies that either or both λ_1 and λ_2 will be positive.
 155 Should the estimated coefficients in Equation (8a) be significantly negative, the forecasts P₂
 156 are not more accurate than P₁. The test that λ_1 and λ_2 is jointly zero is based on the F distribution

157
$$AGS = \frac{(\sum g_i^2 - \sum \hat{u}_i^2)/2}{\sum \hat{u}_i^2/(n-2)} \quad (8c)$$

158 where \hat{u}_i are the observed residuals in Equation (8a) obtained by using the estimated parameter
 159 values of Equation (8b). If the null hypothesis is true, AGS has an F-distribution with 2 and n-
 160 2 degrees of freedom, $AGS \sim F(2, n-2)$. Thus, there is a significant difference between creep
 161 model predictions if $AGS > F_{crit}$, where F_{crit} is the critical value for the F distribution.

162 If creep model predictions are further assumed to be unbiased (so $(\bar{e}_1^2 - \bar{e}_2^2) = 0$), the
 163 AGS test then suggests that the equality of mean percentage squared errors is equivalent to
 164 equality of prediction error variances. This leads to the MGN test. This test starts with the
 165 regression equation

166
$$s_i = \lambda(g_i) + \varepsilon_i \quad (9a)$$

167 The null hypothesis H₀: $\lambda = 0$ states that the sum of the prediction errors from the two
 168 creep models equals random values given by ε_i , i.e., the prediction errors from each model
 169 differ from each other by random unexplained amount ε and so they are essentially the same.
 170 The standard t test can be used to test this hypothesis

171
$$MGN = \frac{\hat{\lambda} - 0}{\sqrt{\frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{(\sum_{i=1}^n \{g_i^2\})(n-1)}}} \quad \text{with} \quad \hat{\lambda} = \frac{\sum_{i=1}^n s_i g_i}{\sum_{i=1}^n g_i^2} \quad (9b)$$

172 where $\hat{\varepsilon}_i$ are calculated from Equation (9a) using the value for λ given by Equation (9b). Under
 173 the null hypothesis that $\lambda = 0$ (i.e. the error in predictions from each creep model are the same),
 174 MGN follows a Student t distribution with $n - 1$ degrees of freedom.

175 The standard assumptions required for the AGS and MGN tests to be valid are that the
 176 residuals u_i and ε_i both have a mean of zero, both have a constant variance, and that these
 177 residuals are uncorrelated with $s_i - \bar{s}$ and g_i respectively. It must be additionally assumed that
 178 the u_i and ε_i are normally distributed. However, Diebold and Mariano provide simulation
 179 evidence showing that the test can be seriously over-sized, even in very large samples, when
 180 the normality assumption does not hold. This drawback could be serious, as heavy-tailed
 181 distributions for prediction errors would seem to be quite plausible in practice.

182 Moreover, this problem persists, and, indeed, becomes worse as the sample size
 183 increases. It is therefore desirable to have a test that is robust to this type of non-normality. The
 184 source of this problem is that the variance of λ in the denominator Equation (9b) is an
 185 inconsistent estimator of the true variance for λ . This in turn is a consequence of the fact that,
 186 although g_i and s_i are uncorrelated under the null hypothesis, they are not independent.

187 Harvey et. al. [12] recommend a modification of the above test to overcome this
 188 problem

$$189 \quad MGN^* = \frac{\hat{\lambda} - 0}{\sqrt{\frac{\sum_{i=1}^n g_i^2 \hat{\varepsilon}_i^2}{(\sum_{i=1}^n g_i^2)^2}}} \quad (9c)$$

190 Although the distribution for MGN^* is no longer exactly a Student t distribution, a
 191 reasonable practical procedure is to compare this test statistic with critical values from the
 192 Student's t distribution with $(n - 1)$ degrees of freedom. Simulations by Harvey et. al. revealed
 193 that the MGN^* test does yield a test of the correct size for large samples. However, in the
 194 smallest samples, the MGN^* test is seriously over-sized even when the error distribution is
 195 normal – more so than the MGN test. Thus, they recommend use of the MGN^* test only when
 196 moderately large samples are available. In the absence of large samples, they recommended a
 197 consideration of non-parametric approaches.

198 *A non-parametric test*

199 Non-parametric procedures are recommended when any of the assumptions behind the
 200 above tests are invalid. Non-parametric tests are said to be distribution free in that they don't
 201 assume a normally distributed population, and so can be used in the presence of non-normality.
 202 To compare the accuracy of two sets of creep property predictions, the sign test (Mendenhall
 203 and Reinmuth [17]) is based on the number of times that predictions from creep model 1 are
 204 better than predictions from creep method 2. Any ties are ignored. As such it can be used to
 205 see if there are difference between two MPSE's or MPAE's. The number of times creep model
 206 1 is better than model 2 is given by the binomial distribution, and under the null hypothesis
 207 that both creep models are equally accurate, the probability parameter p of this distribution is
 208 0.5. That is, if both creep models are equally accurate, the probability of a specific creep model
 209 producing a better prediction than another model is 0.5 for each time a prediction is made. The
 210 binomial distribution with $p = 0.5$ can be written as

$$211 \quad P[X = x] = C_x^n 0.5^x (1 - 0.5)^{n-x} \quad (10a)$$

212 where X is the number of times creep model 1 produces better predictions than model 2, when
 213 n predictions are made. $P[X = x]$ reads the probability that the random variable X will equal a
 214 specific value x (e.g. the probability that creep model 1 produces $x = 10$ better predictions in
 215 the n predictions produced). C_x^n are the number of different ways of getting x superior
 216 predictions from model 1, in the n predictions made.

