



The cardiopulmonary exercise test grey zone; optimising fitness stratification by application of critical difference

Journal:	<i>British Journal of Anaesthesia</i>
Manuscript ID	BJA-2017-00683-JT064.R2
Article Type:	Clinical Investigation
Date Submitted by the Author:	08-Feb-2018
Complete List of Authors:	Rose, George; University of South Wales, Faculty of Life Sciences and Education Davies, Richard; Cardiff and Vale University Health Board, Department of Anaesthetics Davison, Gareth; Ulster University, Faculty of Life and Health Sciences Adams, Richard; Cardiff University, Department of Medicine; Velindre Cancer Centre, Oncology Williams, Ian; Cardiff and Vale University Health Board, Department of Surgery Lewis, Michael; University of South Wales, Faculty of Life Sciences and Education; Cwm Taf University Health Board, Department of Surgery Appadurai, Ian R.; Cardiff and Vale University Health Board, Department of Anaesthetics Bailey, Damian; University of South Wales, Faculty of Life Sciences and Education
Mesh keywords:	Anaerobic threshold, cardiopulmonary exercise test, Risk assessment

1
2
3 **The cardiopulmonary exercise test grey zone; optimising fitness stratification by**
4 **application of critical difference**
5
6
7

8
9 G. A. Rose^{1*}, R. G. Davies², G. W. Davison³, R. A. Adams⁴, I. M. Williams⁵, M. H. Lewis⁶, I.
10 R. Appadurai², and D. M. Bailey^{1*}
11
12
13

14
15 ¹Neurovascular Research Laboratory, Faculty of Life Sciences and Education, University of
16 South Wales, UK; ²Department of Anaesthetics, University Hospital of Wales, Cardiff, UK;
17
18 ³Sport and Exercise Sciences Research Institute, Ulster University, Newtownabbey, NI, UK;
19
20 ⁴School of Medicine, Cardiff University, and Velindre Cancer Centre, Cardiff, UK;
21
22 ⁵Department of Surgery, University Hospital of Wales, Cardiff, UK; ⁶Department of Surgery,
23
24 Royal Glamorgan Hospital, Llantrisant, UK.
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

50 **Running title:** Cardiorespiratory fitness and surgical risk stratification
51
52
53
54
55
56
57
58
59
60

***Correspondence:**

¹Neurovascular Research Laboratory, Faculty of Life Sciences and Education, University of
South Wales, United Kingdom, CF37 4AT. Tel: +44 (0)1443-652296; Fax: +44 (0)1443-
652285; email: damian.bailey@southwales.ac.uk or george.rose@southwales.ac.uk

Abstract

Background. Cardiorespiratory fitness (CRF) can inform patient care, though to what extent natural variation in CRF influences clinical practice remains to be established. We calculated natural variation for cardiopulmonary exercise test (CPET) metrics, which may have implications for **fitness** stratification.

Methods. In a two-armed experiment, critical difference (CD) comprising analytical imprecision and biological variation was calculated for CRF, and thus defined the magnitude of change required to claim a clinically meaningful change. This metric was retrospectively applied to 213 patients scheduled for colorectal surgery. These patients underwent CPET and the potential for misclassification of **fitness** was calculated. We created a model with boundaries inclusive of natural variation (CD applied to oxygen uptake at anaerobic threshold ($\dot{V}O_2$ -AT): 11mL O₂ kg⁻¹ min⁻¹, peak oxygen uptake ($\dot{V}O_2$ peak): 16mL O₂ kg⁻¹ min⁻¹, and ventilatory equivalent for carbon dioxide at AT ($\dot{V}_E/\dot{V}CO_2$ -AT): 36).

Results. The CD for $\dot{V}O_2$ -AT, $\dot{V}O_2$ peak, and $\dot{V}_E/\dot{V}CO_2$ -AT was 19%, 13%, and 10%, resulting in false negative and false positive rates of up to 28 and 32% for **unfit** patients. Our model identified boundaries **for unfit and fit patients**: AT < 9.2 and ≥ 13.6 mL O₂ kg⁻¹ min⁻¹, $\dot{V}O_2$ peak < 14.2 and ≥ 18.3 mL kg⁻¹ min⁻¹, $\dot{V}_E/\dot{V}CO_2$ -AT ≥ 40.1 and < 32.7, between which an area of indeterminate-**fitness** was established. With natural variation considered, up to 60% of patients presented with indeterminate-**fitness**.

Conclusions. These findings support a reappraisal of current clinical interpretation of CRF highlighting the potential for incorrect **fitness** stratification when natural variation is not **accounted for**.

Key words: Anaerobic threshold; cardiopulmonary exercise test; risk assessment.

Introduction

Cardiopulmonary exercise testing (CPET) is a non-invasive procedure to determine the level of cardiorespiratory fitness (CRF) of patients during a progressive exercise challenge to symptom limited maximum. CPET is used as a tool for preoperative assessment of **physical fitness** for intra-abdominal surgery to aid clinical decision-making given its increasingly proven association with post-operative outcome.¹⁻⁷ Furthermore, The American Heart Association recently published a scientific statement promoting cardiorespiratory fitness (CRF) as a clinical vital sign.⁸ Despite increasing support for CPET, the mechanisms underpinning CRF that provide protection require further investigation.

The seminal work of Older and colleagues documented an 18% mortality rate in elderly surgical patients with a pulmonary oxygen uptake at the anaerobic threshold ($\dot{V}O_2$ -AT) of $< 11\text{mL oxygen (O}_2\text{) kg}^{-1}$ (total body mass) min^{-1} compared to 0.8% recorded in patients with a $\dot{V}O_2$ -AT $\geq 11\text{mL O}_2\text{ kg}^{-1}\text{ min}^{-1}$.⁹ Other biomarkers including peak oxygen uptake ($\dot{V}O_2$ peak) $< 15\text{mL O}_2\text{ kg}^{-1}\text{ min}^{-1}$ and ventilatory equivalent for carbon dioxide at AT ($\dot{V}_E/\dot{V}CO_2$ -AT) > 42 have predicted post-operative survival following abdominal aortic aneurysm surgery.² Studies have further attempted to define threshold values in an effort to optimise risk prediction; for example a range of AT values from 9.0 to $11\text{mL O}_2\text{ kg}^{-1}\text{ min}^{-1}$ have been reported,^{4 5 9-12} thus demonstrating that variation is present and that a single cut-point cannot be recommended.

