

Trustworthy performance evaluations: The Performance Outcome Scoring Template (POS-T) for transparent assessments in real-world programs

Benjamin P. Raysmith (✉ ben.raysmith@athletics.org.au)

Linköping University

Toomas Timpka

Linköping University

Jenny Jacobsson

Linköping University

Michael K. Drew

Australian Institute of Sport

Örjan Dahlström

Linköping University

Research Article

Keywords: POS-T, outputs and outcomes, inputs and activities, self-comparison, singular allocation

Posted Date: September 17th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-899494/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Title**

2 Trustworthy performance evaluations: The Performance Outcome Scoring Template (POS-T)
3 for transparent assessments in real-world programs.

4

5 **Short Title**

6 Trustworthy performance evaluation

7

8 **Authors**

9 Benjamin P. Raysmith^{1, 2, 3*}, Toomas Timpka¹, Jenny Jacobsson^{1, 4}, Michael K. Drew^{5, 6,}
10 ⁷, Örjan Dahlström^{1,8}

11

12 **Affiliations**

13 ¹ Athletics Research Centre, Linköping University, Linköping, Sweden

14 ² Western Australian Institute of Sport, Perth, Australia

15 ³ Athletics Australia, Melbourne, Australia

16 ⁴ Swedish Athletics Association, Stockholm, Sweden

17 ⁵ Athlete Availability Program, Australian Institute of Sport, Bruce, Australian Capital
18 Territory, Australia

19 ⁶ Australian Collaboration for Research into Injury in Sport and its Prevention
20 (ACRISP), Australia

21 ⁷ University of Canberra Research Institute for Sport and Exercise (UCRISE), Canberra,
22 Australian Capital Territory, Australia

23 ⁸ Department of Behavioural Sciences and Learning, Linköping University, Linköping,
24 Sweden

25

26 *Corresponding Author: ben.raysmith@athletics.org.au

27

28 **Abstract**

29 In applied program settings, such as in natural environment control and education,
30 performance evaluation is usually conducted by evaluators considering both self-
31 comparison and comparison with peers. We have developed the Performance
32 Outcome Scoring Template (POS-T) for assessments with high face-validity in these
33 settings. POS-T puts achievements of individuals or groups in context, i.e. the
34 resulting performance outcome score (POS) reflects a meaningful measure of
35 performance magnitude with regards to internal and external comparisons.

36 Development of a POS is performed in four steps supported by a statistical
37 framework. Software is supplied for creation of scoring applications in different
38 performance evaluation settings. We demonstrate the POS-T by evaluation of CO₂
39 emissions reduction amongst 36 OECD member countries.

40

41

42

43

44

45

46

47 **MAIN TEXT**

48 **Introduction**

49 Performance evaluation seeks to examine the achievement of predetermined objectives or
50 goals by individuals or groups through the broad assessment of processes (inputs and
51 activities) and results (outputs and outcomes).¹⁻³ These evaluations are used to assess
52 program efficiency and effectiveness and provide accountability to resource allocation,
53 strategy and policy direction.^{3,4} Performance evaluations are usually case-specific and
54 defined by the stakeholders with the authority and responsibility to do so.^{5,6} Furthermore,
55 they are expected to provide a contextual judgement of performance at a moment in time
56 and require the measurement of credible ongoing outputs (performance measures) that
57 relate specifically to the needs of the evaluators.^{3,7-9} The assessment of goal achievement
58 lies with the process of performance evaluation and considers a broad array of factors that
59 include addressing the “How” and “Why” questions of achieving pre-determined objectives.
60 The complexity associated with performance evaluations has led to the development of
61 performance measurement systems to collect ongoing data and to monitor and report
62 progress towards pre-determined goals.^{6,7}

63 Performance evaluations addressed through stakeholder questions require that the
64 measures used to report performance results are relevant and trustworthy.⁹ A performance
65 ‘outcome’ is defined as the resultant effect of a system towards a pre-determined objective,
66 whereas a performance ‘output’ is the data generated by a single unique metric impacting
67 the outcome.^{3,10-12} A balanced performance evaluation includes both internal measurement
68 of the individual units in the program and of the external environment.^{13,14} These

69 measurements permit a reflection on achievement to date as well as what *could* be possible
70 to achieve within an equal context.¹⁵ For instance, internal measures of output include
71 ‘exam scores’ in an academic setting, ‘race finish time’ in a sport setting, ‘greenhouse gas
72 emissions per capita’ in an environment setting. Each of these metrics can be used as a
73 performance output measure that represents ‘self-comparison’ when collected over time.
74 However, internal measures alone may be insufficient in performance evaluation as the
75 appraisal of a performance output or trend can vary when assessed relatively against peer
76 performance under comparable circumstances. External environment measures of output
77 include ‘league tables’, ‘event final rankings’ or ‘comparisons with industry benchmarks’ and
78 have become a popular way to compare peers within industries.¹⁶⁻¹⁹ As with internal
79 measures, external measures of relative performance in isolation equally may lack face-
80 validity or attribution (singular allocation) in performance evaluation. This is seen in
81 circumstances where the appraisal of a ranking against peers can vary when interpreted in
82 the context of individual achievement or progress.²⁰⁻²²

83 Examples of bespoke performance measurement systems that comprise singular and
84 multiple output measures exist across industry and sectors: academic,²³ health,^{24,25} profit
85 and non-profit organisations,^{26,27} sport,^{28,29} natural environment,³⁰ government and private
86 sectors,³¹⁻³³ and finance.^{34,35} The use of a single output measure to reflect performance
87 outcome may provide a too narrow perspective on achievement in a complex program and
88 therefore risk evaluation face-validity in any of these areas.⁹ The selection of performance
89 output measures that enhance attribution and enable meaningful and trustworthy
90 evaluations therefore benefits from appropriate consideration.^{25,36,37} Different
91 measurements of performance can be combined into composite scores to increase the face-

92 validity of a program evaluation appraisal without adding difficulty to its interpretation.^{38,39}
93 Measurements on different data scales need here to be transformed into a common scale
94 before combining, even though the introduction of re-scaling may reduce reliability.⁴⁰ An
95 alignment of scales that reflect a meaningful magnitude change within and between
96 variables can improve face-validity in regards to program outcomes and thereby make a
97 composite score preferred for decision-making.^{38,41} It is therefore essential that the
98 measurements are transparent and comprehensible for all stakeholders involved to avoid
99 misleading program decisions using such composite scores.⁹ A robust metric must be
100 determined to assess objectively what has occurred and how this may influence future
101 outcomes relative to the investment in any program or individual and any third-party
102 interest.

103 The aim with this study was to develop a performance scoring template that combines
104 internal and external measures by alignment of scales that reflect meaningful magnitudes of
105 change in stakeholder defined contexts. The purpose is to provide evaluators of applied
106 programs a means to report performance outcomes with convincing face-validity. The
107 scoring template is exemplified by application to evaluation of CO₂ emissions reduction
108 amongst OECD member countries.

109

110 **Results**

111 Application of the Performance Outcome Scoring Template (POS-T) commences with
112 selection of data sources and concludes with a composite score that is adjusted for optimal
113 face-validity (POS) (Fig. 1). Subject to the evaluation purpose and selected time-point the
114 template can be utilised in comparing a result with a predetermined objective, benchmark

115 standard or appraise change over time. Statistical software (in the R language) is supplied to
116 support the development of a POS (Data S1). Stepwise instructions provided in the software
117 detail file data set-up for application in the code.

118

119 Two generic data sources, achievement (continuous scale) and rank (ordinal scale) are
120 handled in four cardinal and one optional (weighting) index development steps:

121

122 • **Components and Parameters** – Quantifiable domains deemed to have primacy with
123 respect to the face-validity of the performance outcome. Stakeholder selected
124 parameters that frame the evaluation.

125 • **Component Outputs** – Performance metrics collected from each component.
126 Representative of the data collection fields that comprise the performance output
127 measure and a comparative reference output within a distribution.