217 If a large or small value for x is actual observed, but this value for x leads to $P[X = x]$
 218 being very small (and below a selected significance level) when using $p = 0.5$ in Equation
 219 (10a), then this low probability for observing something that has actually happened can only
 220 be explained by p not being equal to 0.5. In which case the null hypothesis is incorrect, and
 221 creep model 1 performs differently to creep model 2 at the selected significance level (whether
 222 its better or worse depends on whether x is small or large). More formally, let x equal the
 223 number of times creep model 1 produces a percentage prediction error smaller than creep model
 224 2

$$225 \quad x = \#e_{1i} < e_{2i} \quad \text{for } i = 1 \text{ to } n \quad (10b)$$

226 No matter how the prediction errors are distributed, this count will always have a
 227 binomial distribution given by Equation (10a) under the null hypothesis that the two creep
 228 models are equally accurate in prediction (assuming independence). In the context of creep
 229 model comparisons, the null hypothesis is always of the form $H_0: p = 0.5$ so that each creep
 230 model is equally likely to produce the best prediction at any one test condition. The alternative
 231 is that $H_a: p > 0.5$. Here H_a says that creep model 1 is more likely to produce the best prediction
 232 at any one test condition. The probability that the null hypothesis is true, the p -value, can be
 233 evaluated using the following version of the binomial distribution

$$234 \quad p\text{-value} = P[X \geq x] = \sum_x^n \frac{n!}{(n-x)!x!} 0.5^x (1 - 0.5)^{n-x} \quad (10c)$$

235 While this sign test has the attraction of being distribution free so that very few
 236 assumptions need to be made, it is important to realise that they ignore some of the information
 237 which is available within the data - the sign test ignores the numerical size of the errors in
 238 prediction for example. This loss of information can result in the power of the sign test being
 239 considerably smaller than those for the parametric tests described above – depending on the
 240 specific circumstances.

241 ***Testing the assumptions behind parametric tests***

242 Normality of u_i and ε_i can be tested using the test statistic proposed by Doornik and
 243 Hansen [18]. This is a more powerful version of the asymptotic version proposed earlier by
 244 Jarque and Bera [19]. Both these approaches involve testing for the presence of skewness and
 245 kurtosis in the least squares residuals. Testing for normality of the u_i in Equation (8a) can be
 246 done using the Jarque and Bera test statistic

$$247 \quad JB = \frac{n(w_1)^2}{6} + \frac{n(w_2-3)^2}{24} \quad (11a)$$

248 where

$$249 \quad w_1 = \frac{\sum_{t=1}^N (\hat{u}_t - \bar{u})^3 / N}{\{\sum_{t=1}^N (\hat{u}_t - \bar{u})^2 / N\}^{2/3}} \quad \text{and} \quad w_2 = \frac{\sum_{t=1}^N (\hat{u}_t - \bar{u})^4 / N}{\{\sum_{t=1}^N (\hat{u}_t - \bar{u})^2 / N\}^2} \quad (11b)$$

250 Under the null hypothesis of no skewness and kurtosis in the residuals, this test statistic
 251 has an asymptotic chi square distribution with 2 degrees of freedom, $JB \sim \chi(2)$. Thus, there is
 252 a significant deviation from normality when $JB > \chi_{crit}$, where χ_{crit} is the critical value for the
 253 chi square distribution. Testing for normality of the ε_i in Equations (9a) can be done by
 254 replacing terms contain u with ε in Equations (11). The Doornik and Hansen [18] version
 255 simply adjusts this test statistic for small samples.

256 Homoscedasticity of the residuals u_i or ε_i can be tested using the test statistics proposed
 257 by White [20] and by Nicholls and Pagan [21] - both of which look to see if the variance in
 258 these residuals, as estimated using the squared residuals, depends on all the explanatory
 259 variables (and their squares) that are in the regression equation. The involves estimating the \hat{u}_t
 260 in Equation (8a) using the λ estimates in Equation (8b). Then \hat{u}_t^2 is regressed on $(s_i - \bar{s})$ and $(s_i - \bar{s})^2$
 261 and the coefficient of determination R^2 extracted. Under the null hypothesis of
 262 homoscedastic residuals, the statistic $W = NR^2$ has a chi square distribution with 1 degree of
 263 freedom, $W_1 \sim \chi(1)$. Thus, there is a significant deviation from homoscedasticity when $W_1 >$
 264 χ_{crit} . Homoscedasticity of the residuals ε_i in Equations (9a) can be tested by first estimating the,
 265 $\hat{\varepsilon}_t$ in Equation (9a) using the λ estimate in Equation (9b). Then $\hat{\varepsilon}_t^2$ is regressed on g_i and $(g_i)^2$
 266 and the coefficient of determination R^2 extracted. Under the null hypothesis of homoscedastic
 267 residuals, the statistic $W = NR^2$ has a chi square distribution with 1 degree of freedom.
 268 $W_2 \sim \chi(1)$. Thus, there is a significant deviation from homoscedasticity when $W > \chi_{crit}$.

269 *Specification of two parametric creep models*

270 The above tests will be applied to two very different creep models to test whether they
 271 produce significantly different predictions of long-term life. The aim is not to establish
 272 superiority of one model over the other – simply to identify whether any statistically significant
 273 differences exist between them. The first model used in this paper is that attributed to Orr-
 274 Sherby-Dorn (OSD) [22] and is given by

$$275 \quad \ln(t_F) = \ln(B) + n \ln(\sigma) + \frac{Q_c}{R} \left(\frac{1}{T} \right) \quad (12a)$$

