



OPEN

Design of the performance outcome scoring template (POS-T) with example application on CO₂ emissions reduction amongst 36 OECD member countries

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

In applied program settings, such as in natural environment control and education, performance evaluation is usually conducted by evaluators considering both self-comparison and comparison with peers. We have developed the performance outcome scoring template (POS-T) for assessments with high face-validity in these settings. POS-T puts achievements of individuals or groups in context, i.e. the resulting performance outcome score (POS) reflects a meaningful measure of performance magnitude with regards to internal and external comparisons. Development of a POS is performed in four steps supported by a statistical framework. Software is supplied for creation of scoring applications in different performance evaluation settings. We demonstrate the POS-T by evaluation of CO₂ emissions reduction amongst 36 OECD member countries.

Performance evaluation seeks to examine the achievement of predetermined objectives or goals by individuals or groups through the broad assessment of processes (inputs and activities) and results (outputs and outcomes)^{1–3}. These evaluations are used to assess program efficiency and effectiveness and provide accountability to resource allocation, strategy and policy direction^{3–5}. Performance evaluations are usually case-specific and defined by the stakeholders with the authority and responsibility to do so^{6,7}. Furthermore, they are expected to provide a contextual judgement of performance at a moment in time and require the measurement of credible ongoing outputs (performance measures) that relate specifically to the needs of the evaluators^{3,8–10}. The assessment of goal achievement lies with the process of performance evaluation and considers a broad array of factors that include addressing the “How” and “Why” questions of achieving pre-determined objectives. The complexity associated with performance evaluations has led to the development of performance measurement systems to collect ongoing data and to monitor and report progress towards pre-determined goals^{7,8}.

Performance evaluations addressed through stakeholder questions require that the measures used to report performance results are relevant and trustworthy¹⁰. A performance ‘outcome’ is defined as the resultant effect of a system towards a pre-determined objective, whereas a performance ‘output’ is the data generated by a single unique metric impacting the outcome^{3,11–13}. A balanced performance evaluation includes both internal measurement of the individual units in the program and of the external environment^{14,15}. These measurements permit a reflection on achievement to date as well as what *could* be possible to achieve within an equal context¹⁶. For instance, internal measures of output include ‘exam scores’ in an academic setting, ‘race finish time’ in a sport setting, ‘greenhouse gas emissions per capita’ in an environment setting. Each of these metrics can be used as a performance output measure that represents ‘self-comparison’ when collected over time. However, internal measures alone may be insufficient in performance evaluation as the appraisal of a performance output or trend can vary when assessed relatively against peer performance under comparable circumstances. External environment

¹Athletics Research Centre, Linköping University, Linköping, Sweden. ²Western Australian Institute of Sport, Perth, Australia. ³Athletics Australia, Melbourne, Australia. ⁴Swedish Athletics Association, Stockholm, Sweden. ⁵Athlete Availability Program, Australian Institute of Sport, Bruce, ACT, Australia. ⁶Australian Collaboration for Research into Injury in Sport and Its Prevention (ACRISP), Perth, Australia. ⁷University of Canberra Research Institute for Sport and Exercise (UCRISE), Canberra, ACT, Australia. ⁸Department of Behavioural Sciences and Learning, Linköping University, Linköping, Sweden. ✉email: ben.raysmith@athletics.org.au