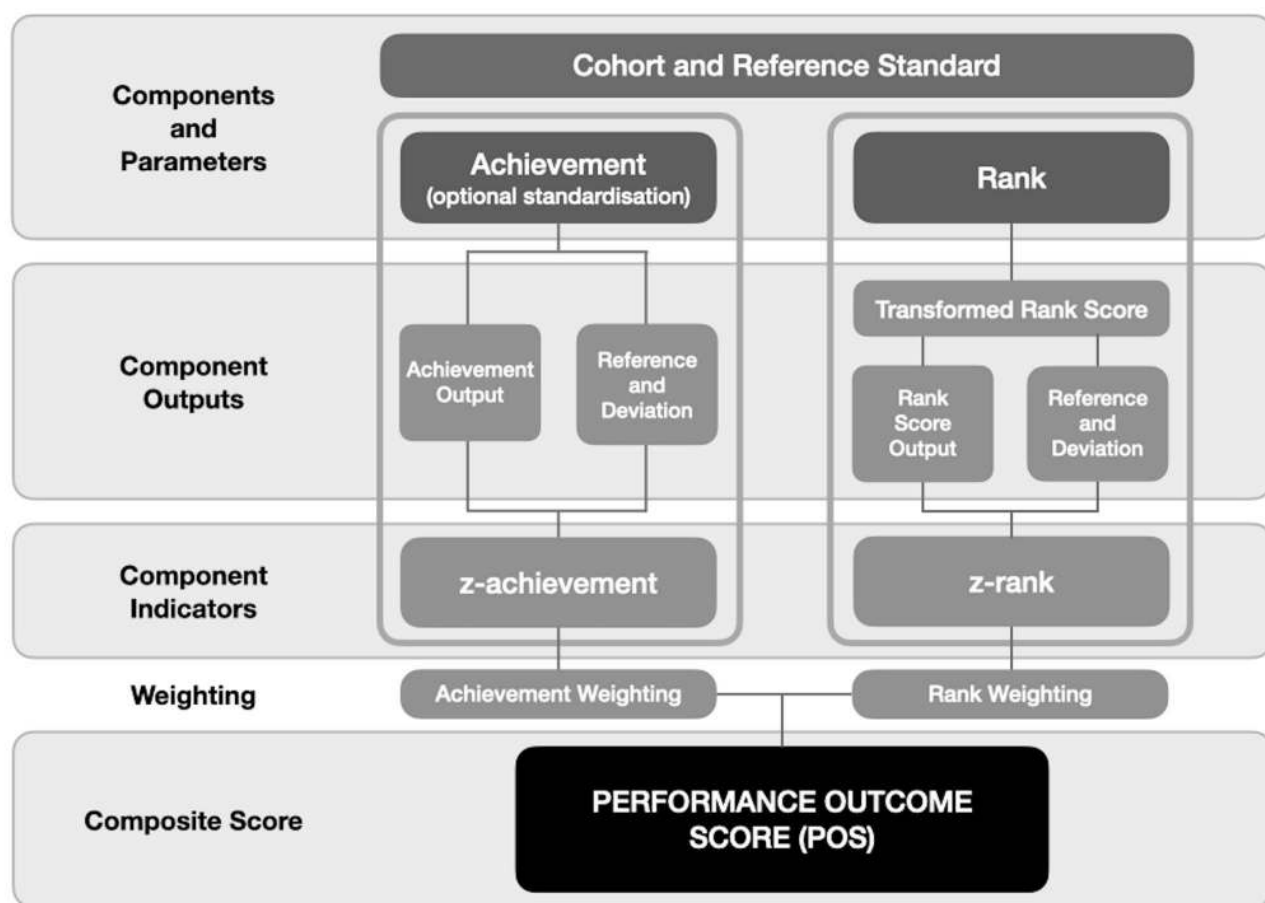
128 • **Component Indicators** – Component outputs transformed to normalised measures.
129 The performance outputs measured against a reference standard given assumptions
130 of both achievement and rank deviations.

131 • **Composite score** – A final viable measure is secured through aggregating and
132 optionally weighting normalised component indicators that meet the desired face-
133 validity.

134

135 Data handling through the POS-T is described in four detailed steps. Following stakeholder
136 selections made in step 1 the outputs from steps 2-4 are produced automatically when
137 applying the performance data to the statistical R-code software provided (Data S1).

138



139

140 Fig. 1. POS-T signifying four cardinal and one optional (weighting) data handling levels.

141

142 **Step 1 - Components (Quantifiable domains) and Parameters**

143 **Actions:** Define evaluation entity (individual, group or population participating in a specified
 144 program), sample of entities to evaluate, comparison cohort, and comparator (frame of
 145 reference for evaluation). Select quantifiable domains representing internal and external
 146 measures that provide face-validity for the performance outcome of an entity (individual or
 147 group). Optional selection of a standardising parameter.

148 **Outcome:** Defined entity(s), cohort, comparator, and quantifiable domains with arguments
 149 for their selection. Optional standardising parameter.

150 **Procedure:** The entities for evaluation are selected within a cohort of interest framed by the
151 context of the comparator. Examples of evaluation comparators include referencing a
152 previous time point to evaluate performance over time, referencing a population mean to
153 evaluate performance of the entity against a population standard, or referencing a single
154 measure like a season average or predetermined objective to evaluate the entity against
155 expectation. Next is selecting quantifiable domains that represent the **components** of
156 'achievement' (measured on a continuous scale) and 'rank' (measured on an ordinal scale).
157 For the entity being evaluated 'achievement' is characterised as the component denoting
158 self-comparison (internal measure), and 'rank' characterises the component denoting
159 comparison with others (external measure). The optional selection of a standardising
160 parameter is applied to the continuous data component. Standardising parameters are
161 measures of exposure applied to the continuous data component and examples include:
162 'per capita' calculations, standardisation by funding, access to resources, or other
163 parameters of exposure. Completion of step 1 is made by recording the arguments behind
164 selecting each of the components and describing the cohort parameters, reference
165 standards and optional standardising parameters.

166

167 **Step 2 - Component Outputs (performance metrics from each component)**

168 **Action:** Collect/calculate performance metrics from each component.

169 **Outcome:** For achievement and rank performance metrics refer to: Output, Reference, and
170 Deviation.

171 **Procedure:** Performance metrics are recorded for both achievement and rank. The
172 'achievement output' and 'rank score output' are established as well as descriptive data for
173 each component parameter, i.e. references and deviations (Table 1). This is completed for

174 each 'entity' (individuals or groups) in the cohort. When referring to separate entity's
175 subscripts '*i*' and '*j*' are used.

176

177 **Achievement performance metrics**

178 Output:

179 The '**achievement output**' (O_A) is quantified directly from the metric of interest crude
180 measure (continuous scale) and reflects the output *being evaluated*.

181 Reference:

182 The '**achievement reference**' (R_A) is the metric of interest crude measure from a previous
183 time point and is used as the metric of comparison as framed by the comparator defined in
184 step 1. The achievement reference and output are measured by the same metric on the
185 same continuous scale.

186 Deviation:

187 The '**achievement deviation**' (D_A) is the deviation of crude measures across the observation
188 period or other collection of entity crude measures that constitute a deviation around the
189 achievement reference. The reference and deviation values are set from the decision to use
190 a cohort pooled achievement standard deviation or individual entity standard deviations.

191 The achievement outputs are assumed to follow a normal distribution, $N(\bar{A}_i, \sigma_{A_i})$, from
192 which the reference and deviation values are set; $R_{A,i} = \bar{A}_i$ and $D_{A,i} = \sigma_{A_i}$.

193

194 **Rank performance metrics**

195 Rank scores (forming the ordinal ranking order) are transformed to a continuous scale value
196 (transformed rank-score) using a pre-defined function, f , to reflect non-equidistance
197 (magnitude difference) between different entities based on the continuous metric that

198 established the ranking order. Lower and upper ‘reference limiters’ are applied to establish
199 the transformed rank-score range. The reference limiters represent the range of minimum
200 and maximum reference scores and establishes a range for future equivalent comparisons.
201 Selection should consider a range beyond current reference score minima and maxima that
202 would account for a realistic range of future reference score possibilities. The lower
203 achievement limiter is attributed a transformed rank-score of 1 and the upper achievement
204 limiter is attributed a transformed rank-score of 100 (Fig. 2a). In different program settings
205 either a higher or lower achievement score may reflect the ‘best’ performance. This
206 directionality is established during the stakeholder selections at the start of the statistical R-
207 code data handling process.

