**Bayesian Multilevel Compositional Data Analysis with the R Package *multilevelcoda***

Flora Le[1], Dorothea Dumuid[2], Tyman E. Stanford[2], and Joshua F. Wiley[1]

[1]School of Psychological Sciences,

Monash University

[2]Alliance for Research in Exercise, Nutrition and Activity,

Allied Health and Human Performance,

University of South Australia

**Author Note**

Reproducible materials for this study are available at:

https://github.com/florale/multilevelcoda-overview.

Correspondence to Flora Le (flora.le@monash.edu) or Joshua F. Wiley (joshua.wiley@monash.edu), School of Psychological Sciences, Monash University, Clayton, VIC, Australia.

# Abstract

Multilevel compositional data, such as data sampled over time that are non-negative and sum to a constant value, are common in various fields. However, there is currently no software specifically built to model compositional data in a multilevel framework. The **R** package *multilevelcoda* implements a collection of tools for modelling compositional data in a Bayesian multivariate, multilevel pipeline. The user-friendly setup only requires the data, model formula, and minimal specification of the analysis. This paper outlines the statistical theory underlying the Bayesian compositional multilevel modelling approach and details the implementation of the functions available in *multilevelcoda*, using an example dataset of compositional daily sleep-wake behaviours. This innovative method can be used to gain robust answers to scientific questions using the increasingly available multilevel compositional data from intensive, longitudinal studies.

*Keywords:* compositional data analysis, multilevel model, Bayesian inference, R

**Bayesian Multilevel Compositional Data Analysis with the R Package *multilevelcoda***

Bayesian approaches have been increasingly employed for multilevel models. Motivations for using Bayesian approaches have been covered extensively in other work, including flexibility to specify complex models (Levy and McNeish, 2023) like non-normal random-effect models, robustness to small sample sizes (Le et al., 2024; Stegmueller, 2013), benefits from incorporating existing empirical information (i.e., priors; van de Schoot et al., 2018), and the ease of quantifying uncertainty around arbitrary calculated quantities using posterior samples (Gelman et al., 2013; Wagenmakers et al., 2016). Bayesian sampling algorithms, including the Markov chain Monte Carlo (MCMC) sampling, or Hamiltonian Monte Carlo (HMC) (Betancourt, 2017; Betancourt et al., 2014) and its extension, the No-U-Turn Sampler (NUTS) (Hoffman and Gelman, 2014). MCMC allows for both model parameter estimation and drawing samples from the posterior predictive distribution, which can be used to assess model fit. Importantly, the individual posterior distribution samples can be used to make inferences about the parameters, such as calculating the mean, standard errors and credible intervals, enabling post-hoc analyses involving any calculated quantities to be directly and intuitively conducted from Bayesian multilevel models. Such features can greatly facilitate the analysis of data with complex structure, such as multilevel compositional data.

Multilevel compositional data exist in various fields, such as time-use epidemiology (e.g., time spent in different sleep-wake behaviours during the 24-hour day), and sleep (e.g., proportion of time spent in different sleep stages during the night), and nutritional epidemiology (e.g., macronutrients like proteins, fats and carbohydrates as proportions of total caloric intake). These data can be classified as compositions, which consist of *compositional parts* that contain relative information about the whole; represented as non-negative values that sum to a constant value. Compositional parts can be expressed as percentages (or proportions) of the composition but also may be in other units that are constrained to a constant total value (e.g., 1440 minutes in a day). They are commonly measured across multiple time points (e.g., across several consecutive days), or nested within clusters (e.g., schools). This means the data are multilevel, with the most

common multilevel data structure having two levels (e.g., consecutive days nested within people). Thus, these data often consist of two sources of variability at each level: between (differences between clusters, such as people) and within (differences within clusters, commonly the deviation of a specific value from the average of that cluster).