276 where B and n are further model parameters. The model implies that there is a linear
 277 relationship between log failure time and log stress, and the role of temperature is to shift this
 278 linear relationship in a parallel fashion. Another approach that is increasingly common in the
 279 literature is to normalise stress using the high temperature tensile strength of the material σ_{TS}
 280 with the additional constraint that when $\sigma = \sigma_{TS}$, t_F will equal zero and as s tends to zero t_F
 281 tends to infinity. This approach is rather appealing because provided the tensile strength is
 282 measure at a high enough constant strain rate, σ_{TS} should represent the stress inducing
 283 instantaneous failure upon the application of a constant load in a creep test. Such a constraint
 284 is introduced using an inverted S shaped curve for the relationship between σ/σ_{TS} and t_F at
 285 constant temperature. Thus, Wilshire and Battenbough [23] introduced the creep model

$$286 \quad t_F = A \left[-\ln \left(\frac{\sigma}{\sigma_{TS}} \right) \right]^m \exp \left\{ \frac{Q_c}{RT} \right\} \quad (12c)$$

287 which at constant temperature is the same equation that describes the inverted Weibull
 288 cumulative density function. Yang et. al. and Wang et. al. [24,25] introduced the model

$$289 \quad t_F = A \left[\frac{\sigma}{\sigma_{TS} - \sigma} \right]^{-m} \exp \left\{ \frac{Q_c}{RT} \right\} \quad (12d)$$

290 which at constant temperature is the same equation that describes the inverted Logistic
 291 cumulative density function. Finally, Evans [26] unified these two models by adding an
 292 additional parameter k , so allowing for other S shaped curves that fall within these two special
 293 cases

$$294 \quad t_F = A\tau^m \exp\left\{\frac{Q_c}{RT}\right\} \quad \text{where} \quad \tau = \left[k \left(\frac{\sigma}{\sigma_{TS}} \right)^{-1/k} - k \right] \quad (12e)$$

295 This unification comes about because as $k \rightarrow \infty$, $\tau^m \rightarrow \left[-\ln \left(\frac{\sigma}{\sigma_{TS}} \right) \right]^m$ and when $k = 1$
 296 $\tau^m = \left[\frac{\sigma}{\sigma_{TS} - \sigma} \right]^{-m}$. As a cumulative probability must range between 0 and 1, all the models
 297 contained within Equations (12e) have the constraint $t_F \rightarrow \infty$ as $\sigma/\sigma_{TS} \rightarrow 0$ and $t_F \rightarrow 0$ as
 298 $\sigma/\sigma_{TS} \rightarrow 1$. So, the transformed stress given in Equation (12e) is not some arbitrarily correction
 299 made to stress – instead it is integral to some new approach in the literature to creep life
 300 prediction. Further, these S shaped creep models have been successfully applied to many
 301 different high temperature materials for example, applications to 2.25Cr-1Mo [27], 316
 302 stainless steels [28], and Inconel 740/740H [29].

303 *Estimating these parametric model*

304
 305 In all applications of these models, there is a need to estimate their unknown
 306 parameters, and this is complicated by the fact that this should be done in a piecewise fashion.
 307 This is required because the model parameters will be different over different stress and
 308 temperature ranges to reflect the different creep mechanism operating in such ranges. This
 309 creates a further complication, namely that the test conditions at which such creep mechanisms
 310 change must also be part of the estimation procedure. The following approach estimates such
 311 break points from the data and allows the model parameters to be different either side of these
 312 break points.

313
 314 Wilshire and Whittaker [28], when studying all the batches of 2.25Cr-1Mo steel in the
 315 NIMS creep data base, and using Equation (12c) identified two distinct breaks and so three
 316 distinctly different creep mechanisms. To explain this, they suggested that at the highest
 317 stresses creep takes place through the generation and movement of new dislocations, formed
 318 at appropriate sources, since the yield stress of the material is exceeded. These new dislocations
 319 that are continuously generated due to the high stress, leads to large net movement of atoms,
 320 and so contributes to high creep rates making failure times very sensitive to changes in stress
 321 (and so large values for m in Equation (12c)). This creep occurs largely because of these
 322 dislocations moving within the grains, and under these circumstances Q_c is expected to be high
 323 and equivalent to that for self-diffusion in bainitic matrices. At intermediates stresses, where
 324 the stress is below the yield stress, these authors suggest creep occurs not by the generation of
 325 new dislocations but by the movement of the dislocations pre-existing in the as received
 326 bainitic microstructure. With fewer atoms moving, the creep rate slows (and so m falls in
 327 value). Provided the movement of preexisting dislocation is predominantly within grains the
 328 activation energy will remain unchanged.

329
 330 However, Whittaker and Harrison [30] have suggested that at these intermediates
 331 stresses creep takes place predominantly within the grain boundary zones rather than in the
 332 grains, and this would reflect itself in a lower activation energy for creep representing diffusion
 333 along dislocations and grain boundaries. They did not formally test for this and so this paper

334 will use an estimation strategy that does. Then at the lowest stresses, Wilshire and Whittaker
 335 state that the bainitic regions transform to ferrite and coarse molybdenum carbide particles due
 336 to the long-term nature of the tests conducted at the very low stresses. This they argued result
 337 in creep once again taking place within the grains where the activation energy would be higher
 338 and equal to that for self-diffusion. In comparison to the Wilshire and Whittaker study, Brear
 339 [31] identified only two distinct ranges of test conditions when using Equation (12a), and he
 340 suggested that this change may be due to the high amounts of oxidized material seen on the
 341 failed creep specimens at the very lowest stresses – the result of prolonged testing. This
 342 reduction in the number of implied creep mechanisms will also be tested in the results section
 343 below.