208

209 Output

210 The ‘rank score output’ (O_{RS}) is formed by, for each i , aggregating the comparison of
211 transformed rank-scores between entity i , $f(\rho_i)$, and each entity j , $f(\rho_j)$, when entity i has
212 a better final rank (R_i) relative to entity j (R_j). This reflects meaningful magnitude
213 differences between the final ranks for entities i and j .

214

215 For entity i the ‘rank score output’ ($O_{RS,i}$) is then defined as:

216

$$O_{RS,i} = \sum_{\forall j \neq i} \delta(i, j) \quad (1)$$

217

218 where

219

$$\delta(i, j) = \begin{cases} \frac{f(\rho_j)}{f(\rho_i)}, & \text{if } R_i < R_j \\ 0, & \text{if } R_i \geq R_j \end{cases} \quad (2)$$

220

221 Reference and Deviation

222 The underlying performance metrics for ‘rank score outputs’ are generated using
 223 a **simulator**. The simulator generates an underlying distribution of rank score outputs based
 224 on the underlying distribution of possible achievement outputs. It randomly selects one
 225 achievement output for each entity, transforms them into ranks, transformed rank-scores,
 226 and finally rank score outputs. The outputs are saved for each entity i and the procedure is
 227 iteratively repeated, resulting in distributions of probable rank-score outputs, Φ_i , for each
 228 entity i (Fig. 2b).

229 For each entity i , the ‘rank score reference’, $R_{RS,i}$, and the ‘rank score deviation’, $D_{RS,i}$, are
 230 chosen as:

$$R_{RS,i} = Median(\Phi_i) \quad (3)$$

231

$$D_{RS,i} = \begin{cases} Median(\Phi_i) - P_{16}(\Phi_i), & \text{if } O_{R,i} < Median(\Phi_i) \\ P_{84}(\Phi_i) - Median(\Phi_i), & \text{if } O_{R,i} \geq Median(\Phi_i) \end{cases} \quad (4)$$

232

233 where P_{16} and P_{84} are the 16th and 84th percentiles, respectively.

234

| COMPONENT | COMPONENT OUTPUT | METRIC QUANTIFICATION DESCRIPTION |
|-------------|-------------------------------------|---|
| Achievement | Achievement output ($O_{A,i}$) | Outcome of interest crude measure. |
| | Achievement reference ($R_{A,i}$) | Achievement output from previous time-point. OR Mean achievement output over a time-period. OR A population standard. OR Other comparator of interest measured by the same metric and continuous scale. |

| | | |
|-------------|-------------------------------------|---|
| | Achievement deviation ($D_{A,i}$) | Deviation of: Individual achievement outputs over a time-period. OR Population achievement outputs from peers in cohort of interest. OR Population deviation for metric of interest. |
| | Lower and Upper reference limiters | Practical lower and upper reference limits. To establish a cohort range for consistent future comparative evaluations. |
| Rank | Final rank (R_i) | Entity rank in order of crude measures at evaluation time-point. |
| | Initial rank (ρ_i) | Entity rank at time-point of comparison prior to evaluation event or period. |
| | Transformed rank-scores $f(\rho_i)$ | Ranks transformed to a continuous value. Reflecting non-equidistance between entity ranks generated from crude measure. |
| | Rank score output ($O_{RS,i}$) | Magnitude difference of final rank position relative to peers. Aggregation of transformed rank-scores from entities with a final rank behind entity i . |
| | Rank score reference ($R_{RS,i}$) | Median of simulated measures Φ_i for each entity. |
| | Rank score deviation ($D_{RS,i}$) | Deviations of simulated rank score outputs. Absolute value of the difference between the median and either of P_{16} or P_{84} of simulated measures Φ_i for each entity. |

235 **Table 1.** Component parameters and how they are quantified for an entity 'i'.

236

237 **Step 3 – Component indicators (normalised component outputs)**

238 **Action:** Component outputs transformed to normalised measures reflecting the magnitude
239 change in achievement and rank.

240 **Outcome:** z-achievement and z-rank.

241 **Procedure:** Normalised z-scores are calculated based on the component outputs. Z-
242 achievement and z-rank indicators are established as proportional measures of output
243 deviation from their corresponding references.

244

245 For the achievement component:

246

$$z_{A,i} = \frac{\text{signum} \cdot (O_{A,i} - R_{A,i})}{D_{A,i}}$$

(5)

247 where

$$248 \quad \text{signum} = \begin{cases} 1, & \text{if lower achievement scores are better} \\ -1, & \text{if higher achievement scores are better} \end{cases}$$

249

250 For the rank component:

251

$$z_{RS,i} = \frac{O_{RS,i} - R_{RS,i}}{D_{RS,i}} \quad (6)$$

252

253 **Step 4 - Composite Score**

254 **Action:** Aggregated component indicators with optional weighting.

255 **Outcome:** Performance outcome score (POS)

256 **Procedure:** The component indicators are combined in a weighting procedure based upon a
257 user pre-defined setup of their respective relevance. An argument for the choice of weights
258 for each respective component is formulated. The POS is established together with an
259 explanation of the component weighting. When applying the option of proportional
260 weighting to the component indicators the larger the importance of the component, the
261 higher the weighting factor: w_A for achievement and w_R for rank. The weighting is
262 performed on centred z-scores for achievement $c_{A,i}$ and rank $c_{R,i}$ respectively:

$$c_{A,i} = z_{A,i} - M_A \quad (7)$$

263

$$c_{R,i} = z_{R,i} - M_R \quad (8)$$

264 where

$$M_A = \text{mean}_{\forall i}(z_{A,i}) \quad (9)$$

$$M_R = \text{mean}_{\forall i}(z_{RS,i}) \quad (10)$$

265

266 After the weighting, the score is shifted back so that the centre of the composite score
 267 reflects the centre of the achievement indicator, so, the composite score for *entity i* is then
 268 defined as

269

$$\begin{aligned} POS_i &= \frac{w_A c_{A,i} + w_R c_{R,i}}{w_A + w_R} + M_A = \\ &= \frac{w_A (Z_{A,i} - \text{mean}_{\forall i}(Z_{A,i})) + w_R (Z_{R,i} - \text{mean}_{\forall i}(Z_{R,i}))}{w_A + w_R} + \text{mean}_{\forall i}(Z_{A,i}) \end{aligned} \quad (11)$$

270

271 **POS-T application example**

272

273 The context is 36 OECD countries (excluding countries with incomplete data publicly
 274 available) that have set out to reduce their CO₂ emissions over the 10-year period 2006 –
 275 2015 in the global reduction of greenhouse effect (Table 2). Four entities (Sweden, Mexico,
 276 Norway, and Luxembourg) were hypothetically to be evaluated with regard to three
 277 performance goals:

278

279 Output goals:

280 Goal 1. “To evaluate reduction in CO₂ emissions per capita over the period 2006 - 2015”.

281 Goal 2. “To evaluate change in international ranking with respect to CO₂ emissions per
 282 capita”

283 Outcome goal:

284 Goal 3. “To evaluate reduction performance relative to peers regarding CO₂ emissions per
285 capita over the period 2006 - 2015”.

286

287 To develop performance outcome evaluation measures for these goals using the POS-T, the
288 comparative evaluation cohort is first described. The four countries (entities) will have their
289 performance evaluated against the cohort of 36 OECD countries. The comparator framing
290 the evaluation was defined as a comparison of performance over time.

291

| CONTEXT | DESCRIPTION |
|-----------------------------------|--|
| Entities to evaluate | Four countries (Sweden, Mexico, Norway, Luxembourg) |
| Comparison cohort | 36 OECD countries |
| Comparator | CO ₂ emissions per capita over the period 2006 - 2015 |
| Quantifiable domains (components) | Internal comparison - Achievement External comparison - Rank |

292 **Table 2.** Contextual features framing the evaluation.