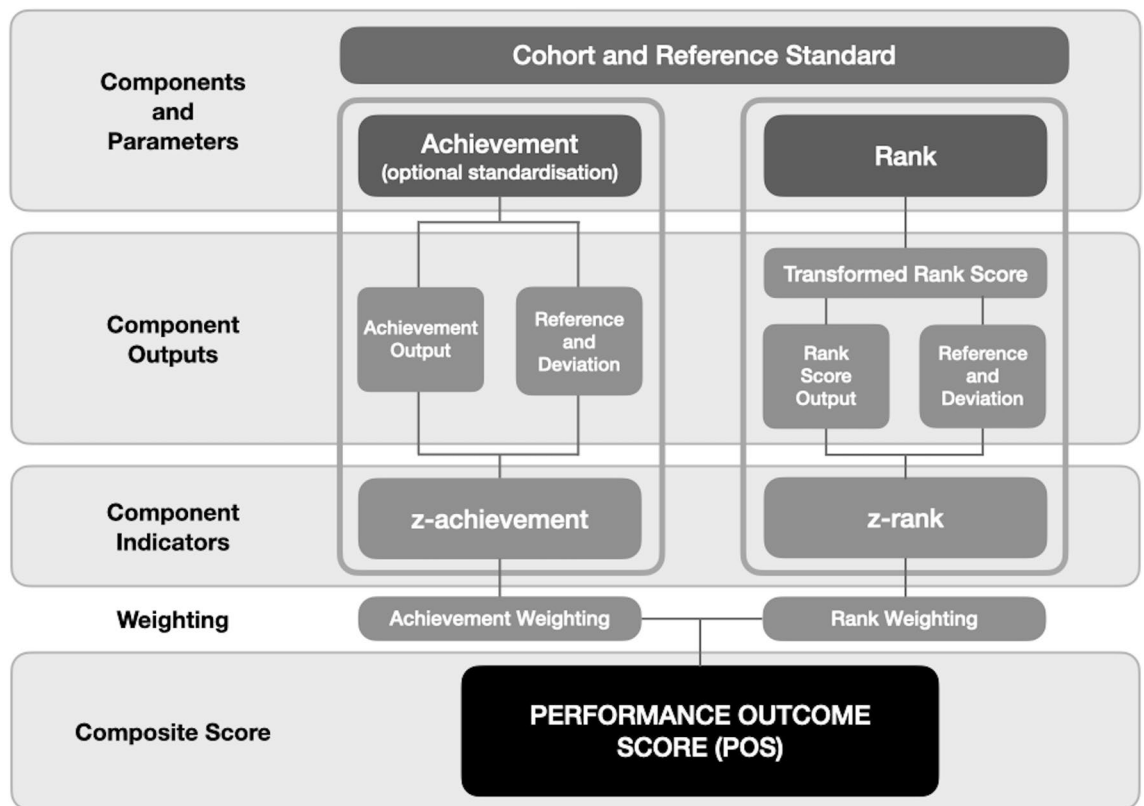


Figure 1. POS-T signifying four cardinal and one optional (weighting) data handling levels.

measures of output include ‘league tables’, ‘event final rankings’ or ‘comparisons with industry benchmarks’ and have become a popular way to compare peers within industries^{17–20}. As with internal measures, external measures of relative performance in isolation equally may lack face-validity or attribution (singular allocation) in performance evaluation. This is seen in circumstances where the appraisal of a ranking against peers can vary when interpreted in the context of individual achievement or progress^{21–23}.

Examples of bespoke performance measurement systems that comprise singular and multiple output measures exist across industry and sectors: academic²⁴, health^{25,26}, profit and non-profit organisations^{27,28}, sport^{29,30}, natural environment^{31,32}, government and private sectors^{33–35}, and finance^{36,37}. The use of a single output measure to reflect performance outcome may provide a too narrow perspective on achievement in a complex program and therefore risk evaluation face-validity in any of these areas¹⁰. The selection of performance output measures that enhance attribution and enable meaningful and trustworthy evaluations therefore benefits from appropriate consideration^{26,38,39}. Different measurements of performance can be combined into composite scores to increase the face-validity of a program evaluation appraisal without adding difficulty to its interpretation^{40,41}. Measurements on different data scales need here to be transformed into a common scale before combining, even though the introduction of re-scaling may reduce reliability⁴². An alignment of scales that reflect a meaningful magnitude change within and between variables can improve face-validity in regards to program outcomes and thereby make a composite score preferred for decision-making^{40,43}. It is therefore essential that the measurements are transparent and comprehensible for all stakeholders involved to avoid misleading program decisions using such composite scores¹⁰. A robust metric must be determined to assess objectively what has occurred and how this may influence future outcomes relative to the investment in any program or individual and any third-party interest.

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

Results

Application of the performance outcome scoring template (POS-T) commences with selection of data sources and concludes with a composite score that is adjusted for optimal face-validity (POS) (Fig. 1). Subject to the evaluation purpose and selected time-point the template can be utilised in comparing a result with a predetermined objective, benchmark standard or appraise change over time. Statistical software (in the R language) is supplied to support the development of a POS (Data S1). Stepwise instructions provided in the software detail file data set-up for application in the code.

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

- *Components and parameters* Quantifiable domains deemed to have primacy with respect to the face-validity of the performance outcome. Stakeholder selected parameters that frame the evaluation.
- *Component outputs* Performance metrics collected from each component. Representative of the data collection fields that comprise the performance output measure and a comparative reference output within a distribution.
- *Component indicators* Component outputs transformed to normalised measures. The performance outputs measured against a reference standard given assumptions of both achievement and rank deviations.
- *Composite score* A final viable measure is secured through aggregating and optionally weighting normalised component indicators that meet the desired face-validity.

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

Step 1: components (quantifiable domains) and parameters. *Actions.* Define evaluation entity (individual, group or population participating in a specified program), sample of entities to evaluate, comparison cohort, and comparator (frame of reference for evaluation). Select quantifiable domains representing internal and external measures that provide face-validity for the performance outcome of an entity (individual or group). Optional selection of a standardising parameter.

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

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

Step 2: component outputs (performance metrics from each component). *Action.* Collect/calculate performance metrics from each component.

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

Procedure. Performance metrics are recorded for both achievement and rank. The ‘achievement output’ and ‘rank score output’ are established as well as descriptive data for each component parameter, i.e. references and deviations (Table 1). This is completed for each ‘entity’ (individuals or groups) in the cohort. When referring to separate entity’s subscripts ‘*i*’ and ‘*j*’ are used.

Achievement performance metrics. *Output.* The ‘achievement output’ (O_A) is quantified directly from the metric of interest crude measure (continuous scale) and reflects the output *being evaluated*.

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

Deviation. The ‘achievement deviation’ (D_A) is the deviation of crude measures across the observation period or other collection of entity crude measures that constitute a deviation around the achievement reference. The reference and deviation values are set from the decision to use a cohort pooled achievement standard deviation or individual entity standard deviations. The achievement outputs are assumed to follow a normal distribution, $N(\bar{A}_i, \sigma_{A_i})$, from which the reference and deviation values are set; $R_{A,i} = \bar{A}_i$ and $D_{A,i} = \sigma_{A_i}$.