Like most biomarkers, CRF is a dynamic metric subject to natural variation and thus needs to be interpreted with caution. Such variation encompasses both analytical and biological components that collectively contribute to the critical difference (CD) as originally described by Fraser and Fogarty.¹³ The CD represents random variation around a homeostatic point indicative of the change that must occur before a true difference of clinical significance can be claimed. The concept of CD, yet to be applied to clinical CPET variables, emanates

1
2
3 from the field of clinical biochemistry and has been applied to metabolic biomarkers of
4
5 exercise stress and clinical patients.^{14 15}
6

7 The current study reflects the first attempt within the clinical setting to quantify the
8
9 CD of established CPET markers of CRF with corresponding implications for patient
10
11 management. We hypothesise that natural variation is present in markers of CRF, and will
12
13 thus impact upon patient fitness stratification.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

Methods

Ethical approval

The University of South Wales Ethics Committee (LSE1636GREO), and Cardiff and Vale University Health Board (15/AIC/6352) approved the study. All procedures were carried out in accordance with the Declaration of Helsinki of the World Medical Association.¹⁶ Written informed consent was obtained from participants in study arm 1. Study arm 2 constituted a retrospective analysis of an anonymized database and thus patient consent was waived.

Design

We conducted a two-armed study. First, in order to determine the CDs of selected CPET variables (reported as independent predictors of post-operative outcome), analytical variation was calculated, and biological variation derived using repeated CPET results from a young apparently healthy population (Arm 1). Subsequently, these CD values were retrospectively applied to an anonymised database of patients who had CPET prior to colorectal surgery, in order to re-appraise **fitness** stratification (Arm 2).

Study arm 1: CD determination

Analytical variation (CV_A); the first component of CD, was determined by repeatedly passing inspired and expired gases through a Medgraphics Ultima metabolic cart (MedGraphicsTM, Gloucester, UK) in a manner that replicated typical ventilatory responses during the latter stages of a patient CPET (ie. pulmonary minute ventilation of $25 \text{ L}\cdot\text{min}^{-1}$). In a series of eight repeated trials each lasting ten respiratory cycles, a 250 L Douglas bag containing saturated expired gas (17% O_2 , 5% CO_2) and an equivalent volume of ambient gas was passed through a pneumotach and gas analyser. Inspiration and expiration were simulated using two-way non-rebreathing valves (2700 Series) connected to two factory calibrated 3 L syringes (Hans

1
2
3 Rudolph, Kansas City, USA) operated simultaneously (Figure 1.). Prior to sampling,
4
5 calibration was undertaken in accordance with manufacturer's guidelines using a 3 L syringe
6
7 and a known precision gas. During data collection the middle five of seven breaths were
8
9 averaged.

10
11 The within participant coefficient of variation (CV_w) from which biological variation
12
13 could be calculated, was determined by completion of three repeat CPETs separated by a
14
15 minimum of 24 hours, for 12 healthy participants (Table 1). Tests were conducted in a
16
17 randomised order at three time points across operating hours for patient CPET clinics (09:00
18
19 to 10:30, 12:00 to 13:30, and 15:00 to 17:00). All CPETs were conducted to volitional fatigue
20
21 using the Wasserman protocol,¹⁷ the same metabolic cart and investigator, and calibration
22
23 undertaken as previously described. Following three minutes of resting data collection,
24
25 participants cycled at 60 revolutions per minute on an electromagnetically braked cycle
26
27 ergometer (Lode, Gronigen, The Netherlands) for three minutes in an unloaded
28
29 "freewheeling" state. A progressively ramped period of exercise (10 to 30 W min^{-1} based on
30
31 stature, age, and predicted $\dot{V}O_2$)¹⁷ was then undertaken to volitional termination and followed
32
33 by three minutes recovery. Heart rate (Polar electro, Oy, Finland) was recorded throughout.
34
35
36

37 Medgraphics BreezeTM software automatically determined $\dot{V}O_2$ peak (defined as the
38
39 highest $\dot{V}O_2$ during the final 30 seconds of exercise reported), oxygen uptake efficiency slope
40
41 (OUES), and peak oxygen pulse (O_2 pulse). The AT was manually interpreted by a clinician
42
43 using the V-slope method,¹⁸ supported by $\dot{V}_E/\dot{V}CO_2$ -AT, and $\dot{V}_E/\dot{V}O_2$ -AT.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Critical Difference

Natural variation is described by the magnitude of CD and determines the difference in CRF required to demonstrate change not simply due to the “noise” associated with analytical imprecision (represented by CV_A) and biological variation (represented by CV_B), in order for it to be considered clinically meaningful.^{13 14} Critical difference uses ANOVA to determine the magnitude of random fluctuation around a homeostatic set point within which there is 95% probability that repeated measures will fall. The 95% probability is represented by a constant k (2.77) in Equation 1 (calculated from $\sqrt{2} * 1.96$ (two standard deviations)). Coefficients of variation were calculated dividing the standard deviation by the mean score and converted into a percentage as shown in the example of CV_A (Equation 2). The coefficient of analytical variation was subtracted from the CV_W determined from the repeated trials to calculate CV_B (Equation 3).

$$CD = k \sqrt{CV_A^2 + CV_B^2} \quad (\text{Eq 1})^{13}$$

Where:

k = constant equal to 2.77 at $P < 0.05$

CV_A = coefficient of analytical variation

CV_B = coefficient of biological variation

CV_A was calculated using the following equation:

$$CV_A = \frac{SD}{\bar{x}} \times 100 (\%) \quad (\text{Eq 2})$$

Where:

SD = standard deviation

\bar{x} = mean

CV_B was calculated from $\dot{V}O_2$ data from each participant, collected at periodic times as described, using the following equation:

$$CV_B = CV_W(\%) - CV_A(\%) \quad (\text{Eq 3})$$

Where:

CV_W = coefficient of within participant variation

Consequently, when interpreting CPET results, and in order to address the presence of natural variation, the CD (applied above and below an observed score) must be considered to determine the range in which a patient can present without any change in CRF (ie. before clinical significance can be claimed).