344
 345 The presence of three different stress regimes suggested by the Wilshire and Whittaker
 346 study can be accommodated for in the Evans model by using the following modification of
 347 Equation (12e)

$$349 \ln[t_F] = \ln[A] + m \ln[\tau] + m_1 \max[0, \ln(\tau) - \ln(\tau_1)] + m_2 \max[0, \ln(\tau) - \ln(\tau_2)] + Q_c \frac{1}{RT} + m_3 \frac{D}{RT}$$

350 (13a)

351 where $\tau_1 < \tau_2$ and these are the values for τ where there is a change in creep mechanism and D
 352 is a dummy variable that equals 1 when $\tau_1 < \tau < \tau_2$ and zero otherwise. m_1 to m_3 are three
 353 additional parameters that require estimation. Thus, in the high stress regime where $\tau < \tau_1$, the
 354 model simplified to

$$356 \ln[t_F] = \{\ln[A]\} + (m) \ln[\tau] + \frac{Q_c}{RT}$$

357 (13b)

358 and in the low stress regime where $\tau > \tau_2$, the model becomes

$$360 \ln[t_F] = \{\ln[A] - m_1 \ln(\tau_1) - m_2 \ln(\tau_2)\} + (m + m_1 + m_2) \ln[\tau] + \frac{Q_c}{RT}$$

361 (13c)

362 and in the intermediate stress regime where $\tau_1 < \tau < \tau_2$, the model becomes

$$364 \ln[t_F] = \{\ln[A] - m_1 \ln(\tau_1)\} + (m + m_1) \ln[\tau] + \frac{Q_c + m_3}{RT}$$

365 (13d)

366 Thus in the intermediate stress regime the activation energy drops from Q_c (which
 367 represents that for self-diffusion) to $Q_c + m_3$, as the expectation is for m_3 to be negative based
 368 on the Wilshire and Harrison [30] paper. In the high and low stress regimes the activation
 369 energy will be higher and estimated by Q_c . But if $m_3 = 0$, that will support the view put forward
 370 by Wilshire and Whittaker [28] that dislocation movement remains predominantly within the
 371 grains in the intermediate stress range. When jumping from the high to intermediate stress
 372 regime parameter m changes to a value of $m + m_1$, with the expectation (based on the analysis
 373 by Wilshire and Whittaker and Whittaker and Harrison) that m_1 will be negative so that the
 374 value for n diminishes as the stress falls (due to fewer moving dislocations). Parameter m
 375 changes to $m + m_1 + m_2$ in the low stress regime. If the value for m_2 is such that $m + m_1 + m_2$
 376 approximates m , that would support the view that microstructural degradation is taking place
 377 so increasing m . But, if the value for m_2 is such that $m + m_1 + m_2$ drops off towards a value of 1,

378 that would be consistent with Nabarro-Herring creep. This then represents a piecewise fitting
 379 procedure where the parameters are different over different stress ranges.

380

381 To estimate values for the parameters in Equations (13) the following procedure can be
 382 used. First k is set equal to 1, with τ_1 and τ_2 set equal to random values within the limits of the
 383 largest and smallest values for τ within the sample of data. This enables a τ value to be
 384 associated with each failure time in the NIMS creep data base so that all the variables on the
 385 right-hand side of Equation (13a) are fully quantified allowing $\ln(t_F)$ to be regressed on $\ln(\tau)$,
 386 $\max[0, \ln(\tau) - \ln(\tau_1)]$, $\max[0, \ln(\tau) - \ln(\tau_2)]$, $1/RT$ and D/RT to obtain least squares
 387 estimates for A , m , $m_1 - m_3$ and Q_c . This regression equation has a residual sum of squares
 388 (RSS) associated with it - that of course is minimised for the chosen value for k , τ_1 and τ_2 by
 389 the linear least squares procedure. Then a generalized reduced gradient non-linear search
 390 technique that uses centralised numerical derivatives is used to search for values of k , τ_1 , and
 391 τ_2 that minimise the RSS associated with the regression line given by Equation (13a). This non-
 392 linear search is carried out using Excel's Solver [32] subroutine.

393

394 The presence of three different stress regimes suggested by Wilshire and Whittaker study
 395 can be accommodated for in the OSD model by using the following modification of Equation
 396 (12a)

397

$$398 \ln[t_F] = \ln[B] + n \ln[\sigma] + n_1 \max[0, \ln(\sigma) - \ln(\sigma_1)] + n_2 \max[0, \ln(\sigma) - \ln(\sigma_2)] + Q_c \frac{1}{RT} + n_3 \frac{D}{RT}$$

399 (14a)

400 where $\tau_1 < \tau_2$ and these are the values for τ where there is a change in creep mechanism and D
 401 is a dummy variable that equals 1 when $\tau_1 < \tau < \tau_2$ and zero otherwise. n_1 to n_3 are three
 402 additional parameters that require estimation. Thus, in the high stress regime where $\tau < \tau_1$, the
 403 model simplified to

404

$$405 \ln[t_F] = \{\ln[B]\} + (n) \ln[\tau] + \frac{Q_c}{RT} \quad (14b)$$

406

407 and in the low stress regime where $\sigma > \sigma_2$, the model becomes

408

$$409 \ln[t_F] = \{\ln[B] - n_1 \ln(\sigma_1) - n_2 \ln(\sigma_2)\} + (n + n_1 + n_2) \ln[\sigma] + \frac{Q_c}{RT} \quad (14c)$$

410

411 and in the intermediate stress regime where $\sigma_1 < \sigma < \sigma_2$, the model becomes

412

$$413 \ln[t_F] = \{\ln[B] - n_1 \ln(\sigma_1)\} + (n + n_1) \ln[\sigma] + \frac{Q_c + n_3}{RT} \quad (14d)$$

414

415 This then represents a piecewise fitting procedure where the parameters are different
 416 over different stress ranges. The parameters of Equation (14a) can be estimated in the same
 417 way as for the Evans model by searching, using Excel's Solver, for the optimal break points
 418 (value for where the creep mechanism changes - σ_1 and σ_2)

419

420 **Results: An application of statistical tests for 2.25Cr-1Mo steel**

421 *Data sources*

422 This paper makes use of the creep failure times contained within Creep Data Sheets 3B
423 & 50A, published by the Japanese National Institute for Materials Science (NIMS) [33-34].
424 This has extensive data on twelve batches of 2.25Cr-1Mo steel where each batch has a different
425 chemical composition that underwent one of four different heat treatments - details of which
426 are given in [33]. This paper makes use of just one of these batches, the MAF batch, which was
427 in tube form that had an outside diameter of 50.8mm, a wall thickness of 8mm and a length of
428 5000 mm with a chemical composition as shown in **Table 1**. Specimens for creep testing were
429 taken longitudinally from this material. Each test specimen had a diameter of 6mm with a gauge
430 length of 30mm. The creep tests were obtained over a wide range of conditions: 400 MPa -
431 22MPa and 723K – 923K. **Figure 1** plots the failure times obtained for this MAF batch at the
432 different stresses and temperatures used.