Rank performance metrics. Rank scores (forming the ordinal ranking order) are transformed to a continuous scale value (transformed rank-score) using a pre-defined function, f , to reflect non-equidistance (magnitude difference) between different entities based on the continuous metric that established the ranking order. Lower and upper ‘reference limiters’ are applied to establish the transformed rank-score range. The reference limiters represent the range of minimum and maximum reference scores and establishes a range for future equivalent comparisons. Selection should consider a range beyond current reference score minima and maxima that would

Component	Component output	Metric quantification description
Achievement	Achievement output ($O_{A,i}$)	Outcome of interest crude measure
	Achievement reference ($R_{A,j}$)	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

Table 1. Component parameters and how they are quantified for an entity ‘ i ’.

account for a realistic range of future reference score possibilities. The lower achievement limiter is attributed a transformed rank-score of 1 and the upper achievement limiter is attributed a transformed rank-score of 100 (Fig. 2A). In different program settings either a higher or lower achievement score may reflect the ‘best’ performance. This directionality is established during the stakeholder selections at the start of the statistical R-code data handling process.

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

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

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

where

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

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

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

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

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

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

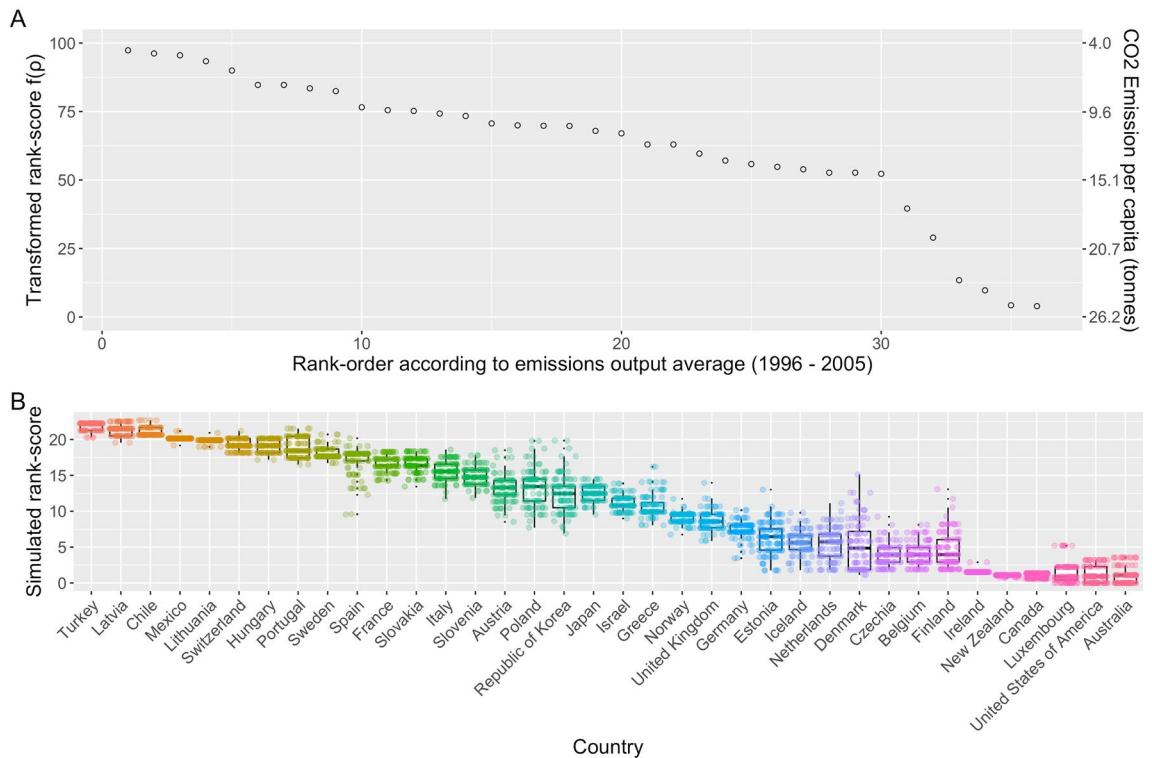


Figure 2. (A) Entity rankings, (B) Simulated rank scores (example data). (A) Initial rank-order (ranked by average annual CO₂ emissions per capita 1996–2005) (ρ_i) versus average CO₂ emission per capita 1996–2005 (tonnes) and transformed rank-score $f(\rho)$. (B) Simulation (example with 100 iterations shown for illustrative purposes) of potential rank score outputs (Φ_i) for each entity based on underlying distribution of achievement output. Box plot showing 16th and 84th percentiles.

Step 3: component indicators (normalised component outputs). *Action.* Component outputs transformed to normalised measures reflecting the magnitude change in achievement and rank.

Outcome. Z-achievement and z-rank.