Despite the abundance of multilevel compositional data, standard statistical methods, including multilevel models, do not produce valid results for raw compositional data. This is due to the perfect multicollinearity present in their constant-sum nature (i.e., compositional parts are linearly dependent). Instead, compositional data analysis (CoDA; Aitchison, 1986; Pawlowsky-Glahn and Buccianti, 2011) utilises the relative information contained in compositional data using log-ratios. Certain log-ratio transformations can remove the linear dependence of compositional parts, while retaining the relative nature of the compositional components (i.e., changes in compositional components but not their total). One such transform is the isometric log-ratio tranform (ilr), that is the most commonly used in the physical activity and sedentary behaviour research. The ilr transformation eliminates the multicollinearity by producing one less ilr coordinate compared to the number of compositional parts (e.g. two ilr coordinates are calculated from a 3-part composition), thus allowing standard statistical methods to be applied on the transformed data. Some software for compositional data analysis exist, mostly in **R**, including ***compositions*** (Van den Boogaart and Tolosana-Delgado, 2008, 2013), ***compositions*** (Palarea-Albaladejo and Martín-Fernández, 2015), ***robCompositions*** (Templ et al., 2011), ***Compositional*** (Tsagris et al., 2023), ***codaredistlm*** (Dumuid et al., 2018; Stanford et al., 2022). These packages offer general tools for manipulating or modelling compositional data so they can be used in standard statistical models. However, they do not easily accommodate the manipulation of compositional data with a multilevel structure or provide functions to fit multilevel models with compositional data (or their corresponding log-ratios) as response or predictor variables.

Further, to facilitate the interpretability of CoDA, isotemporal compositional substitution analysis (Dumuid et al., 2019) is a post-hoc approach to examine the model predicted changes in

an outcome associated with changes to the compositional parts. Substitution analysis provides an opportunity to interrogate the model to answer questions about the predicted change in an outcome when the compositional parts are redistributed. This analysis can answer questions such as whether there is a change in health when people spend time being physically active at the expense of sitting. Increasing evidence from isotemporal compositional substitution analysis shows that reallocating time between behaviours are associated with both physical and mental heath outcomes (Janssen et al., 2020; Grgic et al., 2018; Miatke et al., 2023). Most existing studies are, however, cross-sectional, with less evidence available from longitudinal data. This may be due to the lack of tools to efficiently work with longitudinal compositional data in substitution analysis. To our knowledge, no software currently automates isotemporal compositional substitution analysis, especially in a multilevel framework.

The ***multilevelcoda*** package (Le and Wiley, 2024), presented in this article, aims to address these gaps with tools to automate estimating multilevel models for compositional data (and their associated log-ratios) and isotemporal compositional substitution analysis in a Bayesian framework. Specifically, ***multilevelcoda*** advances the analysis of multilevel compositional data by offering three important contributions:

- Compute multilevel compositions and perform log-ratio transformation. Decompose data into between and within levels if necessary.

- Automatically fit Bayesian multilevel models with compositional predictors and/or outcomes. The Bayesian models are implemented using the ***brms*** package, which supports a variety of generalised (non-)linear multivariate multilevel models.

- Estimate isotemporal compositional substitution analysis for Bayesian multilevel models at both between and within levels. We leverage the posterior draws of Bayesian models to derive reliable estimates of credible intervals of the predicted differences in the expected outcome for the reallocation of compositional parts.

We begin by introducing the fundamental concepts of multilevel composition and the underlying

multilevel models for compositional data. We then describe the functionality of the software using an example of daily 24-hour sleep-wake behaviours. Lastly, we provide a comparison across packages for compositional data analysis and discuss plans for extending the package.

## Multilevel Composition Data

### Properties of Compositional Data

A *composition* is defined as a vector of $D$ positive components, called *compositional parts*, that sum to a constant $\kappa$

$$\boldsymbol{x} = (x_1, x_2, \ldots, x_D), \tag{1}$$

where $\sum_{i=1}^{D} x_i = \kappa$ and $x_i > 0 \ \forall i = 1, 2, \ldots, D$. Consequently, compositions are elements in the $D$-simplex, denoted as $\mathscr{S}^D \subset \mathbb{R}^D$, whose parts are linearly dependent as each part can be deduced with knowledge of the $D-1$ other parts (e.g., $x_j = \kappa - \sum_{\forall i \neq j} x_i$). An important operation on the simplex is perturbation (Aitchison, 1986), or closure operation applied to the element-wise product. Perturbation is the analogous operation to addition in the Euclidean space (Van den Boogaart and Tolosana-Delgado, 2013; Aitchison, 1986), defined as