433 **Figure 1.** Relationship between stress, temperature, and time to failure for the MAF batch of
434 2.25Cr-1Mo steel contained within NIMS creep data sheets 3B &50A [33-34].

435 The same NIMS data sheet also contains results from high temperature tensile testing.
436 The tensile strength measurements shown in the first row of Table 2 were obtained using a
437 constant strain rate of $0.0013s^{-1}$. The models defined by Equations (12c-e) require σ_{TS} to
438 represent that stress that will induce -in practical terms - instantaneous creep rupture when it is
439 instantly applied (so inducing a very high strain rate) to a specimen in a creep test. It follows
440 therefore that σ_{TS} must be obtained from a high temperature tensile test that is conducted at a
441 sufficiently high strain rate – one high enough to make the tensile strength values independent
442 of the chosen strain rate. However, a recent paper by Evans [35] has revealed that the NIMS
443 strain rate of $0.0013s^{-1}$ is not sufficiently high enough to meet this criteria, because data on the
444 same 2.25Cr-1Mo steel obtained by Bueno and Sobrinho [7] shows the tensile strength
445 measurements remain strain dependent above a strain rate of $0.0013s^{-1}$. The steel tensile tested
446 by Bueno and Sobrinho was supplied in plate form with 25.4 mm thickness, according to
447 ASTM A 387, grade 22, in the normalized and tempered condition, with a chemical
448 composition shown in the second row of Table 1. Metallographic analysis indicated the
449 presence of 30% bainite and 70% ferrite in the as-received condition. The specimens for tensile
450 testing were extracted from the rolling direction and a gauge length $L_0 = 25$ mm and an initial
451 diameter $d_0 = 6.25$ mm were used for all specimens. The constant strain rate tests were carried
452 out in a servo-hydraulic 8802 model INSTRON machine.

453 **Table I. Chemical Composition for Different Batches of 2.25Cr- 1Mo Steel (%/wt)**

Source	Cr	Mo	C	Si	Mn	P	S	Ni	Cu	Al
MAF [33-34]	2.46	0.94	0.1	0.23	0.43	0.011	0.009	0.008	0.07	0.005
Bueno & Sobrinho ^[7]	2.09	1.08	0.097	0.32	0.50	0.007	0.002	0.03	0.01	0.05
Balance Fe.										

454

455 Whilst the chemical composition of the specimens tested by Bueno and Sobrinho (Table
456 1) are slightly different from those associated with the NIMS MAF batch, these differences are
457 not enough to explain the extent to which the Bueno and Sobrinho tensile strength
458 measurements obtained at the much higher strain rate of $0.013s^{-1}$ differ from the NIMS data

459 obtained at 0.0013s^{-1} . For this reason, the results shown below for the Evans model make use
 460 the NIMS stresses normalised using the Bueno and Sobrinho tensile strength measurements
 461 shown in Table 2.

462 **Table 2. Tensile strength measurements made on two different batches for 2.25Cr- 1Mo**
 463 **Steel (MPa)**

Data Source	723K	748K	773K	798K	823K	873K	923K
MAF [33-34]	448	423	394	380	365	283	216
Bueno & Sobrinho ^[7]	534*	513*	486	471*	456	413	360

*Interpolated values based on the temperature dependency of the tensile strength

464

465 *Estimating the Evans model*

466 The estimation technique described above is applied to the NIMS failure times with
 467 values of 10,000h or less, i.e. to the solid data points in Figure 2(a). Failure times above 10,000h
 468 are shown as open circles. It was found that the RSS is minimised when $k = 1$ suggesting the
 469 Yang et. al. model [24] represents the creep data better than the Wilshire model (hence the
 470 vertical axis shows $\ln[(\sigma_{\max}/\sigma) - 1]$ – which is what τ in Equation (12e) collapses to when $k =$
 471 1). The model’s predictions are shown by the solid segmented line and this fits the short-term
 472 data (failure times less than 10,000h) very well– with an R^2 value of 98.04%. The model yields
 473 an activation energy of 309 kJmol^{-1} in the high and low stress regions - which is quite close to
 474 that for lattice self-diffusion in this material (350 kJmol^{-1}). This value is therefore consistent
 475 with dislocation creep controlled by lattice self-diffusion in these stress regions. Whilst the
 476 activation energy is slightly different in the intermediate stress region, this difference is not
 477 statistically significant. This suggests that in this region dislocation movement is still
 478 predominantly within grains, but because the stresses are below the yield stress, it is the
 479 movement of preexisting rather than new dislocations – hence the observed increase in the
 480 value for m and mainly unchanged value for Q_c in this region. Hence the Whittaker and
 481 Harrison [30] suggestion that creep is predominantly the result of dislocation movement
 482 between grain is not fully supported by the data.

483 The predictions of failure times beyond 10,000h are summarised in Figure 2(b). Notice
 484 how close the trend line fitted to the data points in this Figure is to the 45° line where actual
 485 and predicted values are identical. The equation of this trend line picks up systematic errors
 486 made by the creep model in predicting time to failure, whilst the scatter of the data points
 487 around the trend line picks up the random prediction errors that the model makes. With the
 488 power term of predicted time being almost unity, the trend line reveals that this creep model
 489 has systematic prediction errors averaging 5.9% ($1/0.944$) above the actual failure times.