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

For the achievement component:

$$z_{A,i} = \frac{\text{signum} \cdot (O_{A,i} - R_{A,i})}{D_{A,i}} \tag{5}$$

where

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

For the rank component:

$$z_{RS,i} = \frac{O_{RS,i} - R_{RS,i}}{D_{RS,i}} \tag{6}$$

Step 4: composite score. *Action.* Aggregated component indicators with optional weighting.

Outcome. Performance outcome score (POS).

Procedure. The component indicators are combined in a weighting procedure based upon a user pre-defined setup of their respective relevance. An argument for the choice of weights for each respective component is formulated. The POS is established together with an explanation of the component weighting. When applying the option of proportional weighting to the component indicators the larger the importance of the component, the

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

Table 2. Contextual features framing the evaluation.

higher the weighting factor: w_A for achievement and w_R for rank. The weighting is performed on centred z-scores for achievement $c_{A,i}$ and rank $c_{R,i}$ respectively:

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

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

where

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

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

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

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

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

Output goals:

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

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

Outcome goal:

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

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

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

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

Table 3. Components, parameters, and metrics used to populate the component outputs.

Component	Component output	Sweden	Mexico	Norway	Lux
Achievement	Achievement output ($O_{A,i}$)	5.47	5.74	10.47	18.17
	Achievement reference ($R_{A,i}$)	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 to 2.27)

Table 4. Component outputs for the four stakeholders derived from Step 2 of the POS-T. (Units of achievement = tonnes of CO₂ equivalent per capita. Rank score output = proportional score gained from ranking ahead of other countries. Rank score reference = median score from simulation based on achievement descriptive statistics. Rank score deviation = simulation outputs based on 16th and 84th percentiles. LUX = Luxembourg).

Example step 2: component outputs (performance metrics from each component). Data were managed following the stepwise process detailed in the R-code software provided (Data S1). Data for each component were collected⁴⁴ and presented in file format (Data S2 and S3). Lower and upper reference limiters were set beyond the minimum and maximum CO₂ emissions per capita outputs. Individual entity standard deviations were chosen for use in the simulations. Application of the POS-T R-code to the performance data produced automated performance outputs and descriptive statistics. The following component outputs were generated (Fig. 2B and Table 4).

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

Component output: Three of the four countries reduced their raw CO₂ emissions per capita. Luxembourg saw the largest reduction of the four example countries and largest reduction compared to the full OECD cohort (–5.90 tonnes per capita), followed by Sweden (13th overall; –2.43 tonnes per capita), and Norway (21st overall; –1.76 tonnes per capita). Mexico saw an increase in CO₂ emissions (30th overall; +0.26 tonnes per capita).

Goal 2. “To evaluate change in international ranking with respect to CO₂ emissions per capita over the period 2006–2015”.

Component output: Three of the four countries improved their ranking. Sweden (= 3rd largest rank shift overall; +8 places), Luxembourg and Mexico (= 10th largest rank shift overall; +1 place) relative to all 36 comparison countries. Norway fell in ranking (= 27th largest rank shift overall; –3 places) reflecting having not reduced their emissions per capita to the same level over the observation period as the comparison cohort.

Example step 3: component indicators (normalised component outputs). The component outputs were transformed to normalised measures demonstrating internal and external magnitude change relative to the achievement and rank reference standards respectively (Table 5). The three countries that saw a reduction in raw CO₂ emissions per capita demonstrated positive internal magnitude change (z-achievement: Sweden: +6.85, Norway: +12.92, Luxembourg: +2.12). One country saw an increase in raw CO₂ emissions per capita demonstrating negative internal magnitude change (z-achievement: Mexico: –2.56). Three of four countries improved their ranking demonstrating positive external magnitude change (z-rank: Sweden: +8.69, Mexico: +0.92, and Luxembourg: +1.00). One country regressed in ranking demonstrating a negative external magnitude change (Norway: –27.73).

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

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

Example step 4: composite score (POS). The selected weighting ratio for achievement and rank was set at 1:1. The relative magnitude of change (component indicators) for both component outputs were combined, resulting in a composite score (POS) (Table 6).

Outcome goal:

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

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

Summary information gained for performance evaluation (4 stakeholder countries highlighted). See Table 6.

Discussion

This study set out to develop a scoring template that combines internal and external measures of performance by alignment of measurement scales which represent meaningful magnitudes of change. The resulting POS-T adheres to the principle of providing the stakeholders governing applied programs a means to report performance outcomes with convincing face-validity^{39,45,46}. We exemplified application of POS-T in an evaluation of CO₂ emission reduction amongst OECD member countries. Flexible and transparent evaluation methods oriented towards stakeholders and usefulness have repeatedly been asked for in the environmental sciences^{47,48}. To ensure that POS-T produces scores useful for stakeholders, an inductive (discovery) approach was found best suited^{49,50}. This approach aligns with the principles of design thinking⁵¹ where the emphasis is placed on defining the problem to be solved through the needs of the stakeholders involved^{45,46}. In the following, the main features of the POS-T are discussed considering the CO₂ emission example and directions are outlined for future research.