Study arm 2: application of CD metrics to patients

A consecutive sample of 213 patients (Table 1) scheduled for elective colorectal surgery who had undergone CPET testing was retrospectively examined. CPETs were conducted in accordance with the American Thoracic Society/ American College of Chest Physician Statement on Cardiopulmonary Exercise Testing,¹⁹ using identical equipment, investigators, and protocols as outlined in Study arm 1.

Calculated CD metrics were subsequently applied to CPET metrics with established evidence to independently identify **unfit** patients during pre-surgical assessment.^{1-4 6 11 20 21} Reference CRF threshold values were established from the European Association for Cardiovascular Prevention & Rehabilitation (EACPR)/American Heart Association (AHA) Scientific Statement: $\dot{V}O_2\text{-AT} < 11\text{mL O}_2 \text{ kg}^{-1} \text{ min}^{-1}$, $\dot{V}O_2 \text{ peak} < 16\text{mL O}_2 \text{ kg}^{-1} \text{ min}^{-1}$, and $\dot{V}_E/\dot{V}CO_2\text{-AT} \geq 36$.²² The CD for additional CPET metrics was calculated for $\dot{V}_E/\dot{V}O_2\text{-AT}$,^{3 20} and peak O₂ pulse.^{5 7 10 24}

To determine the impact of natural variation on **fitness** stratification, patient counts were calculated for uncorrected (observed) **fit** and **unfit** categories according to EACPR/AHA

1
2
3 threshold values, positively corrected (+CD), and negatively corrected (-CD) values. A
4 revised **fitness** stratification model for each CPET metric was created by applying \pm CD to
5 threshold values, thus creating upper and lower boundaries associated with natural variation,
6
7 and the area in-between the newly defined boundaries classified as indeterminate-**fitness**.
8
9
10
11 Finally, patient counts were compared for current versus newly revised models.
12
13
14

15 **Statistics**

16
17 Statistical analyses were conducted using IBM SPSS Statistics for Windows (Version 23.0
18 Armonk, NY). Distribution normality was confirmed using Shapiro-Wilk W tests. Within-
19 subject time of day difference in CPET performance was assessed using Bonferroni corrected
20 repeated measures analysis of variance. Patient counts were analysed using Chi-Square tests.
21
22 Continuous data are presented as mean (standard deviation) or median (range), and
23 categorical data as absolute values (%). Significance for all two-tailed tests was established at
24
25 $P < 0.05$. Retrospective sample size calculations were conducted attaining 80% power at the P
26
27 < 0.05 level with the minimum effect of clinical importance represented by the calculated CD
28
29 (from study arm 1, Table 2) and between-patient standard deviations (from study arm 2,
30
31 Table 1).²⁵
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Results

Natural variation

Study arm 1 identified a CD of 19% for $\dot{V}O_2$ -AT (CV_A 2.2%, CV_B 6.5%), 13% for $\dot{V}O_2$ peak (CV_A 2.2%, CV_B 3.9%), and 10% for $\dot{V}_E/\dot{V}CO_2$ -AT (CV_A 0.6%, CV_B 3.6%) (Table 2.). The time of day that CPET was conducted had no effect in measured metrics ($\dot{V}O_2$ -AT: $P = 0.40$, $\dot{V}O_2$ peak: $P = 0.81$, and $\dot{V}_E/\dot{V}CO_2$ -AT: $P = 0.75$). When CD was applied to current CPET **fitness** threshold values of $\dot{V}O_2$ -AT: 11mL O₂ kg⁻¹ min⁻¹, $\dot{V}O_2$ peak: 16mL kg⁻¹ min⁻¹, and $\dot{V}_E/\dot{V}CO_2$ -AT: 36, a variation of ± 2.1 mL O₂ kg⁻¹ min⁻¹, ± 2.0 mL kg⁻¹ min⁻¹, and ± 3.7 respectively was observed.

Potential for incorrect **fitness** stratification

We applied CD to positively and negatively correct (the range of) patient CPET scores around their observed (single-point estimate) scores, and subsequently calculated the number of “false positive” and “false negative” results. While these terms are not technically correct given the unavoidable uncertainty associated with biological variation and corresponding inability to determine an individual’s “true” level of CRF at any given point in time, it nonetheless provides a conceptual framework to illustrate how blunt application of current thresholds has the potential to affect perioperative planning for a large proportion of patients undergoing major elective surgery.

The application of natural variation (\pm CD) presented a mathematical possibility for patient results to transcend current **fitness** stratification boundaries thus demonstrating potential for misclassification (Figure 2) using $\dot{V}O_2$ -AT, $\dot{V}O_2$ peak, and $\dot{V}_E/\dot{V}CO_2$ -AT ($P < 0.001$ in all cases). Differences in patient counts assigned to a given **fitness** category resulted in false negatives (whereby patients were stratified as **fit** with variation positively corrected when they were originally **unfit**), and false positives (whereby patients were stratified as **unfit**

1
2
3 with variation negatively corrected when they were originally **fit**). Thus, natural variation
4 may have caused up to 59 (28%) false negatives and 69 (32%) false positives at the AT, 33
5 (15%) false negatives and 35 (16%) false positives at peak $\dot{V}O_2$, and 37 (17%) false negatives
6 and 43 (20%) false positives at the $\dot{V}_E/\dot{V}CO_2$ -AT.
7
8
9
10

11 12 13 **Revised model**

14
15 A revised **fitness** stratification model (Figure 3) was created with CD defining asymmetrical
16 upper and lower boundaries for absolute values (13.6 and 9.2mL O₂ kg⁻¹ min⁻¹ for AT, 18.3
17 and 14.2mL kg⁻¹ min⁻¹ for $\dot{V}O_2$ peak, 40.1 and 32.7 for $\dot{V}_E/\dot{V}CO_2$ -AT) that were independent
18 of **fitness** misclassification based on natural variation. The resultant area between the upper
19 and lower boundaries represented a newly defined and additional category labelled
20 “Indeterminate-**fitness**”. The indeterminate-**fitness** category accounted for 60, 32, and 40% of
21 patients for the AT, $\dot{V}O_2$ peak and $\dot{V}_E/\dot{V}CO_2$ -AT metrics respectively (Figure 4), and thus
22 fewer patients were stratified **as unfit or fit**.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Discussion