293

294 **Example step 1 - Components (Quantifiable domains) and Parameters**

295 Tonnes of CO₂ equivalent (“CO₂ emissions”) was chosen as the achievement output measure
296 for the evaluation of Goal 1 (Table 3). A standardising parameter of ‘per capita’ was applied
297 to the achievement component to permit direct comparison between cohort countries on a
298 per capita basis. The achievement reference was defined as the average annual CO₂
299 emissions per capita during the 10-year period 1996-2005 and the achievement output for
300 comparison was defined as CO₂ emissions per capita in the year 2015. The OECD world
301 ranking table position was chosen as the rank output measure when evaluating Goal 2. To
302 evaluate Goal 3, the POS-T was used to develop a POS that depicts the performance
303 outcome regarding emission change over time for evaluation in the context of self-
304 comparison and comparison with peers.

305

| COMPONENT | PARAMETER | METRIC |
|--------------------|-------------------------|---|
| Achievement | Output measure | Tonnes of CO ₂ equivalent 2015 |
| | Standardising parameter | Per Capita |
| | Reference | 10-year annual output average 1996 - 2005 |
| | Distribution | Individual entity annual variations 1996 - 2005 |
| Rank | Output measure | OECD ranking table 2015 |
| | Reference | Rank of median achievement references |

306 **Table 3.** Components, Parameters, and metrics used to populate the component outputs.

307

308 **Example step 2 - Component Outputs (performance metrics from each component)**

309 Data were managed following the stepwise process detailed in the R-code software

310 provided (Data S1). Data for each component were collected ⁴² and presented in file format

311 (Data S2 and S3). Lower and upper reference limiters were set beyond the minimum and

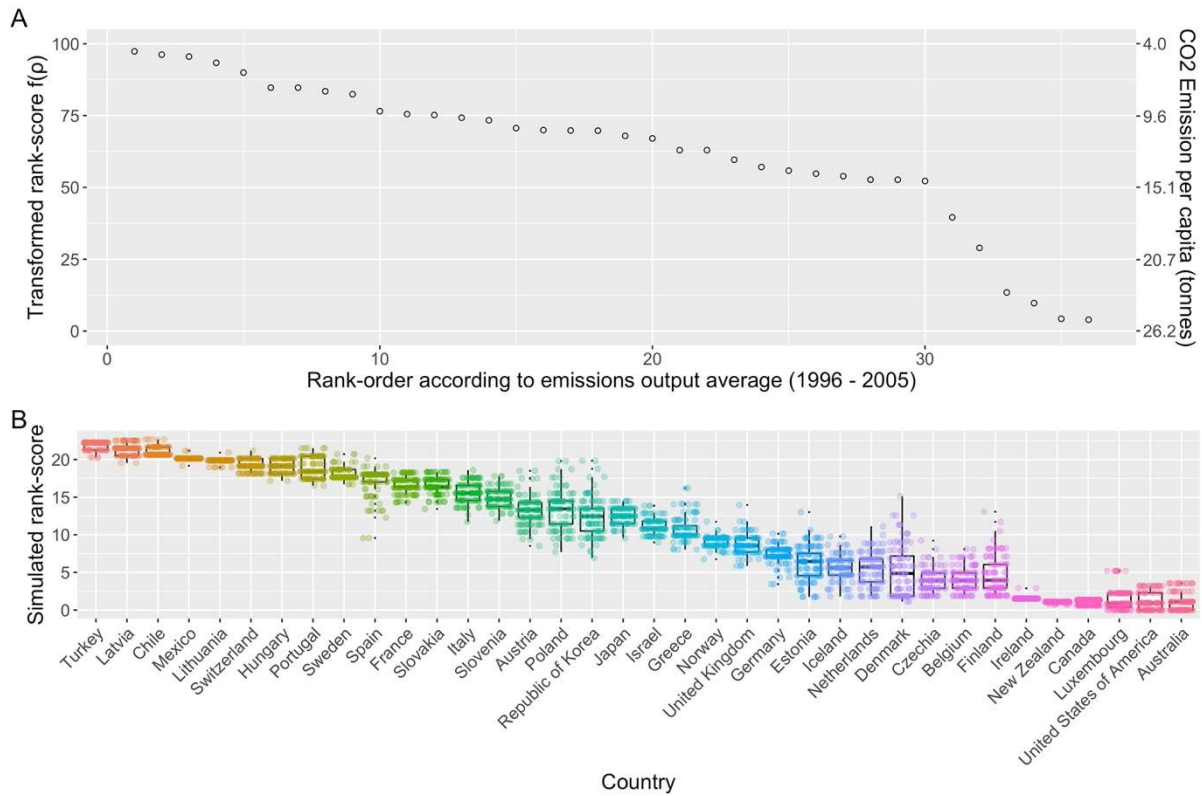
312 maximum CO₂ emissions per capita outputs. Individual entity standard deviations were

313 chosen for use in the simulations. Application of the POS-T R-code to the performance data

314 produced automated performance outputs and descriptive statistics. The following

315 component outputs were generated (Fig. 2b and Table 4).

316



317

318 **Fig 2. A. Entity rankings, B. Simulated rank scores (example data).** A. Initial rank-order (ranked by average

319 annual CO₂ emissions per capita 1996 - 2005) (ρ_i) vs average CO₂ emission per capita 1996 - 2005 (tonnes) and

320 transformed rank-score $f(\rho)$. B. Simulation (example with 100 iterations shown for illustrative purposes) of

321 potential rank score outputs (Φ_i) for each entity j based on underlying distribution of achievement output. Box

322 plot showing 16th and 84th percentiles.

323

| COMPONENT T | COMPONENT OUTPUT | SWEDEN | MEXICO | NORWAY | LUX. |
|--------------------|---|---|-------------------------|-----------------------|------------------------|
| Achievement | Achievement output ($O_{A,i}$) | 5.47 | 5.74 | 10.47 | 18.17 |
| | Achievement reference ($R_{A,j}$) | 7.90 | 5.48 | 12.23 | 24.06 |
| | Achievement deviation ($D_{A,i}$) | 0.34 | 0.10 | 0.14 | 2.78 |
| | Lower and upper reference limiters | Lower reference limiter 4 Upper reference limiter 26 | | | |
| Rank | Final rank (R_i) | 1 | 3 | 24 | 33 |
| | Initial rank (ρ_i) | 9 | 4 | 21 | 34 |
| | Transformed rank-score out of 100 ($f(\rho_i)$) | 82.44 | 93.36 | 62.97 | 9.72 |
| | Rank score output ($O_{RS,i}$) | 26.44 | 21.32 | 7.14 | 2.23 |
| | Rank score reference ($R_{RS,i}$) | 17.65 | 20.14 | 8.61 | 0.86 |
| | Rank score deviation ($D_{RS,i}$) | 1.01 (16.63 - 18.66) | 1.28 (18.86 - 21.42) | 0.05 (8.56 - 8.66) | 1.38 (-0.53 - 2.27) |

324 **Table 4.** Component outputs for the four stakeholders derived from Step 2 of the POS-T. (Units of achievement

325 = tonnes of CO₂ equivalent per capita. Rank score output = proportional score gained from ranking ahead of

326 other countries. Rank score reference = median score from simulation based on achievement descriptive
327 statistics. Rank score deviation = simulation outputs based on 16th and 84th percentiles. LUX. = Luxembourg.)
328

329 Goal 1. "To evaluate reduction in CO₂ emissions per capita over the period 2006 - 2015".

330 Component output: Three of the four countries reduced their raw CO₂ emissions per capita.

331 Luxembourg saw the largest reduction of the four example countries and largest reduction

332 compared to the full OECD cohort (-5.90 tonnes per capita), followed by Sweden (13th

333 overall; -2.43 tonnes per capita), and Norway (21st overall; -1.76 tonnes per capita). Mexico

334 saw an increase in CO₂ emissions (30th overall; +0.26 tonnes per capita).

335

336 Goal 2. "To evaluate change in international ranking with respect to CO₂ emissions per

337 capita over the period 2006 - 2015

338 Component output: Three of the four countries improved their ranking. Sweden (=3rd

339 largest rank shift overall: +8 places), Luxembourg and Mexico (=10th largest rank shift

340 overall: +1 place) relative to all 36 comparison countries. Norway fell in ranking (=27th

341 largest rank shift overall: -3 places) reflecting having not reduced their emissions per capita

342 to the same level over the observation period as the comparison cohort.