The measures and methods available to report performance results delimit a stakeholder’s capacity to evaluate applied programs⁴⁰. In their governance, accomplishment has been described as the gap between expected and actual output or the deviation of the output from an industry standard¹². The POS-T provides in Step 2 methodology to establish a *magnitude* of this gap for both the achievement and rank components, i.e. for internal and external comparisons. In step 3, ready-to-combine component indicators are formed by normalisation of the component outputs produced. The resulting component indicators describe magnitudes of change, i.e. the gaps between expected and actual results or change over time for internal and external comparisons. Producing the relative magnitude of change from a point of reference for both components (achievement and rank) rather than binary measures alone provides those evaluating the outcome with greater context. For example, on a binary scale, Sweden, Luxembourg, and Norway each demonstrated a *reduction* in CO₂ emissions per capita (positive achievement outcome) relative to their own reference standard in 2005. However, when accounting for the external context, Norway during the evaluation period slipped down the ranking table from 21st to 24th (negative rank outcome) due to that other countries reduced their relative emissions by greater amounts. Conversely, Mexico gained a ranking place from 4th to 3rd (positive rank outcome) even having had a small increase in crude CO₂ emissions due to that other countries close in rank had relatively larger increases in crude CO₂ emissions. Proportional output measures relative to self-comparison and comparison with others in a chosen cohort provides a context for stakeholders to better frame and evaluate an outcome against expectation and describe the overall accomplishment.

Integration of self-comparison and rank change magnitudes adds complexity to program evaluation indicators. Maintenance of face-validity in such composite scores requires measurement system transparency⁵². The POS-T supports transparency and face-validity by offering evaluators semantic clarity regarding the components of the integrated composite score. Stakeholders evaluating performance using the POS-T will base their assessments on normalised indicators of any measured or pre-existing method of reporting achievement output for internal (within-individual) comparisons. The crude achievement outputs can in any circumstance be ranked⁵³ and the normalised indicators computed for rank changes and external comparisons (between individuals). The normalised indicators are calculated in a standardised manner, i.e. for ‘achievement’ by subtracting the mean from an individual raw score and then dividing the difference by the standard deviation. The mean and standard deviation are based on individual, or population standards as chosen in step 1 by the stakeholder evaluating the performance. Normalised indicators are calculated for ‘rank’ by transforming the ordinal scale to continuous relative values before applying the same normalisation process to a series of simulated rank outputs. Such normalised indicators are well-known and are broadly used in, for instance, global health settings for comparative evaluations of development processes at individual and population levels, e.g. in the child growth area^{54,55}. In the example application of POS-T on CO₂ emissions, the normalised ‘self-comparison’ indicator showcases

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	<i>10.19</i>
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	<i>3.98</i>
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	<i>1.60</i>
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	<i>-4.99</i>

Table 6. OECD countries ordered highest to lowest by the POS including component crude output variation and component magnitude of change over the observation period. (Crude CO₂ emissions = tonnes of CO₂ equivalent per capita). Significant values are in bold and italics.

that Luxembourg's crude emissions reduction from the reference of 2005 is achieved in the context of a broader distribution of the annual fluctuation in emissions by Luxembourg compared to Sweden. In essence, Sweden's emissions reduction achievement is at face value more substantial in the context of self-comparison due to the narrower distribution of annual emissions fluctuations. The magnitude of this achievement is demonstrated by a higher achievement indicator. Regarding external comparisons, both Luxembourg and Mexico gained one place on the ranking table, yet the normalised rank indicator calculated in step 3 shows that the magnitude of Luxembourg's gain is greater than Mexico's. This is due to that Mexico's emissions per capita were very close in volume to those countries with similar rank, the effect being that a change in rank may occur even from small changes in emissions output resulting in a smaller rank indicator for the same crude rank change. By using the continuous data that formulates the ranking order, context and magnitude is apportioned to the component indicator in step 3 representing the rank change. Internal and external measures presented as a magnitude of

change against a reference and accompanying distribution provide meaningful context to the performance. The selection of the lower and upper reference limiters provides an important step when applying the statistical code in establishing consistent comparators and context for future equivalent evaluations. Performance outcomes presented this way can be used to observe longitudinal performance trends within an individual entity or relative to a population as well as measuring a single performance against expectation.

Some limitations to the use of the POS-T are important to consider. In experimental evaluation research, influence from external factors is controlled in the study design. Emulating an experimental design in observational performance evaluations in practice settings would require information on all confounding factors⁵⁶. Application of the POS-T does not per se assure that the POS reflects causal effects of the program, and consideration of confounding factors is always needed when interpreting POS scores in practice settings. Moreover, it should be taken into regard that the simulation process used in POS-T to determine the rank score deviations for each entity uses the reference and its standard deviation as assumptions in the calculation. The outcome from each simulation may thus vary slightly. This effect is minimised by always running an adequate number of simulations on each occasion. Furthermore, when the data available to calculate the achievement deviation is limited, a decision must be made regarding what to use as the achievement deviation for comparison with the achievement output. The preferred option is to use the achievement deviation unique to each entity. However, an option is to use the cohort population standard deviation as this broadens the dataset to calculate deviations and improves its reliability. This may be a satisfactory solution when the comparative data sets between entities in the cohort have similar deviations. If this is not the case, an option may be to use the largest or smallest deviation in the cohort. The flexibility in selecting components, standardising parameters and weighting of the component indicators opens the composite score to variability in its robustness. Testing for robustness is recommended and aided by the level of transparency described by the evaluator in selecting optional features in corresponding steps of the POS-T framework.