The present findings highlight the potential for incorrect patient **fitness** stratification when natural variation is not taken into account. We formulated a revised model (accounting for natural variation) which established that many patients were stratified with indeterminate-**fitness**. We therefore encourage clinicians to be aware of natural variation and its implications for **fitness** stratification and suggest this concept be applied to markers of CRF to further optimise patient management. Whilst this investigation aims to improve the prognostic interpretation of CPET results, we acknowledge and advocate that clinical decision making does not rely on the application of threshold values alone. There are clear dangers of just using a single point estimate, even if it may be a better number when natural variation is considered. A multitude of additional variables such as work rate, heart rate, duration of exercise, reason for stopping the exercise all go into a composite estimate of functional capacity to be considered alongside other clinical measures when planning perioperative care.

Potential for incorrect patient **fitness** stratification

The mean CPET score for patients undergoing colorectal surgery was identical to the threshold marker value for AT, within $0.3\text{mL O}_2 \text{ kg}^{-1} \text{ min}^{-1}$ for $\dot{V}\text{O}_2$ peak, and 2.4 lower for $\dot{V}_E/\dot{V}\text{CO}_2\text{-AT}$. Thus, when patient scores were positively or negatively corrected with CD, large numbers of patients transcended the EACPR/AHA threshold CRF boundaries demonstrating that natural variation may cause significant rates of incorrect **fitness** stratification. Of the three primary CPET metrics reported, the AT demonstrated the most incorrectly stratified patients, closely followed by peak $\dot{V}\text{O}_2$, and to a lesser albeit significant extent $\dot{V}_E/\dot{V}\text{CO}_2\text{-AT}$ in line with magnitudes of reported CD values and close proximity of patient scores to threshold boundaries. Furthermore, a valid and reliable identification of

1
2
3 $\dot{V}O_2$ -AT is not always possible and has been well documented in patients with heart failure,²⁶
4 and thus may contribute to greater variance in AT.
5
6
7
8

9 **Revised fitness stratification**

10
11 Our revised model (with its wider boundaries accounting for natural variation) excluded
12 many patients from both **unfit** and **fit** categories, and thus large numbers were stratified in the
13 indeterminate-**fitness** category (Figure 4). Not only does this occurrence confirm the impact
14 of natural variation, but consequently presents the challenge of planning perioperative care
15 for patients within this additional **fitness** category. Concerns may be associated with the
16 introduction of an additional **fitness** category. For example, patients undergoing colorectal
17 surgery who fell into an intermediate-**fitness** group (albeit not comparable with our
18 indeterminate-**fitness** category) have reported a higher rate of serious complications if
19 admitted to the ward rather than HDU.²⁷
20
21
22
23
24
25
26
27
28
29
30

31 The most effective way to assess patient risk is likely a combined approach using
32 clinical variables, biomarkers of susceptibility to disease, and physiological testing (CPET).²⁸
33 We suggest further development of our model by inclusion of known risk factors independent
34 of CRF to optimise perioperative care.
35
36
37
38
39
40
41

42 **Limitations**

43
44 We recognise that this study has limitations and simply reflects a “proof of principle”
45 concept. Measures of CD were derived from young healthy participants and applied to a
46 cohort of older patients. Comparative values for older controls were not available and would
47 present considerable ethical challenges to determine given that repeat CPET to volitional
48 exhaustion would be required. Our CV_W (given by $CV_A + CV_B$ from Table 2) of 6.1% for
49 $\dot{V}O_2$ peak is comparable with chronic obstructive pulmonary disease (6.6%) and congestive
50
51
52
53
54
55
56
57
58
59
60

1
2
3 heart failure patients (5.7% and 6.0%).²⁹⁻³¹ Furthermore, our CV_W for AT (8.7%) is consistent
4
5 with patient data (6.8%, 9.2% and 10%),^{32 30 33} and in excess of CV_W values for $\dot{V}O_2$ peak, the
6
7 probable consequence of observer error when determining AT via the V-slope method.¹⁸
8
9 Thus, our method has potential application to clinical populations. However, reported metrics
10
11 for CD may reflect a best-case scenario (ie. lowest CD) if natural variation increases with age
12
13 and/or pathology.
14

15
16 Study arm 1 comprised of men only, whilst the calculated CD was subsequently
17
18 applied to a population of whom 41% were women. For the $\dot{V}O_2$ peak and $\dot{V}O_2$ -AT metrics,
19
20 our coefficients of variation were comparable with the studies previously stated which also
21
22 included female data. Metrics represented by ventilatory equivalents however must be treated
23
24 with caution (for female comparison) as any disparity between the sexes is not accounted for.
25

26
27 Many CPET metrics are scaled to body mass. Further investigation is required to
28
29 determine if there are any effects on the magnitude of asymmetry for absolute values reported
30
31 around our zones of indeterminate-fitness resulting from scaling to body mass.
32

33
34 Data were collected on a single system in both arms of this study. We are aware that
35
36 analytical precision is likely to vary widely between different manufacturers thus affecting
37
38 CV_A and consequently CD. Therefore, our results can only be applied with certainty to
39
40 clinical tests using Medgraphics equipment. At the time of conducting the study the authors
41
42 did not have access to a metabolic calibrator used to calculate CV_A however we are confident
43
44 that our findings (up to 2.2%) are comparable with data produced from such devices which
45
46 typically report with accuracy of $\pm 2\%$.³⁴
47
48
49

50 **Prospective sample size calculations**

51
52 From an experimental design perspective, our observations have implications when
53
54 prospectively determining sample sizes for future randomised controlled exercise trials. We
55
56
57
58
59
60