490 **Figure 2.** Showing (a) the Evans representation of failure times using σ_{TS} from Bueno and
 491 Sobrinho, where the model is estimated using $t_F < 10,000\text{h}$ and, (b) actual v predicted t_F values
 492 beyond 10,000h.

493 *Estimating the OSD model*

494 Like for the Evans model above, the parameters of Equation (14a) were estimated using
 495 only failure times of 10,000h or less. The model’s predictions are shown by the solid segmented
 496 line and the fit to the short-term data (failure times less than 10,000h) is again good – with an

497 R^2 value of 97.30% - slightly below that for the Evans model. There appears to be two breaks
498 in this prediction line at a log stress of 4.44 and 4. Both breaks are statistically significant and
499 so the data does not support Brear's [31] claim of a single break. The model yields an activation
500 energy of 383 kJmol^{-1} in the high and low stress regions - which is again quite close to that for
501 lattice self-diffusion in this material (350 kJmol^{-1}). Whilst the activation energy is slightly
502 different in the intermediate stress region, this difference is not statistically significant. In the
503 low stress regime n continues to fall in absolute terms to a value of around 2. This decrease is
504 not consistent with Brear's oxidation argument nor with Wilshire and Whittaker's particle
505 coarsening argument, It is more consistent with a move towards diffusion creep within grains
506 at lower stresses. The predictions of failure times beyond 10,000h are summarised in [Figure](#)
507 [3\(b\)](#). Notice the fitted trend line in this Figure has a slope different from the 45° line, implying
508 a larger degree of systematic error compared to the Evans model. That is, at low failure times
509 the model consistently underpredicts with the opposite occurring at higher failure times.

510 **Figure 3.** Showing (a) the OSD representation of failure times, where the model is estimated
511 using $t_F < 10,000\text{h}$ and, (b) actual v predicted t_F values beyond 10,000 h

512 [Figure 4](#) shows the predictions from the Evans and OSD models in the more familiar
513 stress v time space. The tendency for the OSD model to produce the larger predicted times at
514 a given stress and temperature is clearly seen in this figure.

515 **Figure 4.** Relationship between stress, temperature, and time to failure for the MAF batch of
516 2.25Cr-1Mo steel contained within NIMS creep data sheets 3B &50A [31-32], together with
517 the predicted values from the OSD and Evans models.

518

519 *Statistical tests*

520

521 [Figures 2\(b\)](#) and [3\(b\)](#) reveal differences in the characteristics of the predicted failure
522 times made by each creep model. The Evans model slightly overpredicts all the failure times
523 beyond 10,000h, whilst the OSD model under predicts at the lowest times and over predicts
524 and the higher failure times. If the Evans model is labelled creep model 1 and the OSD model
525 creep model 2 then the data in [Figures 2\(b\)](#) and [3\(b\)](#) produce $\text{RMPSE}_1 = 30.67\%$, $\text{RMPSE}_2 =$
526 46.23% with $\text{MPAE}_1 = 25.31\%$ and $\text{MPAE}_2 = 42.55\%$. Clearly then, the Evans model is
527 producing the most accurate predicted lifetimes, but the question remains as to whether this
528 result is a peculiarity of this sample of data on 2.25Cr-1Mo steel or whether these differences
529 are real. Not also how squaring the prediction errors (leading to the RMPSE) underestimates
530 the predictive accuracy of each model as $\text{RMSPE} > \text{MPAE}$.

531

532 Using $d_i = (e_{1i})^2 - (e_{2i})^2$, gave $\text{DM} = -26.14$ and $\text{HLN} = -26.68$, whilst using $d_i = |e_{1i}| - |e_{2i}|$
533 gave $\text{DM} = -19.71$ and $\text{HLN} = -20.12$. Using a 5% significance level, the value for $Z_{\text{crit}} = 1.96$
534 and the value for $T_{\text{crit}} = 2.06$. So, $|\text{DM}| > Z_{\text{crit}}$ and $|\text{HLN}| > T_{\text{crit}}$. Consequently, the Evans model
535 produces statistically significantly better mean squared percentage errors and mean percentage
536 absolute errors, irrespective of whether the original or modified DM test is used. That is, the
537 chances of the population average value for d being zero is less than 5%. Further, given the
538 sample of 25, and based on the results of Diebold-Mariano [11], these tests reject the null
539 hypothesis only about 0.3 percentage points more than they should.

540

541 Estimation of Equation (8a) yielded $g_i = 0.0801 - 0.3634(s_i - \bar{s})$ with an R^2 value of 16.03%.
542 The AGS statistic came out at 4.39, and with $F_{crit} = 3.42$ at the 5% significance level, it can
543 again be concluded that the Evans model produces a statistically significantly lower mean
544 percentage square error (i.e. the null hypothesis of $\lambda_1 = \lambda_2 = 0$ is rejected at the 5% significance
545 level). The Student t statistic on λ_1 is 0.80 and -2.10 on λ_2 . Consequently, λ_1 is not significantly
546 different from zero, whilst λ_2 is significantly less than 1 - which together imply the predictions
547 from the OSD model are not more accurate than those from the Evans model.

548
549 Recall that this test assumes that the u_i in Equation (8a) are normally distributed with a
550 constant variance. The JB statistic comes out at 4.12, and with $\chi_{crit} = 5.99$ at the 5% significance
551 level it can be concluded that the null hypothesis of normality cannot be rejected at this
552 significance level. Further, the White test for homoscedastic residuals produced a test statistic
553 value of $W = 4.15$, and with $\chi_{crit} = 3.84$ at the 5% significance level it can be concluded that
554 the null hypothesis that the u_i have constant variance can be rejected at this 5% significance
555 level. But with $\chi_{crit} = 6.63$ at the 10% significance level it can be concluded that the null
556 hypothesis that the u_i have constant variance cannot be rejected at this higher 10% significance
557 level. Perhaps then, some degree of caution should be attributed to the results of the above AGS
558 test statistic.