The POS-T in its current form can be applied to any program governance setting. A POS can be determined for single entities at multiple time points to assess performance trends or for multiple entities at a single time point to assess performances relative to a population standard or to peers. The component indicators can also be used to evaluate each component in isolation. The rank simulation process in isolation may furthermore be utilised to determine probabilities of performance outcome. In a sports context, where both personal achievement and comparative rank are considered in a performance evaluation, the POS-T may provide valid comparison of performance outcomes by individual athletes across the span of a career. In this example the evaluator has flexibility in selecting the achievement reference, e.g. a population standard, or the athletes own season averages. The POS-T can also provide a point in time comparison between individual competitors within an event or funded individual sports programs within a country. In an education context the POS-T may be applied to the performance evaluation of students across semesters or to evaluate the performance of education institutions over time wherever a ranking score is calculated and ranking table produced. The POS-T can be applied to any evaluation setting where a continuous data metric could be used to rank entities. Further development of the POS-T will include development of the statistical code to include the evaluation of performance in settings where achievement is not readily quantified, e.g. when it mainly is established through head-to-head contests.

Conclusion

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

Materials and methods

Construction principles. Construction of the POS-T employed an iterative approach to solution design that prioritises application of the final template in real-world program governance settings. Its practical use was further supported through the parallel construction of a statistical framework and software for score development in applications⁵⁰. A design panel was composed for the construction consisting of scientists and practitioners (n = 5) with backgrounds in epidemiology, public health and sport settings, organisational development, statistical methods, and experimental design. A composite score development model was used to guide the construction process (Fig. 3).

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

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

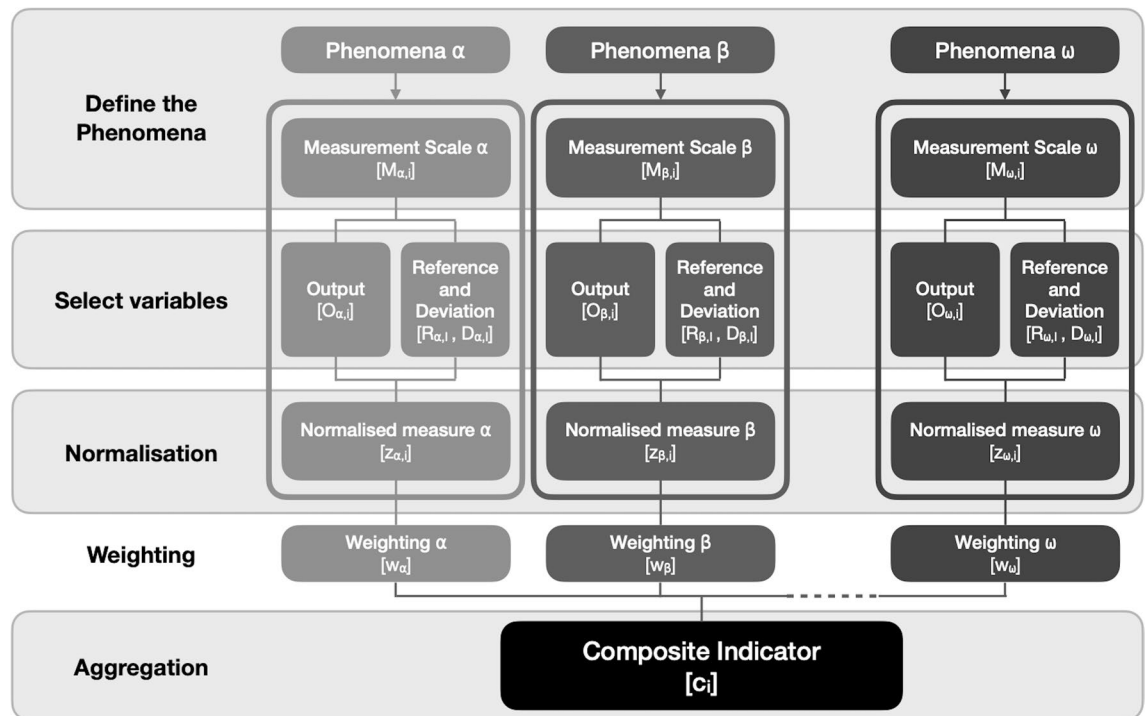


Figure 3. Composite score development model with five development steps and three example phenomena.

POS-T application example. The resulting POS-T was finally applied in an evaluation of reduction of greenhouse gas emission among 36 OECD countries. The OECD has collected and reported annual greenhouse gas emissions between 1996 and 2015 for a cohort of 36 countries producing an annual emissions ranking table⁴³. The OECD data set was used to showcase performance evaluations based on a composite score for comparison within either of the singular reporting metrics; tonnes of CO₂ emissions equivalent per capita, or the emissions ranking table. To showcase the template and the statistical framework, the design panel took on the virtual role of an international stakeholder commission in the environmental protection area. The final statistical framework was exemplified by applying the R-code to the 36 country OECD data set and analysing the performance outcomes of four countries: Sweden, Luxembourg, Mexico, and Norway.