$$\boldsymbol{x} \oplus \boldsymbol{x^*} = \mathscr{C}\left(x_1 \cdot x_1^*, x_2 \cdot x_2^*, \ldots, x_D \cdot x_D^*\right) \tag{2}$$

where

$$\mathscr{C}(\boldsymbol{x}) = \frac{\kappa}{\sum_{i=1}^{D} x_i} \boldsymbol{x}$$

is the closure operation that normalises the compositional parts of a vector $\boldsymbol{x}$ to the constant $\kappa$ (Aitchison, 1986), and $\boldsymbol{x}, \boldsymbol{x^*} \subset \mathscr{S}^D$. Such properties are incompatible with many standard mathematical operations (e.g., $\mathscr{S}^D$ is not closed under addition) and statistical methods that assume independence (e.g., multiple linear regression) that are developed in the real space $\mathbb{R}^{D-1}$ (for detailed discussion on the properties of compositional data and their consequences, see Aitchison, 1994, 1986).

**Log-ratio Approach for Multilevel Compositional Data Analysis**

Two modelling strategies can be used for compositional data: a) directly working on the simplex, and b) transforming compositional data from the simplex to the real space then modelling the transformed data using a multivariate normal distribution, which is referred to as principle of working in coordinates (Mateu-Figueras et al., 2011). The most widely studied distribution on the simplex is the Dirichlet, which can be used to model composition directly (Gueorguieva et al., 2008). However, its use in applications is quite limited, as the strong independence structure of Dirichlet distribution (i.e., independent, equally scaled gamma-distributed variables) poorly models the dependence between compositional components (Aitchison, 1982).

The alternative approach to modelling compositional data, originating from Aitchison, 1982, focuses on the relative magnitudes and variations of components, rather than their absolute values. The prototype of a distribution in the simplex is the logistic-normal distribution (additive log-ratio transformation; Atchison and Shen, 1980; Aitchison, 1982), which is also referred to as normal distribution on the simplex. Using log-ratios, a composition in the simplex ($\mathscr{S}^D$) can be expressed in terms of ratios of the components of the composition in the Euclidean space ($\mathbb{R}^{D-1}$) where standard mathematical operations and statistical methods are valid. A family of log-ratio transformations for modelling compositional data includes additive log-ratio (alr; Aitchison, 1982), centered log-ratio (clr; Aitchison, 1982), and isometric log-ratio (ilr; Egozcue et al., 2003). Non-orthonomal transformations (e.g., alr, clr, or simple log-ratio transformations) have limitations in modelling compositional data (Mateu-Figueras et al., 2011). Specifically, the alr transformation is not isometric, thus, does not preserve the properties of compositional data (angles and distances). The clr transformation preserves the distance but does not break the sum constraint, which results in a singular covariance matrix. In contrast, the ilr-transform maps the $D$-part compositional data from the simplex to non-overlapping subgroups in the $(D-1)$-dimension Euclidean space isometrically by using an orthonormal basis, thereby preserving the compositional properties and yielding a full-rank covariance matrix. We describe

ilr transform in detail in the following.

Consider a composition $x \in \mathscr{S}^D$ and a corresponding set of $(D-1)$ ilr coordinates $(z_1, z_2, \ldots, z_{D-1}) = z \in \mathbb{R}^{D-1}$. The individual $z_k$ coordinate is constructed as normalised log-ratio of the geometric mean of compositional parts in the numerator (a mutually exclusive set of subcompositions denoted as $R_k$) to the geometric mean of compositional parts in the denominator (a set of subcompositions denoted as $S_k$). The $k^{\text{th}}$ $(k = 1, 2, \ldots, D-1)$ ilr coordinate can then be written as

$$z_k = \sqrt{\frac{r_k s_k}{r_k + s_k}} \ln \left( \frac{\tilde{x}_{R_k}}{\tilde{x}_{S_k}} \right), \quad k = 1, 2, \ldots, D-1 \tag{3}$$

where

$$\tilde{x}_{R_k} = \left( \prod_{x_d \in R_k} x_d \right)^{\frac{1}{r_k}} \quad \text{and} \quad \tilde{x}_{S_k} = \left( \prod_{x_d \in S_k} x_d \right)^{\frac{1}{s_k}}$$

with $r_k$ and $s_k$ being the size of the sets $R_k$ and $S_k$, respectively, and $\sqrt{\frac{r_k s_k}{r_k + s_k}}$ being a normalising constant.