559
560 Estimation of Equation (9a) yielded $s_i = -0.43(g_i)$, with an R^2 value of 15.23%. The MGN
561 statistic associated with this regression comes out at -2.08 and for the modified version MGN^*
562 $= -2.77$. Both $|MGN|$ and $|MGN^*|$ exceed the value of T_{crit} of 2.06 and so once again the null
563 hypothesis that both models produce the same MPSE can be rejected at the 5% significance
564 level. Indeed, the negative value for λ means that the Evans model produces a significantly
565 small MPSE.

566
567 Again, this test assumes that the ε_i in Equation (9a) are normally distributed with a
568 constant variance. The JB statistic comes out at 0.89, and with $\chi_{crit} = 5.99$ at the 5% significance
569 level it can be concluded that the null hypothesis of normality cannot be rejected at this
570 significance level. Further, the White test for homoscedastic residuals produced a test statistic
571 value of $W = 3.99$, and with $\chi_{crit} = 3.84$ at the 5% significance level it can be concluded that
572 the null hypothesis that the ε_i have constant variance can be rejected at this 5% significance
573 level. But with $\chi_{crit} = 6.63$ at the 10% significance level it can be concluded that the null
574 hypothesis that the ε_i have constant variance cannot be rejected at this higher 10% significance
575 level. Perhaps then, some degree of caution should be attributed to the results of the above MGN
576 test statistic.

577
578 As there is an element of uncertainty surrounding the constancy of the variance for u_i and
579 ε_i , it is worth looking at the results from the sign test that do not require this assumption. For
580 the sign test, the Evans model produced a smaller percentage prediction error (measured either
581 as an absolute or squared error) 20 times out of the 25 predictions made. The probability of
582 getting such a count or higher under the assumption that both models are equally accurate ($p =$
583 0.5) is, from the Binomial distribution, 0.046%. That is, such counts are highly unlikely, yet
584 given that such a result has occurred, the only way of explaining this is to conclude that $p >$

585 0.5. Consequently, the Evans model has a higher probability of producing a better prediction
586 of time to failure at any stated test condition compared with the OSD model.

587 **Conclusions**

588 This paper has put forward several statistical tests to assess whether one creep model
589 produced better predictions of creep properties than another. These tests were obtained from a
590 review of the literature – mainly within the field of economic time series forecasting where
591 such approaches are well developed. The tests reviewed were based on either the MPSE or the
592 MPAE, or indeed both. These statistics were then applied to predicted failure times for 2.25Cr-
593 1Mo steel obtained from two competing creep models - the well-known OSD model and a
594 recently proposed model by Evans. The purpose being, not to identify a model that will always
595 be superior for all materials or for different batches of this material, but to suggest a statistical
596 testing procedure that will be useful for practitioners to analyse their own results (it is unlikely
597 that the superiority of the Evans model displayed in this paper generalises to other materials).

598 More specifically, it was found that the Evans model produced both a lower RMPSE
599 and a lower MPAE. Furthermore, all the proposed statistical tests concluded that this observed
600 difference in the MPSE and the MPAE was real and not something that had occurred by chance.
601 Using the 5% significance level all tests concluded the Evans model produced statistically
602 significantly lower MPSE's and MPAE's (although the assumption of homogenous variances
603 required for these tests was only accepted at the 10% significance level). However, the
604 Binomial test, that does not require this assumption, came to the same conclusions as the other
605 tests. These statistical tests can be used by future practitioners to identify the best creep models
606 to work with for a particular high temperature material. Consequently, one area for future
607 research is to apply these methods to different materials and to different creep models.

608

609 Disclosure statement: The authors report there are no competing interests to declare.

610 Data: All data used in this publication are in the public domain: References [26-27]

611 Conflict of Interest: This research was not funded by research council grants or private sponsors
612 and as such there are no financial relationships to declare.

613 **References**

- 614 [1] Evans M. Creep in Materials - Origins, Analysis and Prediction Methodologies. *Journal*
615 *of Materials Education*, 2008.
- 616 [2] Holdsworth SR, Askins B, Baker MA, Gariboldi E, Holmström S, Klenkf A, Ringelf 597
617 M, Merckling Z, Sandstromh GR, Schwienheeri N, Spigarellij S. Factors influencing creep 598
618 model equation selection. *International Journal of pressure Vessels and Piping*, 2008; 5: 80–
619 88.
- 620 [4] Evans M. Semi-parametric estimation of the Wilshire creep life prediction model: an
621 application to 2.25Cr-1Mo steel. *Materials Science and Technology*, 2019; 35(16): 1977-1987
- 622 [5] Holdsworth SR, Merckling G. ECCC developments in the assessment of creep-rupture
623 data. In: Proceedings of sixth international Charles Parsons Conference on engineering issues
624 in turbine machinery, power plant and renewables, Trinity College, Dublin, 16–18 September,
625 (200).
- 626 [6] Abdallah Z, Gray V, Whittaker M, Perkins K. A Critical analysis of the conventionally
627 employed creep lifing methods. *Materials*, 2014; 7(5): 3371-3398.
- 628 [7] José Francisco dos Reis Sobrinho, Levi de Oliveira Bueno, Correlation between creep
629 and hot tensile behaviour for 2.25Cr-1Mo steel from 500°C to 700°C. Part 2: An assessment
630 according to different parameterization methodologies, *Revista Matéria*, 10(3) (2005), 463 –
631 471.
632
- 633 [9] Evans M. Estimating threshold stresses using parametric equations for creep: application to
634 low-alloy steels, *Materials Science and Technology*, 2023; 1-16.
- 635 [10] Evans M. Assessing the predictive performance of creep models using absolute rather than
636 squared prediction errors: An application to 2.25Cr-1Mo steel. *Materials at High*
637 *Temperatures*, Published Online: 16 Oct 2023.
- 638 [11] Diebold FX, Mariano, RS. Comparing predictive accuracy. *Journal of Business and*
639 *Economic Statistics*, 1995; 13: 253-63.
- 640 [12] Harvey D, Leybourne S, Newbold P. Testing the equality of prediction mean squared
641 errors. *International Journal of Forecasting*, 1997; 13: 281-91.
- 642
- 643
- 644
- 645 [13] Ashley LR, Granger CWJ, Schmalensee, R. Advertising and aggregate consumption: an
646 analysis of causality. *Econometrica*, 1980; 48: 1149-68.
- 647
- 648 [14] Morgan WA. A test for significance of the difference between the two variances in a
649 sample from a normal bivariate population. *Biometrika*, 1939-1940; 31: 13-19.