Data availability

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

Received: 12 September 2021; Accepted: 21 February 2022

Published online: 15 March 2022

References

1. *Performance Evaluation Methods: Measurement and Attribution of Program Results* (Treasury Board of Canada, Secretariat, Toronto, 1998).
2. DeGroff, A., Schooley, M., Chapel, T. & Poister, T. H. Challenges and strategies in applying performance measurement to federal public health programs. *Eval. Program Plan.* **33**, 365–372 (2010).
3. *Performance Measurement and Evaluation: Definitions and Relationships* (U.S. Government Accountability Office, GAO-11-646SP, Washington DC, 2011).
4. de Lancer Julnes, P. Performance measurement: An effective tool for government accountability? The debate goes on. *Evaluation* **12**, 219–235 (2006).
5. AbuJbara, N. K. & Worley, J. A. Performance measurement indicators in the healthcare industry: A systematic review. *Glob. Bus. Econ. Rev.* **21**, 43–68 (2019).
6. Lebas, M. Performance measurement and performance management. *Int. J. Prod. Econ.* **41**, 23–35 (1995).
7. Newcomer, K. E. Using performance measurement to improve programs. *New Dir. Eval.* **1997**, 5–14 (1997).
8. T. Rantala. *Operational Level Performance Measurement in University-Industry Collaboration*. Thesis, LUT University, Finland (2019).
9. Milstein, B., Wetterhall, S. & Group, C. E. W. A framework featuring steps and standards for program evaluation. *Health Promot. Pract.* **1**, 221–228 (2000).
10. Koplan, J. P., Milstein, R. & Wetterhall, S. Framework for program evaluation in public health. *MMWR Recomm. Rep.* **48**, 1–40 (1999).
11. Suter, L. G. *et al.* American college of rheumatology white paper on performance outcome measures in rheumatology. *Arthritis Care Res.* **68**, 1390–1401 (2016).
12. W. D. Savedoff, *Governance in the Health Sector: A Strategy for Measuring Determinants and Performance* (Policy Research working paper; no. WPS 5655, The World Bank, 2011).