One method for ilr transformation employs a sequential binary partition (SBP) process (Egozcue and Pawlowsky-Glahn, 2005), which produces ilr coordinates that are interpretable depending on the application. A SBP is obtained by first partitioning the compositional parts into two non-empty sets, where one set corresponds to the first ilr coordinate's numerator and the other set corresponds to the first ilr coordinate's denominator. Using the same principle, each of the previously constructed sets are recursively partitioned into two non-empty sets until no further non-empty partitions of the subcompositional parts are possible (after $D-1$ steps). This SBP process can be coded via a $D \times (D-1)$ matrix corresponding to the $D$ compositional parts and their membership in the $(D-1)$ ilr coordinates; +1 if the compositional part is the ilr numerator, -1 if the compositional part is the ilr denominator, or 0 if the compositional part is uninvolved in the ilr coordinate. The ilr coordinates can be interpreted as the log ratio of the subcomposition in the numerator in relation to the subcomposition in the denominator.

SBP can be constructed to form conceptually meaningful contrasts (e.g., time spent in

sleeping behaviours all relative to waking behaviours). In some cases, it is not straightforward to correctly interpret ilr coordinates, as they are expressed in terms of the log-ratios of groups of parts. Another choice of SBP for ilr transformation involves constructing the ilr coordinate $z_k$ to capture all information of the compositional part $x_k$ relative to the remaining parts of the composition of $\boldsymbol{x}$, termed pivot balance coordinate (Fišerová and Hron, 2011; Hron et al., 2012), is defined as

$$z_k = \sqrt{\frac{D-k}{D-k+1}} \ln \frac{x_k}{\sqrt[D-k]{\prod_{i=k+1}^{D} x_i}}, \quad k = 1, \ldots, D-1. \tag{4}$$

Table 1 gives an example of a complete SBP used to construct pivot coordinates from a five-part composition $\boldsymbol{x}_{ij} = (x_{1ij}, x_{2ij}, x_{3ij}, x_{4ij}, x_{5ij})$. With this specific choice of coordinates, all relative information about the first part $x_{1ij}$ (pivot element) is contained exclusively in the coordinate $z_{1ij}$, but not in the other coordinates. If one were interested in an interpretation about another part, for example $x_{2ij}$, the role of $x_{1ij}$ and $x_{2ij}$ is exchanged by placing $x_{2ij}$ to the first position in the compositional vector, and the same type of coordinate is constructed. The resulting coordinates are, thus, rotations of the original coordinates. In this way, from a $D$-part composition, we can construct $D$ pivot coordinates for the compositional parts of interest, which are all rotations of each other, and where only the first coordinate (pivot balance) is used for an interpretation of the respective part.

Regardless of the choice of SBP, the ilr coordinates are linearly independent multivariate real values (if the compositional parts are strictly positive) and overcome multicollinearity. Therefore, they can be entered in conventional statistical models (e.g., multilevel models), making them tractable and easy to understand. Importantly, the ilr transformation is invertible, such that the ilr coordinates can be back-transformed via their 1 - 1 relationship to the original composition, as required (Egozcue and Pawlowsky-Glahn, 2005).