- 650 [15] Granger CWJ, Newbold P., *Forecasting Economic Time Series*, 1977, Academic Press,
651 Orlando, FL.
- 652
- 653 [16] Granger CWJ, Newbold P. Some comments on the evaluation of economic
654 forecasts. *Applied Economics*, 1973; 5: 35-47.
- 655 [17] Mendenhall W, Reinmuth J. E. *Statistics for Management and Economics*, 1982, Boston,
656 Mass.: Duxbury Press.
- 657
- 658 [18] Doornik JA, Hansen H. An omnibus test for univariate and multivariate normality. *Oxford*
659 *Bulletin of Economics and Statistics*, 2008; 70: 927–939.
- 660
- 661 [19] Jarque CM, Bera AK. Efficient tests for normality, homoscedasticity and serial
662 independence of regression residuals. *Economics Letters*, 1980; 6: 255–259.
- 663
- 664 [20] White H. A heteroskedastic-consistent covariance matrix estimator and a direct test for
665 heteroskedasticity. *Econometrica*, 1980; 48: 817–838.
- 666
- 667 [21] Nicholls DF, Pagan AR. Heteroscedasticity in models with lagged dependent variables.
668 *Econometrica*, 1983; 51: 1233–1242.
- 669
- 670 [22] Dorn JE, Shepherd LA. What We Need to Know About Creep. In Proceedings of the
671 STP 165 Symposium on The Effect of Cyclical Heating and Stressing on Metals at Elevated
672 Temperatures, Chicago, IL, USA, 17 June 1954.
- 685 [23] Wilshire B, Battenbough, AJ. Creep and creep fracture of polycrystalline copper *Mater.*
686 *Sci. Eng. A*, 2007, vol. 443, pp. 156–66.
- 687 [24] Yang M, Wang Q, Song XL, Jia J, Xiang ZD. On the prediction of long term creep strength
688 of creep resistant steels. *International Journal of Materials Research*, 2016, 107(2).
- 689 [25] Wang Q, Yang M, Song XL, Jia J, Xiang ZD. Rationalization of Creep Data of Creep-
690 Resistant Steels on the Basis of the New Power Law Creep Equation. *Metall. Mater. Trans. A*,
691 2016, vol. 47A (7), pp. 3479–87.
- 692 [26] Evans M. Testing Model Structure Through a Unification of Some Modern Parametric
693 Models of Creep: An Application to 316H Stainless Steel. *Metallurgical and Materials*
694 *Transactions A, Metall. Mater. Trans. A*, 2020, vol. 51 (2), pp. 697–707.
- 695 [27] Wilshire B, Whittaker M. Long term creep life prediction for Grade 22 (2·25Cr-1Mo)
696 steels. *Materials Science and Technology*, 2011; 27(3): 642-647.
- 697
- 698 [28] Whittaker MT, Evans M, Wilshire B, Long-term creep data prediction for type 316H
699 stainless steel. *Materials Science and Engineering: A*, 2016 ;552: 145-150.
- 700 [29] Chen L, Dong Z, Song XL, Jia, J, Xiang ZD. Determination of Activation Energy and
701 Prediction of Long-Term Strength of Creep Rupture for Alloy Inconel 740/740H: A Method
702 Based on a New Tensile Creep Rupture Model. *Metallurgical and Materials Transactions A*,
703 2022, Vol. 53, pp. 1–5.
- 704

705 [30] Whittaker MT, Harrison W. The evolution of the Wilshire Equations for creep life
706 prediction. *Materials at High Temperatures*, 2014, 31(3), pp. 233-238.

707 [31] Brear JM. A perspective on the Wilshire creep equations. *Strength, Fracture and*
708 *complexity*, 2022; 15(1): 79-98.

709 [32] Microsoft Corporation. Microsoft Excel [Internet]. 2018. Available
710 from: <https://office.microsoft.com/excel>.

711

712 [33] NIMS Creep Data Sheet No. 3B: Data Sheets on the Elevated-Temperature Properties
713 of 2.25Cr-1Mo Steel for Boiler and Heat Exchanger Seamless Tubes (STBA 24), National
714 Research Institute for Metals, Tokyo, Japan.

715

716 [34] NIMS Creep Data Sheet No.50A: Long-Term Creep Rupture Data obtained after
717 Publishing the Final Edition of the Creep Data Sheets, National Research Institute for Metals,
718 Tokyo, Japan.

719

720 [35] Evans M. The Important Role Played by High-Temperature Tensile Testing in the
721 Representation of Minimum Creep Rates Using S-Shaped Curve Models. *Metallurgical and*
722 *Materials Transactions A*, 2023, 54, pp. 4796–4805.

723

724