13. Laurian, L. *et al.* Evaluating the outcomes of plans: Theory, practice, and methodology. *Environ. Plan. B Plan. Des.* **37**, 740–757 (2010).
14. Neely, A., Gregory, M. & Platts, K. Performance measurement system design: A literature review and research agenda. *Int. J. Oper. Prod. Manag.* **15**, 80–116 (1995).
15. Raysmith, B. P., Jacobsson, J., Drew, M. K. & Timpka, T. What is performance? A scoping review of performance outcomes as study endpoints in athletics. *Sports* **7**, 66 (2019).
16. World Health Organization. *The World Health Report 2000: Health Systems: Improving Performance* (World Health Organization, 2000).
17. Dill, D. D. & Soo, M. Academic quality, league tables, and public policy: A cross-national analysis of university ranking systems. *High. Educ.* **49**, 495–533 (2005).
18. Derrien, F. & Dessaint, O. The effects of investment bank rankings: Evidence from M&A league tables. *Rev. Finance* **22**, 1375–1411 (2018).
19. Armesto, S. G., Lapetra, M. L. G., Wei, L. & Kelley, E. *Health Care Quality Indicators Project 2006 Data Collection Update Report. OECD Health Working Papers, No. 29* (OECD Publishing, Paris, 2007).
20. Huang, M.-J., Chen, M.-Y. & Yieh, K. Comparing with your main competitor: The single most important task of knowledge management performance measurement. *J. Inf. Sci.* **33**, 416–434 (2007).
21. Mehrpouya, A. & Samiolo, R. Performance measurement in global governance: Ranking and the politics of variability. *Account. Organ. Soc.* **55**, 12–31 (2016).
22. Oliver, T. R. Peer reviewed: Population health rankings as policy indicators and performance measures. *Prev. Chronic Dis.* **7**, A101 (2010).
23. Keasey, K., Moon, P. & Duxbury, D. Performance measurement and the use of league tables: Some experimental evidence of dysfunctional consequences. *Account. Bus. Res.* **30**, 275–286 (2000).
24. Huang, M.-H. A comparison of three major academic rankings for world universities: From a research evaluation perspective. *J. Libr. Inf. Stud.* **9**, 1–25 (2011).
25. Schütte, S., Acevedo, P. N. M. & Flahault, A. Health systems around the world—a comparison of existing health system rankings. *J. Glob. Health* **8**, 010407 (2018).
26. Hurst, J. & Jee-Hughes, M. *Performance Measurement and Performance Management in OECD Health Systems* (OECD Publishing, 2001).
27. Kaplan, R. S. Strategic performance measurement and management in nonprofit organizations. *Nonprofit Manag. Leadersh.* **11**, 353–370 (2001).
28. Kaplan, R. & Norton, D. The balanced scorecard—Measures that drive performance. *Harv. Bus. Rev.* **83**, 71–79 (1992).
29. Fahlén, J. The trust–mistrust dynamic in the public governance of sport: Exploring the legitimacy of performance measurement systems through end-users’ perceptions. *Int. J. Sport Policy Polit.* **9**, 707–722 (2017).
30. O’Boyle, I. & Hassan, D. Performance management and measurement in national-level non-profit sport organisations. *Eur. Sport Manag. Q.* **14**, 299–314 (2014).
31. Pham, H., Sutton, B. G., Brown, P. J. & Brown, D. A. Moving towards sustainability: A theoretical design of environmental performance measurement systems. *J. Clean. Prod.* **269**, 122273 (2020).
32. N. James, M. Menzies, Global and regional changes in carbon dioxide emissions. Pp. 1970–2019. arXiv preprint <https://arxiv.org/abs/2201.13075> (2022).
33. Carter, N., Day, P. & Klein, R. *How Organisations Measure Success: the Use of Performance Indicators in Government* (Psychology Press, 1995).
34. Jacobs, R., Goddard, M. & Smith, P. C. Composite performance measures in the public sector. *Cent. Health Econ. Res. Pap.* **16**, 1–8 (2007).
35. Davies, I. C. Evaluation and performance management in government. *Evaluation* **5**, 150–159 (1999).
36. *Beyond Roe-How to Measure Bank Performance. Appendix to the report on EU banking structures* (European Central Bank, Frankfurt, 2010).
37. White, L. J. *The Credit Rating Agencies and Their Role in the Financial System* (Oxford University Press, 2018).
38. Seeley, T. *Cumminuty Based Organization (CBO) Survey Results: Outcome Evaluation in Voluntary and Not-for-Profit Organizations* (The Muttard Fellowship, 2003).
39. Guyadeen, D. & Seasons, M. Evaluation theory and practice: Comparing program evaluation and evaluation in planning. *J. Plan. Educ. Res.* **38**, 98–110 (2018).
40. Peterson, E. D. *et al.* ACCF/AHA 2010 position statement on composite measures for healthcare performance assessment: American College of Cardiology Foundation/American Heart Association Task Force on performance measures (writing committee to develop a position statement on composite measures). *J. Am. Coll. Cardiol.* **55**, 1755–1766 (2010).
41. Jacobs, R., Goddard, M. & Smith, P. C. How robust are hospital ranks based on composite performance measures?. *Med. Care* **43**, 1177–1184 (2005).
42. Nardo, M. *et al.* *Handbook on Constructing Composite Indicators: Methodology and User Guide* (OECD Publishing, 2005).
43. Song, M.-K., Lin, F.-C., Ward, S. E. & Fine, J. P. Composite variables: When and how. *Nurs. Res.* **62**, 45 (2013).
44. *Greenhouse Gas Emissions by Source* Internet. (OECD Publishing, 2014). Available from: <https://www.oecd-ilibrary.org/content/data/data-00594-en>.
45. A.-M. R. McGowan, C. Bakula, & R. S. Castner, Lessons learned from applying design thinking in a NASA rapid design study in aeronautics. In *58th AIAA/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference*, 0976 (2017).
46. M. Greene, *Systems Design Thinking: Identification and Measurement of Attitudes for Systems Engineering, Systems Thinking, and Design Thinking*. Thesis, University of Michigan (2019).
47. Schwanitz, V. J. Evaluating integrated assessment models of global climate change. *Environ. Model. Softw.* **50**, 120–131 (2013).
48. Haasnoot, M. *et al.* Fit for purpose? Building and evaluating a fast, integrated model for exploring water policy pathways. *Environ. Model. Softw.* **60**, 99–120 (2014).
49. Gould, J. D. & Lewis, C. Designing for usability: Key principles and what designers think. *Commun. ACM* **28**, 300–311 (1985).
50. Dorst, K. The core of ‘design thinking’ and its application. *Des. Stud.* **32**, 521–532 (2011).
51. Brown, T. Design thinking. *Harv. Bus. Rev.* **86**, 84 (2008).
52. Rae, A. & Wong, C. Monitoring spatial planning policies: Towards an analytical, adaptive, and spatial approach to a ‘wicked problem’. *Environ. Plan. B Plan. Des.* **39**, 880–896 (2012).
53. Mohsin, M. *et al.* Developing low carbon economies: An aggregated composite index based on carbon emissions. *Sustain. Energy Technol. Assess.* **35**, 365–374 (2019).
54. Wit, J. M., Himes, J. H., Van Buuren, S., Denno, D. M. & Suchdev, P. S. Practical application of linear growth measurements in clinical research in low-and middle-income countries. *Horm. Res. Paediatr.* **88**, 79–90 (2017).
55. Leung, M. *et al.* Metrics of early childhood growth in recent epidemiological research: A scoping review. *PLoS ONE* **13**, e0194565 (2018).
56. Hernán, M. A. & Robins, J. M. Using big data to emulate a target trial when a randomized trial is not available. *Am. J. Epidemiol.* **183**, 758–764 (2016).

Author contributions

B.P.R. provided conceptual basis and primary author, O.D. developed the statistical code and refined conceptual thinking to practical outputs, T.T. guided the projects and provided major edits to manuscript, M.K.D. and J.J. contributed to conceptual development and manuscript edits.

Funding

Open access funding provided by Linköping University. The authors acknowledge that they received no funding in support for this research.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-08368-w>.

Correspondence and requests for materials should be addressed to B.P.R.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022