343

344 **Example step 3 – Component indicators (normalised component outputs)**

345 The component outputs were transformed to normalised measures demonstrating internal

346 and external magnitude change relative to the achievement and rank reference standards

347 respectively (Table 5.). The three countries that saw a reduction in raw CO₂ emissions per

348 capita demonstrated positive internal magnitude change (z-achievement: Sweden: +6.85,

349 Norway: +12.92, Luxembourg: +2.12). One country saw an increase in raw CO₂ emissions per

350 capita demonstrating negative internal magnitude change (z-achievement: Mexico: -2.56).
 351 Three of four countries improved their ranking demonstrating positive external magnitude
 352 change (z-rank: Sweden: +8.69, Mexico: +0.92, and Luxembourg: +1.00). One country
 353 regressed in ranking demonstrating a negative external magnitude change (Norway: -27.73).
 354

| INDICATOR | SWEDEN | MEXICO | NORWAY | LUXEMBOURG |
|---|--------|--------|--------|------------|
| z-achievement ($z_{A,i}$) | 6.85 | -2.56 | 12.92 | 2.12 |
| z-rank ($z_{RS,i}$) | 8.69 | 0.92 | -27.73 | 1.00 |

355 **Table 5.** Component indicators (normalised component outputs) derived from Step 3 of the POS-T.

356

357 **Example step 4 - Composite score (POS)**

358 The selected weighting ratio for achievement and rank was set at 1:1. The relative
 359 magnitude of change (component indicators) for both component outputs were combined,
 360 resulting in a composite score (POS) (Table 6).

361

362 Outcome goal:

363 Goal 3: “To evaluate reduction performance relative to peers regarding CO₂
 364 emissions per capita over the period 2006 - 2015”.

365 Composite score: The combined relative magnitude of change for the component
 366 indicators showed a positive development of performance for Sweden (+10.19),
 367 Luxembourg (+3.98), and Mexico (+1.60), and a negative development of
 368 performance for Norway (-4.99).

369

370 **Summary information gained for performance evaluation (4 stakeholder countries
 371 highlighted)**

372

| COUNTRY | CRUDE CO ₂ EMISSIONS CHANGE (<0 = reduction in CO ₂ emissions) | CRUDE RANK CHANGE (>0 = improved rank) | MAGNITUDE OF EMISSIONS CHANGE (z-achievement) | MAGNITUDE OF RANK CHANGE (z-rank) | PERFORMANCE OUTCOME SCORE (POS) |
|--------------------------|--|--|---|-----------------------------------|---------------------------------|
| Switzerland | -1.63 | 2 | 17.32 | 3.37 | <i>12.76</i> |
| Sweden | -2.43 | 8 | 6.85 | 8.69 | 10.19 |
| United Kingdom | -4.46 | 8 | 7.89 | 7.48 | <i>10.10</i> |
| United States of America | -4.63 | 0 | 12.72 | 0.00 | <i>8.78</i> |
| Ireland | -4.70 | 3 | 8.57 | 3.61 | <i>8.51</i> |
| Belgium | -4.14 | 6 | 8.34 | 2.86 | <i>8.02</i> |
| Italy | -2.56 | 3 | 9.29 | 1.51 | <i>7.81</i> |
| Hungary | -1.18 | 0 | 9.73 | 0.09 | <i>7.33</i> |
| France | -2.21 | -1 | 10.27 | -0.52 | <i>7.29</i> |
| Canada | -3.20 | -1 | 9.48 | -1.00 | <i>6.66</i> |
| Finland | -4.55 | 11 | 4.75 | 3.45 | <i>6.52</i> |
| Denmark | -5.59 | 11 | 3.93 | 3.12 | <i>5.94</i> |
| Greece | -2.38 | 3 | 4.73 | 2.25 | <i>5.91</i> |
| Slovakia | -1.85 | -1 | 7.65 | -1.00 | <i>5.74</i> |
| Slovenia | -1.80 | -1 | 6.82 | -0.90 | <i>5.38</i> |
| Czechia | -2.46 | 1 | 4.98 | 0.71 | <i>5.26</i> |
| New Zealand | -2.19 | 0 | 5.23 | 0.00 | <i>5.03</i> |
| Australia | -2.98 | 0 | 5.93 | -1.00 | <i>4.88</i> |
| Netherlands | -2.50 | 0 | 3.13 | 0.21 | <i>4.09</i> |
| Luxembourg | -5.90 | 1 | 2.12 | 1.00 | 3.98 |
| Israel | -1.02 | -1 | 3.44 | -0.42 | <i>3.93</i> |
| Germany | -1.88 | -2 | 3.47 | -0.81 | <i>3.75</i> |
| Spain | -1.98 | -1 | 2.95 | -0.33 | <i>3.73</i> |
| Austria | -1.43 | -3 | 2.88 | -0.82 | <i>3.45</i> |
| Portugal | -1.14 | 0 | 2.01 | 0.04 | <i>3.45</i> |
| Poland | -0.41 | -5 | 0.66 | -1.70 | <i>1.90</i> |
| Japan | -0.40 | -4 | 2.83 | -4.35 | <i>1.66</i> |
| Mexico | 0.26 | 1 | -2.56 | 0.92 | 1.60 |
| Estonia | 0.24 | -6 | -0.26 | -1.57 | <i>1.50</i> |
| Latvia | 0.79 | 0 | -3.28 | 0.15 | <i>0.85</i> |
| Lithuania | 0.66 | -4 | -1.65 | -2.91 | <i>0.14</i> |
| Iceland | 0.46 | -6 | -0.93 | -4.36 | <i>-0.22</i> |
| Chile | 1.03 | -3 | -5.31 | -2.13 | <i>-1.30</i> |
| Republic of Korea | 2.92 | -12 | -4.05 | -3.40 | <i>-1.31</i> |
| Turkey | 1.42 | -4 | -7.94 | -3.71 | <i>-3.40</i> |
| Norway | -1.76 | -3 | 12.92 | -27.73 | -4.99 |

373 **Table 6:** OECD countries ordered highest to lowest by the POS including component crude output variation

374 and component magnitude of change over the observation period. (Crude CO₂ emissions = tonnes of CO₂

375 equivalent per capita).

376

377 Discussion

378 This study set out to develop a scoring template that combines internal and external
379 measures of performance by alignment of measurement scales which represent meaningful
380 magnitudes of change. The resulting POS-T adheres to the principle of providing the
381 stakeholders governing applied programs a means to report performance outcomes with
382 convincing face-validity^{37,43,44}. We exemplified application of POS-T in an evaluation of CO₂
383 emission reduction amongst OECD member countries. Flexible and transparent evaluation
384 methods oriented towards stakeholders and usefulness have repeatedly been asked for in
385 the environmental sciences^{45,46}. To ensure that POS-T produces scores useful for
386 stakeholders, an inductive (discovery) approach was found best suited.^{47,48} This approach
387 aligns with the principles of design thinking⁴⁹ where the emphasis is placed on defining the
388 problem to be solved through the needs of the stakeholders involved.^{43,44} In the following,
389 the main features of the POS-T are discussed considering the CO₂ emission example and
390 directions are outlined for future research.

391 The measures and methods available to report performance results delimit a stakeholder's
392 capacity to evaluate applied programs.³⁸ In their governance, accomplishment has been
393 described as the gap between expected and actual output or the deviation of the output
394 from an industry standard.¹¹ The POS-T provides in Step 2 methodology to establish a
395 *magnitude* of this gap for both the achievement and rank components, i.e. for internal and
396 external comparisons. In step 3, ready-to-combine component indicators are formed by
397 normalisation of the component outputs produced. The resulting component indicators
398 describe magnitudes of change, i.e. the gaps between expected and actual results or change
399 over time for internal and external comparisons. Producing the relative magnitude of

400 change from a point of reference for both components (achievement and rank) rather than
401 binary measures alone provides those evaluating the outcome with greater context. For
402 example, on a binary scale, Sweden, Luxembourg, and Norway each demonstrated a
403 *reduction* in CO₂ emissions per capita (positive achievement outcome) relative to their own
404 reference standard in 2005. However, when accounting for the external context, Norway
405 during the evaluation period slipped down the ranking table from 21st to 24th (negative rank
406 outcome) due to that other countries reduced their relative emissions by greater amounts.
407 Conversely, Mexico gained a ranking place from 4th to 3rd (positive rank outcome) even
408 having had a small increase in crude CO₂ emissions due to that other countries close in rank
409 had relatively larger increases in crude CO₂ emissions. Proportional output measures
410 relative to self-comparison and comparison with others in a chosen cohort provides a
411 context for stakeholders to better frame and evaluate an outcome against expectation and
412 describe the overall accomplishment.