1
2
3 suggest that CD be used to determine the minimal clinically important difference (MCID) for
4 any given metric of CRF. Until now, studies often rely on MCID values that appear to lack a
5 well-established scientific basis, such as a $\dot{V}O_2$ -AT of $2\text{mL kg}^{-1} \text{min}^{-1}$ for example.³⁵ This
6 (arbitrarily) defined MCID of $2\text{mL O}_2 \text{kg}^{-1} \text{min}^{-1}$ is in fact incorrect because it falls within our
7 calculated CD of $2.1\text{mL O}_2 \text{kg}^{-1} \text{min}^{-1}$ (i.e. this is part of normal variation). In a worked
8 example using the arbitrary metric of $2\text{mL O}_2 \text{kg}^{-1} \text{min}^{-1}$, a prospective power calculation
9 indicates that a two-armed exercise intervention study would require a minimum of 36
10 patients per group (excluding potential dropout) to detect a treatment effect with 80% power
11 at the $P < 0.05$ level. However, considering natural variation (using our calculated CD of
12 $2.1\text{mL O}_2 \text{kg}^{-1} \text{min}^{-1}$ in place of $2\text{mL O}_2 \text{kg}^{-1} \text{min}^{-1}$) would further inflate the sample size (to
13 39 patients per group) highlighting the potential for a type II error.

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29 We recognise that the sample size calculation is based upon a CD determined from a sample
30 of 12 subjects and is limited to a single (Medgraphics) system. Further research (with larger
31 sample sizes, additional metabolic carts, and calculations across the spectrum of age, health
32 and CRF) is encouraged to better support our prospective calculation of sample sizes.

33 34 35 36 37 38 39 **Conclusions**

40
41 These findings demonstrate the extent of natural variation in CPET data. Natural variation
42 also has potential to influence patient fitness stratification. Therefore, clinicians should not
43 consider fitness as a single point estimate, but instead as a dynamic range of values defined
44 by natural variation and calculated using critical difference. We suggest the use of CRF
45 threshold values inclusive of natural variation to optimise risk prediction models, and
46 encourage clinicians to be aware of natural variation and its implications when determining
47 the appropriate level of post-operative care following major surgery.

48
49
50
51
52
53
54
55
56
57
58
59
60

Author's contributions

All authors were involved in the conception and design of study. R.G.D, I.R.A, G.A.R performed the CPET tests and collated the data. G.A.R. performed the analysis with input from D.M.B, M.H.L, R.G.D, I.R.A. The manuscript was drafted by G.A.R and D.M.B. All authors provided revisions and approved the final version for submission.

Declaration of interest

The authors declare no conflict of interest.

Funding

This work was supported by the Higher Education Funding Council for Wales (to DM Bailey). DM Bailey is a Royal Society Wolfson Research Fellow (#WM170007).

Acknowledgements

Sean Cutler, James O'Flaherty, and Trevor Harris contributed to data collection and technical input in study arm 1.

References

- 1 West MA, Asher R, Browning M, et al. Validation of preoperative cardiopulmonary exercise testing-derived variables to predict in-hospital morbidity after major colorectal surgery. *Br J Surg* 2016; **103**: 744-52
- 2 Grant SW, Hickey GL, Wisely NA, et al. Cardiopulmonary exercise testing and survival after elective abdominal aortic aneurysm repair. *Br J Anaesth* 2015; **114**: 430-6
- 3 Carlisle J, Swart M. Mid-term survival after abdominal aortic aneurysm surgery predicted by cardiopulmonary exercise testing. *Br J Surg* 2007; **94**: 966-9
- 4 Lai CW, Minto G, Challand CP, et al. Patients' inability to perform a preoperative cardiopulmonary exercise test or demonstrate an anaerobic threshold is associated with inferior outcomes after major colorectal surgery. *Br J Anaesth* 2013; **111**: 607-11
- 5 Prentis JM, Trenell MI, Jones DJ, Lees T, Clarke M, Snowden CP. Submaximal exercise testing predicts perioperative hospitalization after aortic aneurysm repair. *J Vasc Surg* 2012; **56**: 1564-70
- 6 Snowden CP, Prentis J, Jacques B, et al. Cardiorespiratory fitness predicts mortality and hospital length of stay after major elective surgery in older people. *Annals of surgery* 2013; **257**: 999-1004
- 7 West MA, Parry MG, Lythgoe D, et al. Cardiopulmonary exercise testing for the prediction of morbidity risk after rectal cancer surgery. *Br J Surg* 2014; **101**: 1166-72
- 8 Ross R, Blair SN, Arena R, et al. Importance of Assessing Cardiorespiratory Fitness in Clinical Practice: A Case for Fitness as a Clinical Vital Sign: A Scientific Statement From the American Heart Association. *Circulation* 2016; **134**: e653-99
- 9 Older R, Smith R, Courtney B, Hone R. Preoperative Evaluation of Cardiac Failure and Ischemia in Elderly Patients by Cardiopulmonary Exercise Testing. *Chest* 1993; **104**: 701-4

- 1
2
3 10 Junejo MA, Mason JM, Sheen AJ, et al. Cardiopulmonary exercise testing for
4 preoperative risk assessment before hepatic resection. *Br J Surg* 2012; **99**: 1097-104
5
6
7 11 Hartley RA, Pichel AC, Grant SW, et al. Preoperative cardiopulmonary exercise testing
8 and risk of early mortality following abdominal aortic aneurysm repair. *Br J Surg* 2012; **99**:
9
10 1539-46
11
12 12 Moran J, Wilson F, Guinan E, McCormick P, Hussey J, Moriarty J. Role of
13 cardiopulmonary exercise testing as a risk-assessment method in patients undergoing intra-
14 abdominal surgery: a systematic review. *Br J Anaesth* 2016; **116**: 177-91
15
16
17 13 Fraser CG, Fogarty Y. Interpreting laboratory results. *BMJ* 1989; **298**: 1659-60
18
19
20 14 Davison GW, Ashton T, McEneny J, Young IS, Davies B, Bailey DM. Critical difference
21 applied to exercise-induced oxidative stress: the dilemma of distinguishing biological from
22 statistical change. *J Physiol Biochem* 2012; **68**: 377-84
23
24
25 15 Bailey DM, Evans TG, Gower Thomas K. Intervisceral artery origins in patients with
26 abdominal aortic aneurysmal disease; evidence for systemic vascular remodelling. *Exp*
27 *Physiol* 2016; **101**: 1143-53
28
29
30
31
32
33
34 16 Williams JR. The Declaration of Helsinki and public health. *Bull World Health Organ*
35 2008; **86**: 650-2
36
37
38 17 Wasserman K. *Principles of exercise testing and interpretation: including*
39 *pathophysiology and clinical applications*. 5th ed Edn. London: Wolters Kluwer/Lippincott
40 Williams & Wilkins, 2012
41
42
43
44
45
46 18 Beaver WL, Wasserman K, Whipp BJ. A new method for detecting anaerobic threshold
47 by gas exchange. *J Appl Physiol* 1986; **60**: 2020-7
48
49
50 19 American Thoracic S, American College of Chest P. ATS/ACCP Statement on
51 cardiopulmonary exercise testing. *Am J Respir Crit Care Med* 2003; **167**: 211-77
52
53
54
55
56
57
58
59
60