413 Integration of self-comparison and rank change magnitudes adds complexity to program
414 evaluation indicators. Maintenance of face-validity in such composite scores requires
415 measurement system transparency.⁵⁰ The POS-T supports transparency and face-validity by
416 offering evaluators semantic clarity regarding the components of the integrated composite
417 score. Stakeholders evaluating performance using the POS-T will base their assessments on
418 normalised indicators of any measured or pre-existing method of reporting achievement
419 output for internal (within-individual) comparisons. The crude achievement outputs can in
420 any circumstance be ranked ⁵¹ and the normalised indicators computed for rank changes
421 and external comparisons (between individuals). The normalised indicators are calculated in
422 a standardised manner, i.e. for 'achievement' by subtracting the mean from an individual

423 raw score and then dividing the difference by the standard deviation. The mean and
424 standard deviation are based on individual, or population standards as chosen in step 1 by
425 the stakeholder evaluating the performance. Normalised indicators are calculated for 'rank'
426 by transforming the ordinal scale to continuous relative values before applying the same
427 normalisation process to a series of simulated rank outputs. Such normalised indicators are
428 well-known and are broadly used in, for instance, global health settings for comparative
429 evaluations of development processes at individual and population levels, e.g. in the child
430 growth area.^{52,53} In the example application of POS-T on CO₂ emissions, the normalised 'self-
431 comparison' indicator showcases that Luxembourg's crude emissions reduction from the
432 reference of 2005 is achieved in the context of a broader distribution of the annual
433 fluctuation in emissions by Luxembourg compared to Sweden. In essence, Sweden's
434 emissions reduction achievement is at face value more substantial in the context of self-
435 comparison due to the narrower distribution of annual emissions fluctuations. The
436 magnitude of this achievement is demonstrated by a higher achievement indicator.
437 Regarding external comparisons, both Luxembourg and Mexico gained one place on the
438 ranking table, yet the normalised rank indicator calculated in step 3 shows that the
439 magnitude of Luxembourg's gain is greater than Mexico's. This is due to that Mexico's
440 emissions per capita were very close in volume to those with similar rank, the effect being
441 that a change in rank may occur even from small changes in emissions output resulting in a
442 smaller rank indicator for the same crude rank change. By using the continuous data that
443 formulates the ranking order, context and magnitude is apportioned to the component
444 indicator in step 3 representing the rank change. Internal and external measures presented
445 as a magnitude of change against a reference and accompanying distribution provide
446 meaningful context to the performance. The selection of the lower and upper reference

447 limiters provides an important step when applying the statistical code in establishing
448 consistent comparators and context for future equivalent evaluations. Performance
449 outcomes presented this way can be used to observe longitudinal performance trends
450 within an individual entity or relative to a population as well as measuring a single
451 performance against expectation.

452 Some limitations to the use of the POS-T are important to consider. In experimental
453 evaluation research, influence from external factors is controlled in the study design.
454 Emulating an experimental design in observational performance evaluations in practice
455 settings would require information on all confounding factors.⁵⁴ Application of the POS-T
456 does not per se assure that the POS reflects causal effects of the program, and
457 consideration of confounding factors is always needed when interpreting POS scores in
458 practice settings. Moreover, it should be taken into regard that the simulation process used
459 in POS-T to determine the rank score deviations for each entity uses the reference and its
460 standard deviation as assumptions in the calculation. The outcome from each simulation
461 may thus vary slightly. This effect is minimised by always running an adequate number of
462 simulations on each occasion. Furthermore, when the data available to calculate the
463 achievement deviation is limited, a decision must be made regarding what to use as the
464 achievement deviation for comparison with the achievement output. The preferred option
465 is to use the achievement deviation unique to each entity. However, an option is to use the
466 cohort population standard deviation as this broadens the dataset to calculate deviations
467 and improves its reliability. This may be a satisfactory solution when the comparative data
468 sets between entities in the cohort have similar deviations. If this is not the case, an option
469 may be to use the largest or smallest deviation in the cohort. The flexibility in selecting

470 components, standardising parameters and weighting of the component indicators opens
471 the composite score to variability in its robustness. Testing for robustness is recommended
472 and aided by the level of transparency described by the evaluator in selecting optional
473 features in corresponding steps of the POS-T framework.

474

475 The POS-T in its current form can be applied to any program governance setting. A POS can
476 be determined for single entities at multiple time points to assess performance trends or for
477 multiple entities at a single time point to assess performances relative to a population
478 standard or to peers. The component indicators can also be used to evaluate each
479 component in isolation. The rank simulation process in isolation may furthermore be utilised
480 to determine probabilities of performance outcome. Further development of the POS-T will
481 include development of the statistical code to include the evaluation of performance in
482 settings where achievement is not readily quantified, e.g. when it mainly is established
483 through head-to-head contests.

484

485 **Conclusion**

486 The POS-T endorses face-validity in real-world program evaluations by that the resulting
487 POS reflects a meaningful magnitude of performance outcome with regards to self-
488 comparison and comparison with peers. The template is presented with statistical software
489 for creating scoring systems and is exemplified by evaluation of CO₂ emissions reduction
490 amongst 36 OECD member countries. Forthcoming research will involve application of the
491 POS in different applied performance evaluation settings.

492

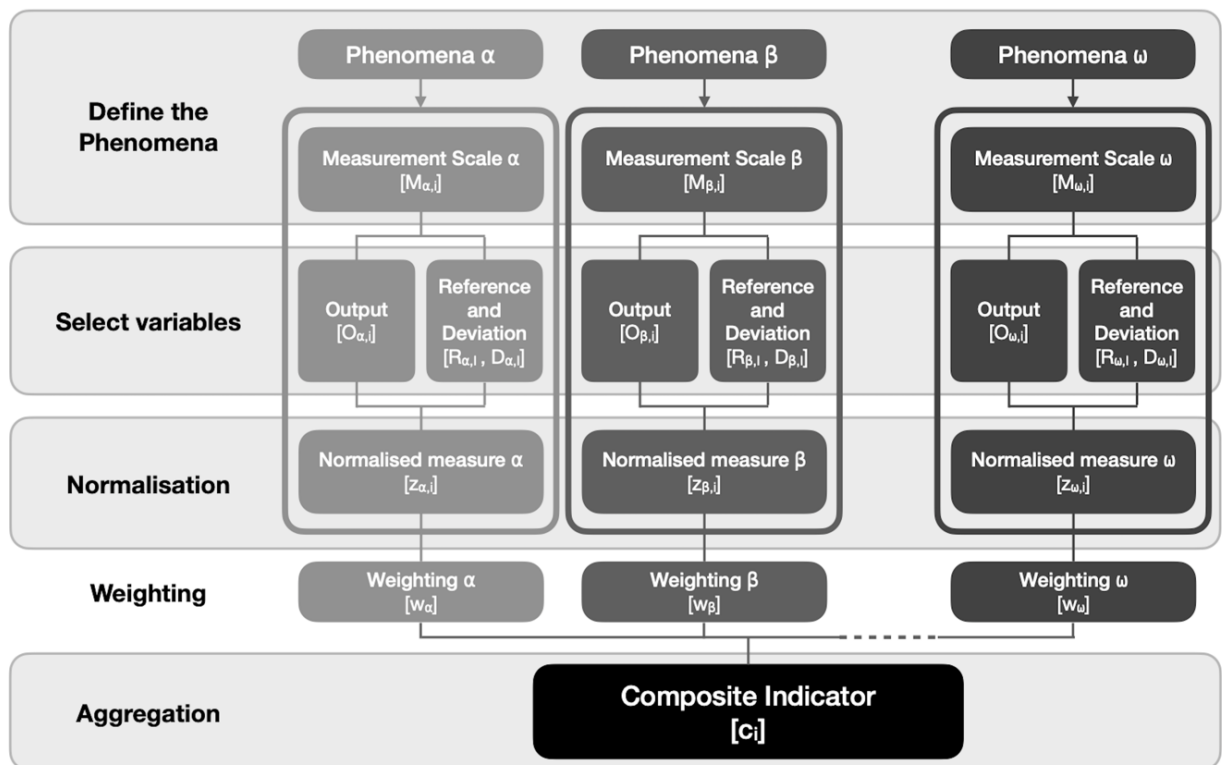
493 **Materials and Methods**

494

495 **Construction principles**

496 Construction of the POS-T employed an iterative approach to solution design that prioritises
497 application of the final template in real-world program governance settings. Its practical use
498 was further supported through the parallel construction of a statistical framework and
499 software for score development in applications.⁴⁹ A design panel was composed for the
500 construction consisting of scientists and practitioners (n=5) with backgrounds in
501 epidemiology, public health and sport settings, organisational development, statistical
502 methods, and experimental design. A composite score development model was used to
503 guide the construction process (Fig. 3).