1
2
3 20 West MA, Loughney L, Barben CP, et al. The effects of neoadjuvant chemoradiotherapy
4 on physical fitness and morbidity in rectal cancer surgery patients. *Eur J Surg Oncol* 2014;
5 **40**: 1421-8
6
7

8
9 21 Wilson RJ, Davies S, Yates D, Redman J, Stone M. Impaired functional capacity is
10 associated with all-cause mortality after major elective intra-abdominal surgery. *Br J Anaesth*
11 2010; **105**: 297-303
12
13

14
15 22 Guazzi M, Arena R, Halle M, Piepoli MF, Myers J, Lavie CJ. 2016 Focused Update:
16 Clinical Recommendations for Cardiopulmonary Exercise Testing Data Assessment in
17 Specific Patient Populations. *Circulation* 2016; **133**: e694-711
18
19

20
21 23 Junejo MA, Mason JM, Sheen AJ, et al. Cardiopulmonary exercise testing for
22 preoperative risk assessment before pancreaticoduodenectomy for cancer. *Ann surg oncol*
23 2014; **21**: 1929-36
24
25
26
27

28
29 24 Epstein SK, Freeman RB, Khayat A, Unterborn JN, Pratt DS, Kaplan MM. Aerobic
30 capacity is associated with 100-day outcome after hepatic transplantation. *Liver transpl* 2004;
31 **10**: 418-24
32
33
34

35
36 25 Altman DG. Statistics and ethics in medical research: III How large a sample? *BMJ* 1980;
37 **281**: 1336-8
38

39
40 26 Arena R, Myers J, Williams MA, et al. Assessment of functional capacity in clinical and
41 research settings: a scientific statement from the American Heart Association Committee on
42 Exercise, Rehabilitation, and Prevention of the Council on Clinical Cardiology and the
43 Council on Cardiovascular Nursing. *Circulation* 2007; **116**: 329-43
44
45
46
47

48
49 27 Swart M, Carlisle JB, Goddard J. Using predicted 30 day mortality to plan postoperative
50 colorectal surgery care: a cohort study. *Br J Anaesth* 2017; **118**: 100-4
51

52
53 28 Grocott MPW. Improving outcomes after surgery. *BMJ* 2009; **339**: b5173
54
55
56
57
58
59
60

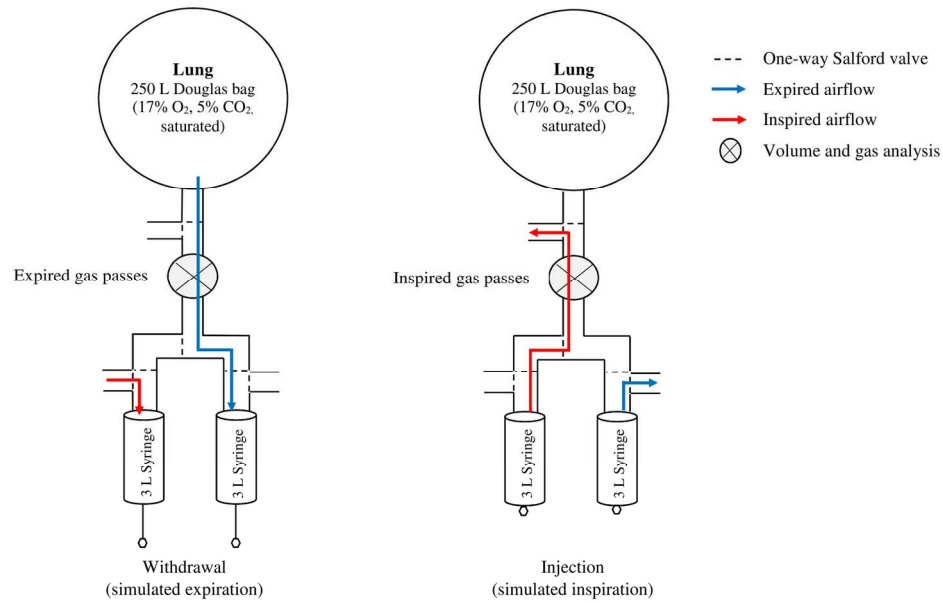
- 1
2
3 29 Owens MW, Kinasewitz GT, Strain DS. Evaluating the Effects of Chronic Therapy in
4 Patients with Irreversible Air-Flow Obstruction. *Am Rev Respir Dis* 1986; **134**: 935-7
5
6
7 30 Janicki JS, Gupta S, Ferris ST, McElroy PA. Long-term Reproducibility of Respiratory
8 Gas Exchange Measurements during Exercise in Patients with Stable Cardiac Failure. *Chest*
9 1990; **97**: 12-7
10
11
12
13 31 Elborn JS, Stanford CF, Nicholls DP. Reproducibility of cardiopulmonary parameters
14 during exercise in patients with chronic cardiac failure. The need for a preliminary test. *Eur*
15 *Heart J* 1990; **11**: 75-81
16
17
18
19 32 Keteyian SJ, Brawner CA, Ehrman JK, Ivanhoe R, Boehmer JP, Abraham WT.
20 Reproducibility of peak oxygen uptake and other cardiopulmonary exercise parameters:
21 implications for clinical trials and clinical practice. *Chest* 2010; **138**: 950-5
22
23
24
25 33 Kothmann E, Danjoux G, Owen SJ, Parry A, Turley AJ, Batterham AM. Reliability of the
26 anaerobic threshold in cardiopulmonary exercise testing of patients with abdominal aortic
27 aneurysms. *Anaesthesia* 2009; **64**: 9-13
28
29
30
31
32 34 Huszczuk A, Whipp B, Wasserman K. A respiratory gas exchange simulator for routine
33 calibration in metabolic studies. *Eur Respir J* 1990; **3**: 465-8
34
35
36
37 35 Kothmann E, Batterham AM, Owen SJ, et al. Effect of short-term exercise training on
38 aerobic fitness in patients with abdominal aortic aneurysms: a pilot study. *Br J Anaesth* 2009;
39 **103**: 505-10
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Participant and patient characteristics. Data are shown as mean (\pm standard deviation) or \sim (range), and *n (%). n, number; IHD, ischaemic heart disease; COPD, chronic obstructive pulmonary disease; $\dot{V}O_2$ peak, peak oxygen consumption; RER, respiratory exchange ratio; AT, estimated anaerobic threshold; $\dot{V}_E/\dot{V}CO_2$, ventilatory equivalent for carbon dioxide; $\dot{V}_E/\dot{V}O_2$, ventilatory equivalent for oxygen; O₂ pulse, oxygen pulse at peak exercise; Work load at AT, work load at estimated anaerobic threshold; Workload at peak, work load at peak exercise.