504



505

506 **Fig. 3. Composite score development model with five development steps and three example phenomena.**

507

508 **Template construction process**

509 The design panel met via an online meeting platform weekly over a twelve-month period
 510 and discussed the POS-T development in the context of four cardinal steps and one optional
 511 step depicted in the construction model (Fig. 3). In each development model step, the
 512 design panel employed an iterative process applying varying methods of data analysis and
 513 representation in the template to identify potential inconsistencies or errors in the
 514 composite score. These were identified, discussed, and addressed at each online meeting
 515 until panel consensus on user application was reached. Consensus required agreement on
 516 the stepwise process necessary for the user when defining the context of POS-T evaluation.

517 Once the process was established and incorporated into the POS-T consensus on the
518 maintenance of the aggregated composite score face-validity was obligatory.

519

520 Statistical code was written using R programming language to automate the methodological
521 outputs of steps 2-4 (Data S1). Data comprising the variables outlined in step 2 of the
522 development model were systematically organised in data files using Microsoft Excel 2016
523 (Data S2 and S3). The current version of the R-code was written to use with discrete cohorts
524 that comprise all entities in sequential ranking order for analysis.

525

526 **POS-T application example**

527 The resulting POS-T was finally applied in an evaluation of reduction of greenhouse gas
528 emission among 36 OECD countries. The OECD has collected and reported annual
529 greenhouse gas emissions between 1996 and 2015 for a cohort of 36 countries producing an
530 annual emissions ranking table.⁴² The OECD data set was used to showcase performance
531 evaluations based on a composite score for comparison within either of the singular
532 reporting metrics; tonnes of CO₂ emissions equivalent per capita, or the emissions ranking
533 table. To showcase the template and the statistical framework, the design panel took on the
534 virtual role of an international stakeholder commission in the environmental protection
535 area. The final statistical framework was exemplified by applying the R-code to the 36
536 country OECD data set and analysing the performance outcomes of four countries: Sweden,
537 Luxembourg, Mexico, and Norway.

538

539 **References**

- 540 1. "Performance Evaluation Methods: Measurement and Attribution of Program Results", (Treasury
541 Board of Canada, Secretariat, Toronto, 1998).
- 542 2. A. DeGroff, M. Schooley, T. Chapel, T. H. Poister, Challenges and strategies in applying performance
543 measurement to federal public health programs. *Evaluation and program planning* **33**, 365-372
544 (2010).
- 545 3. "Performance measurement and evaluation: Definitions and relationships." (U.S. Government
546 Accountability Office, GAO-11-646SP, Washington DC, 2011).
- 547 4. P. de Lancer Julnes, Performance measurement: An effective tool for government accountability? The
548 debate goes on. *Evaluation* **12**, 219-235 (2006).
- 549 5. M. Lebas, Performance measurement and performance management. *International Journal of*
550 *Production Economics* **41**, 23-35 (1995).
- 551 6. K. E. Newcomer, Using performance measurement to improve programs. *New directions for*
552 *evaluation* **1997**, 5-14 (1997).
- 553 7. T. Rantala. "Operational level performance measurement in university-industry collaboration", thesis,
554 LUT University, Finland (2019).
- 555 8. B. Milstein, S. Wetterhall, C. E. W. Group, A framework featuring steps and standards for program
556 evaluation. *Health Promotion Practice* **1**, 221-228 (2000).
- 557 9. J. P. Koplan, R. Milstein, S. Wetterhall, Framework for program evaluation in public health. *MMWR:*
558 *Recommendations and Reports* **48**, 1-40 (1999).
- 559 10. L. G. Suter, C. E. Barber, J. Herrin, A. Leong, E. Losina, A. Miller, E. Newman, M. Robbins, H. Tory, J.
560 Yazdany, American College of Rheumatology white paper on performance outcome measures in
561 rheumatology. *Arthritis care & research* **68**, 1390-1401 (2016).
- 562 11. W. D. Savedoff, "Governance in the health sector: a strategy for measuring determinants and
563 performance", (Policy Research working paper; no. WPS 5655, The World Bank, 2011).
- 564 12. L. Laurian, J. Crawford, M. Day, P. Kouwenhoven, G. Mason, N. Ericksen, L. Beattie, Evaluating the
565 outcomes of plans: Theory, practice, and methodology. *Environment and Planning B: Planning and*
566 *Design* **37**, 740-757 (2010).

- 567 13. A. Neely, M. Gregory, K. Platts, Performance measurement system design: a literature review and
568 research agenda. *International journal of operations & production management* **15**, 80-116 (1995).
- 569 14. B. P. Raysmith, J. Jacobsson, M. K. Drew, T. Timpka, What Is Performance? A Scoping Review of
570 Performance Outcomes as Study Endpoints in Athletics. *Sports* **7**, 66 (2019).
- 571 15. "The world health report 2000: health systems: improving performance", (World Health Organization,
572 2000).
- 573 16. D. D. Dill, M. Soo, Academic quality, league tables, and public policy: A cross-national analysis of
574 university ranking systems. *Higher education* **49**, 495-533 (2005).
- 575 17. F. Derrien, O. Dessaint, The effects of investment bank rankings: Evidence from M&A league tables.
576 *Review of Finance* **22**, 1375-1411 (2018).
- 577 18. S. G. Armesto, M. L. G. Lapetra, L. Wei, E. Kelley, "Health care quality indicators project 2006 data
578 collection update report.", (OECD Health Working Papers, No. 29, OECD Publishing, Paris, 2007).
- 579 19. M.-J. Huang, M.-Y. Chen, K. Yieh, Comparing with your main competitor: the single most important
580 task of knowledge management performance measurement. *Journal of Information Science* **33**, 416-
581 434 (2007).
- 582 20. A. Mehrpouya, R. Samiolo, Performance measurement in global governance: Ranking and the politics
583 of variability. *Accounting, organizations and society* **55**, 12-31 (2016).
- 584 21. T. R. Oliver, Peer Reviewed: Population Health Rankings as Policy Indicators and Performance
585 Measures. *Preventing chronic disease* **7**, (2010).
- 586 22. K. Keasey, P. Moon, D. Duxbury, Performance measurement and the use of league tables: some
587 experimental evidence of dysfunctional consequences. *Accounting and business research* **30**, 275-286
588 (2000).
- 589 23. M.-H. Huang, A comparison of three major academic rankings for world universities: From a research
590 evaluation perspective. *Journal of Library & Information Studies* **9**, (2011).
- 591 24. S. Schütte, P. N. M. Acevedo, A. Flahault, Health systems around the world—a comparison of existing
592 health system rankings. *Journal of global health* **8**, (2018).
- 593 25. J. Hurst, M. Jee-Hughes, "Performance measurement and performance management in OECD health
594 systems.", (OECD Publishing, Paris, 2001).

- 595 26. R. S. Kaplan, Strategic performance measurement and management in nonprofit organizations.
596 *Nonprofit management and Leadership* **11**, 353-370 (2001).
- 597 27. R. Kaplan, D. Norton, The Balanced Scorecard—Measures that Drive Performance. *Harvard Business*
598 *Review*, 71-79 (1992).
- 599 28. J. Fahlén, The trust–mistrust dynamic in the public governance of sport: exploring the legitimacy of
600 performance measurement systems through end-users’ perceptions. *International journal of sport*
601 *policy and politics* **9**, 707-722 (2017).
- 602 29. I. O'Boyle, D. Hassan, Performance management and measurement in national-level non-profit sport
603 organisations. *European Sport Management Quarterly* **14**, 299-314 (2014).
- 604 30. H. Pham, B. G. Sutton, P. J. Brown, D. A. Brown, Moving towards sustainability: A theoretical design of
605 environmental performance measurement systems. *Journal of Cleaner Production* **269**, 122273
606 (2020).
- 607 31. N. Carter, P. Day, R. Klein, *How organisations measure success: the use of performance indicators in*
608 *government* (Psychology Press, 1995).
- 609 32. R. Jacobs, M. Goddard, P. C. Smith, Composite performance measures in the public sector. *Center for*
610 *Health Economics Research Paper* **16**, (2007).
- 611 33. I. C. Davies, Evaluation and performance management in government. *Evaluation* **5**, 150-159 (1999).
- 612 34. “Beyond roe-how to measure bank performance.”, (*Appendix to the report on EU banking structures*,
613 European Central Bank, Frankfurt, 2010).
- 614 35. L. J. White, *The credit rating agencies and their role in the financial system*. (Oxford University Press,
615 London, 2018).
- 616 36. T. Seeley, "Cumminuty based organization (CBO) survey results: Outcome evaluation in voluntary and
617 not-for-profit organizations.", (The Muttard Fellowship, Edmonton, 2003).
- 618 37. D. Guyadeen, M. Seasons, Evaluation theory and practice: Comparing program evaluation and
619 evaluation in planning. *Journal of Planning Education and Research* **38**, 98-110 (2018).
- 620 38. E. D. Peterson, E. R. DeLong, F. A. Masoudi, S. M. O'Brien, P. N. Peterson, J. S. Rumsfeld, D. M.
621 Shahian, R. E. Shaw, ACCF/AHA 2010 position statement on composite measures for healthcare
622 performance assessment: American College of Cardiology Foundation/American Heart Association

- 623 Task Force on performance measures (writing committee to develop a position statement on
624 composite measures). *Journal of the American College of Cardiology* **55**, 1755-1766 (2010).
- 625 39. R. Jacobs, M. Goddard, P. C. Smith, How robust are hospital ranks based on composite performance
626 measures? *Medical care*, 1177-1184 (2005).
- 627 40. M. Nardo, M. Saisana, A. Saltelli, S. Tarantola, A. Hoffman, E. Giovannini, "Handbook on Constructing
628 Composite Indicators: Methodology and user guide.", (OECD Publishing, Paris, 2005).
- 629 41. M.-K. Song, F.-C. Lin, S. E. Ward, J. P. Fine, Composite variables: when and how. *Nursing research* **62**,
630 45 (2013).
- 631 42. "Greenhouse gas emissions by source" Internet. OECD Publishing; 2014. Available from:
632 <https://www.oecd-ilibrary.org/content/data/data-00594-en>.
- 633 43. A-M. R. McGowan, C. Bakula, R. S. Castner, "Lessons Learned from Applying Design Thinking in a NASA
634 Rapid Design Study in Aeronautics ", in *58th AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and*
635 *Materials Conference*. (2017), pp. 0976.
- 636 44. M. Greene, "Systems Design Thinking: Identification and Measurement of Attitudes for Systems
637 Engineering, Systems Thinking, and Design Thinking", thesis, University of Michigan, (2019).
- 638 45. V. J. Schwanitz, Evaluating integrated assessment models of global climate change. *Environmental*
639 *modelling & software* **50**, 120-131 (2013).
- 640 46. M. Haasnoot, W. Van Deursen, J. H. Guillaume, J. H. Kwakkel, E. van Beek, H. Middelkoop, Fit for
641 purpose? Building and evaluating a fast, integrated model for exploring water policy pathways.
642 *Environmental modelling & software* **60**, 99-120 (2014).
- 643 47. J. D. Gould, C. Lewis, Designing for usability: key principles and what designers think. *Communications*
644 *of the ACM* **28**, 300-311 (1985).
- 645 48. K. Dorst, The core of 'design thinking' and its application. *Design studies* **32**, 521-532 (2011).
- 646 49. T. Brown, Design thinking. *Harvard business review* **86**, 84 (2008).
- 647 50. A. Rae, C. Wong, Monitoring spatial planning policies: towards an analytical, adaptive, and spatial
648 approach to a 'wicked problem'. *Environment and Planning B: Planning and Design* **39**, 880-896
649 (2012).

- 650 51. M. Mohsin, A. Rasheed, H. Sun, J. Zhang, R. Iram, N. Iqbal, Q. Abbas, Developing low carbon
651 economies: An aggregated composite index based on carbon emissions. *Sustainable Energy*
652 *Technologies and Assessments* **35**, 365-374 (2019).
- 653 52. J. M. Wit, J. H. Himes, S. Van Buuren, D. M. Denno, P. S. Suchdev, Practical application of linear
654 growth measurements in clinical research in low-and middle-income countries. *Hormone research in*
655 *paediatrics* **88**, 79-90 (2017).
- 656 53. M. Leung, N. Perumal, E. Mesfin, A. Krishna, S. Yang, W. Johnson, D. G. Bassani, D. E. Roth, Metrics of
657 early childhood growth in recent epidemiological research: a scoping review. *PloS one* **13**, e0194565
658 (2018).
- 659 54. M. A. Hernán, J. M. Robins, Using big data to emulate a target trial when a randomized trial is not
660 available. *American journal of epidemiology* **183**, 758-764 (2016).
- 661

662 **Acknowledgments**

663 **Funding**

664 The authors acknowledge that they received no funding in support for this research.

665

666 **Author contribution**

667 BPR provided conceptual basis and primary author, OD developed the statistical code and
668 refined conceptual thinking to practical outputs, TT guided the projects and provided major
669 edits to manuscript, MKD and JJ contributed to conceptual development and manuscript
670 edits.

671

672 **Competing interests**

673 The authors declare that they have no competing interests.

674

675 **Data and material availability**

676 All data are available in the main text or the supplementary materials.

677

678 **Figure and table legends**

679

680 **Fig. 1. POS-T signifying four cardinal and one optional (weighting) data handling levels.**

681 **Table 1.** Component parameters and how they are quantified for an entity 'i'.

682 **Table 2.** Contextual features framing the evaluation.

683 **Table 3.** Components, Parameters, and metrics used to populate the component outputs.

684 **Fig 2. A. Entity rankings, B. Simulated rank scores (example data).** A. Initial rank-order (ranked by average

685 annual CO₂ emissions per capita 1996 - 2005) (ρ_i) vs average CO₂ emission per capita 1996 - 2005 (tonnes) and

686 transformed rank-score $f(\rho)$. B. Simulation (example with 100 iterations shown for illustrative purposes) of

687 potential rank score outputs (Φ_i) for each entity based on underlying distribution of achievement output. Box

688 plot showing 16th and 84th percentiles.

689 **Table 4.** Component outputs for the four stakeholders derived from Step 2 of the POS-T. (Units of achievement

690 = tonnes of CO₂ equivalent per capita. Rank score output = proportional score gained from ranking ahead of

691 other countries. Rank score reference = median score from simulation based on achievement descriptive

692 statistics. Rank score deviation = simulation outputs based on 16th and 84th percentiles. LUX. = Luxembourg.)

693 **Table 5.** Component indicators (normalised component outputs) derived from Step 3 of the POS-T.

694 **Table 6:** OECD countries ordered highest to lowest by the POS including component crude output variation

695 and component magnitude of change over the observation period. (Crude CO₂ emissions = tonnes of CO₂

696 equivalent per capita).

697 **Fig. 3. Composite score development model with five development steps and three example phenomena.**

698

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [S1POST.r](#)
- [S2OECDfoundationdata.xlsx](#)
- [S3OECDdrinkingdata.xlsx](#)