	<i>Study arm 1</i> Apparently healthy participants (n = 12)	<i>Study arm 2</i> Colorectal patients (n = 213)
Demographics:		
Age (years) [~]	22 (20-26)	69 (32-90)
BMI	26 (3.1)	28.3 (5.8)
Sex*		
male	12 (100)	126 (59)
female	0 (0)	87 (41)
Risk factors:		
Smoking*		
no	12 (100)	71 (33)
yes (active/former)	0 (0)	142 (67)
Hypertension*	0 (0)	79 (37)
Diabetes*	0 (0)	34 (16)
IHD*	0 (0)	37 (17)
COPD*	0 (0)	21 (10)
Haemoglobin (g L ⁻¹)	-	12.7 (1.9)
Creatinine (μ mol L ⁻¹)	-	79.2 (19.7)
Cardiopulmonary function:		
Baseline heart rate (beats min ⁻¹)	65 (5)	83 (19)
Peak heart rate (beats min ⁻¹)	178 (5)	124 (28)
$\dot{V}O_2$ peak (mL kg ⁻¹ min ⁻¹)	43.8 (6.0)	16.3 (4.9)
RER at peak $\dot{V}O_2$	1.3 (0.1)	1.1 (0.1)
AT (mL O ₂ kg ⁻¹ min ⁻¹)	23.8 (3.6)	11.0 (3.0)
$\dot{V}_E/\dot{V}CO_2$ -AT	23.5 (1.4)	33.6 (5.3)
$\dot{V}_E/\dot{V}O_2$ -AT	23.5 (4.7)	30.6 (5.9)
O ₂ pulse (mL beat ⁻¹)	20.7 (0.9)	10.5 (3.8)
Work load at AT (W)	160 (28)	52 (28)
Work load at peak (W)	300 (45)	91 (47)

Table 2. Biological variation and critical difference for cardiopulmonary exercise test variables (Study arm 1, n=12). CV_A , coefficient of analytical variation; CV_B , coefficient of biological variation; AT, anaerobic threshold; $\dot{V}O_2$ peak, peak oxygen consumption; $\dot{V}_E/\dot{V}CO_2$, ventilatory equivalent for carbon dioxide; $\dot{V}_E/\dot{V}O_2$, ventilatory equivalent for oxygen; O_2 pulse, oxygen pulse at peak exercise; OUES, oxygen uptake efficiency slope; RER, respiratory exchange ratio.

Parameter	CV_A (%)	CV_B (%)	Critical difference (%)
AT (mL O_2 kg^{-1} min^{-1})	2.2	6.5	19.1
$\dot{V}O_2$ peak (mL kg^{-1} min^{-1})	2.2	3.9	12.5
$\dot{V}_E/\dot{V}CO_2$ -AT	0.6	3.6	10.2
$\dot{V}_E/\dot{V}O_2$ -AT	1.7	3.0	9.6
O_2 pulse (mL $beat^{-1}$)	2.2	2.3	8.9
OUES	2.2	3.8	12.1
RER at peak exercise	1.4	5.3	15.2



28 **Figure 1. The determination of CV_A for CPET metrics using simulated expiration and**
 29 **inspiration.** CV_A , analytical coefficient of variation; CPET, cardiopulmonary exercise test. Simulated oxygen
 30 uptake for trials $\sim 13\text{ mL kg}^{-1} \text{ min}^{-1}$.

31 153x98mm (300 x 300 DPI)

32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

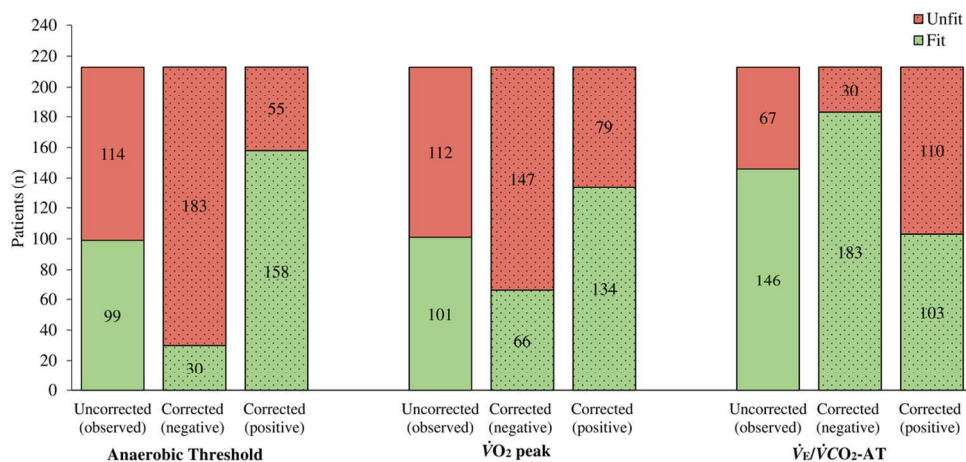


Figure 2. Potential for incorrect patient fitness stratification if natural variation is not taken into account. Patient counts are presented for unfit ($AT < 11\text{ mL O}_2 \text{ kg}^{-1} \text{ min}^{-1}$, $\dot{V}O_2 \text{ peak} < 16\text{ mL kg}^{-1} \text{ min}^{-1}$, $\dot{V}_E/\dot{V}CO_2 \geq 36$) and fit ($AT \geq 11\text{ mL O}_2 \text{ kg}^{-1} \text{ min}^{-1}$, $\dot{V}O_2 \text{ peak} \geq 16\text{ mL kg}^{-1} \text{ min}^{-1}$, $\dot{V}_E/\dot{V}CO_2 < 36$) categories. AT, anaerobic threshold; $\dot{V}O_2$ peak, peak oxygen consumption; $\dot{V}_E/\dot{V}CO_2$, ventilatory equivalent for carbon dioxide; Observed, uncorrected scores indicative of current risk stratification; Positive, corrected scores by addition of CD; Negative, corrected scores by subtraction of CD. $P < 0.001$ across all pairwise comparisons for corrected scores. Natural variation caused 59 (28%) false negatives and 69 (32%) false positives at the AT, 33 (15%) false negatives and 35 (16%) false positives at $\dot{V}O_2$ peak, and 37 (17%) false negatives and 43 (20%) false positives at the $\dot{V}_E/\dot{V}CO_2$ -AT.

114x53mm (300 x 300 DPI)

review

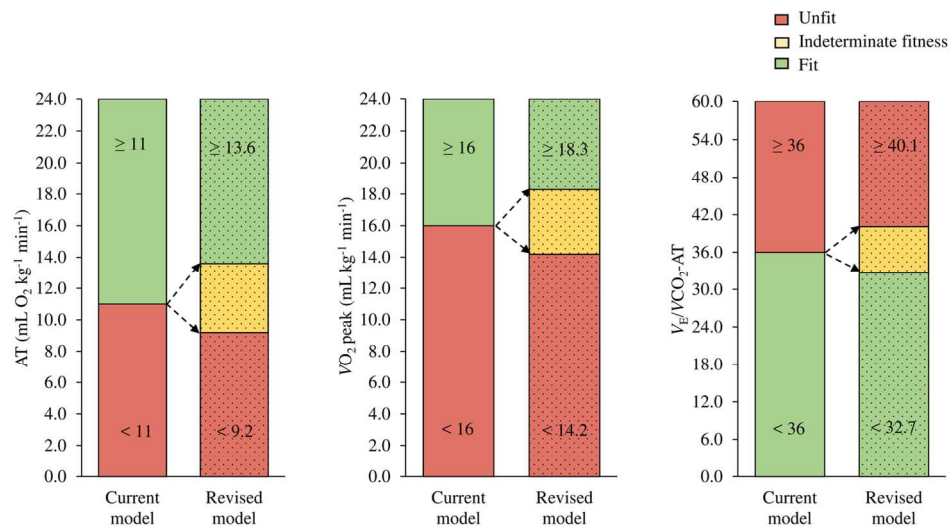


Figure 3. Revised fitness stratification model following incorporation of the critical difference for the anaerobic threshold, VO_2 peak, and V_E/VCO_2-AT . AT, anaerobic threshold; VO_2 peak, peak oxygen consumption; V_E/VCO_2 , ventilatory equivalent for carbon dioxide. Natural variation demonstrates the magnitude of variation present. The lower and upper boundaries define clinically meaningful boundaries not affected by natural variation whilst the area in-between is classified as indeterminate fitness.

131x72mm (300 x 300 DPI)

Review

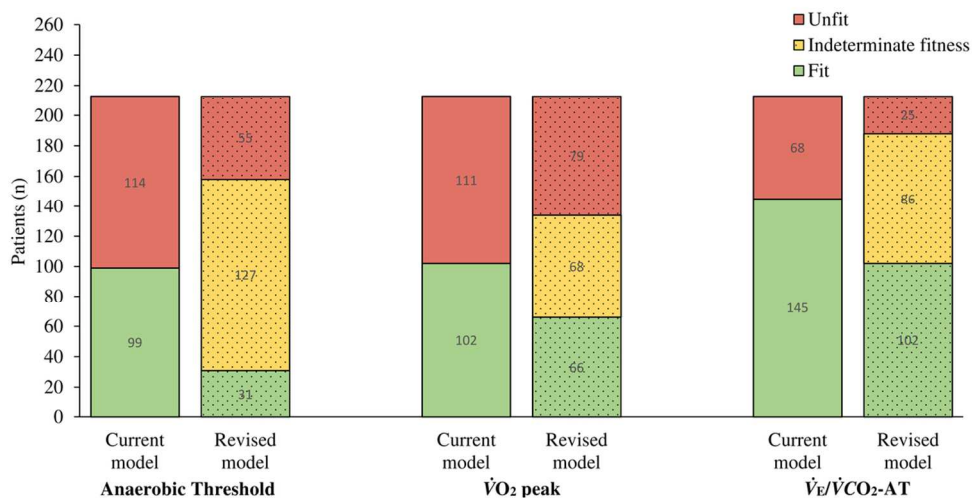


Figure 4. Current versus revised model identification of patient counts by fitness category. AT, anaerobic threshold; $\dot{V}O_2$ peak, peak oxygen consumption; $\dot{V}_E/\dot{V}CO_2$, ventilatory equivalent for carbon dioxide. The revised model demonstrates large numbers of patients that are classified with indeterminate fitness.

111x56mm (300 x 300 DPI)