**Multilevel Compositional Data and Transformations**

Compositional data (e.g., activity, diet) may be measured on multiple people $j = 1, 2, \ldots, J$ at multiple time points $i = 1, 2, \ldots, I$. We denote such data as $\boldsymbol{x}_{ij} = (x_{1ij}, x_{2ij}, \ldots, x_{Dij})$, which is a vector of compositional data observed at the $i^{\text{th}}$ time point for the $j^{\text{th}}$ person. Therefore, $\boldsymbol{x}_{ij}$ can vary between individuals and across time points within an individual, containing both between-person and within-person variability (Curran and Bauer, 2011). We express the $D$-part, time-varying multilevel composition $\boldsymbol{x}_{ij}$ as

$$
\begin{aligned}
\boldsymbol{x}_{ij} &= \left( x_{1ij}, x_{2ij}, \ldots, x_{Dij} \right) \\
&= \mathscr{C} \left( x_{1 \cdot j}^{(b)} \cdot x_{1ij}^{(w)}, x_{2 \cdot j}^{(b)} \cdot x_{2ij}^{(w)}, \ldots, x_{D \cdot j}^{(b)} \cdot x_{Dij}^{(w)} \right) \\
&= \boldsymbol{x}_{\cdot j}^{(b)} \oplus \boldsymbol{x}_{ij}^{(w)}
\end{aligned}
\tag{5}
$$

where

- $x_{d \cdot j}^{(b)}$ is the person-specific mean of the $d^{\text{th}}$ compositional part over time, which contains only between-person variance and no within-person variance. The subscript $\cdot j$ denotes the average across $i$ observations for the individual $j$ and superscript $(b)$ denotes the *between* component of the compositional parts.

- $x_{dij}^{(w)}$ is the time-specific deviation of the $d^{\text{th}}$ compositional part from the person $j$ specific mean (i.e., compositional mean-centered deviate), which has within-person variance and no between-person variance. The superscript $(w)$ denotes the *within* component of the composition parts.

- $\mathscr{C}$ is the closure operation, and

- $\oplus$ is the perturbation operation on the simplex.

The between- and within-person subcompositions can also be expressed as compositions themselves as

$$
\begin{aligned}
\boldsymbol{x}_{\cdot j}^{(b)} &= \mathscr{C}\left(x_{1 \cdot j}^{(b)}, x_{2 \cdot j}^{(b)}, \ldots, x_{D \cdot j}^{(b)}\right) \text{ and} \\
\boldsymbol{x}_{ij}^{(w)} &= \mathscr{C}\left(x_{1ij}^{(w)}, x_{2ij}^{(w)}, \ldots, x_{Dij}^{(w)}\right)
\end{aligned}
\tag{6}
$$

The ilr transformed coordinates $\boldsymbol{z}_{ij} \in \mathbb{R}^{D-1}$ corresponding to the composition $\boldsymbol{x}_{ij} \in \mathscr{S}^D$ can also be uniquely (with respect to the specific ilr transformation) decomposed into its between- and within-person components in a more familiar additive way

$$
\begin{aligned}
\boldsymbol{z}_{ij} &= \left(z_{1ij}, z_{2ij}, \ldots, z_{(D-1)ij}\right) \\
&= \left(z_{1 \cdot j}^{(b)} + z_{1ij}^{(w)}, z_{2 \cdot j}^{(b)} + z_{2ij}^{(w)}, \ldots, z_{(D-1) \cdot j}^{(b)} + z_{(D-1)ij}^{(w)}\right) \\
&= \boldsymbol{z}_{\cdot j}^{(b)} + \boldsymbol{z}_{ij}^{(w)}
\end{aligned}
\tag{7}
$$

in which $z_{kij}$ is the value of the $k^{\text{th}}$ ($k = 1, 2, \ldots, D-1$) ilr coordinate at time point $i$ for individual $j$ and superscripts $(b)$ and $(w)$ denote the between and within components, respectively, of the ilr coordinates. Although we focused on longitudinal data (repeated measures are nested within person) here, the same principles can also be used to distinguish within- and between-person effects in hierarchical data (individuals are nested within groups). In any applications, $i$ index the elementary "level 1" units and $j$ index the clusters or "level 2" units. It should be noted that the separation of within-person and between-person effects only works to two-level data structure, wherein between-person level is person-mean at level 2, and within-person level is the mean-centered deviate at level 1.

Disaggregating effects for more-than-two-level models are not currently supported in *multilevelcoda*. Methods research has generally not explored how to dissaggregate multilevel models beyond two levels. For now, we recommend keeping the data at the aggregate level (that is, not separated by between and within-person effects), while considering appropriate interpretation (see Curran and Bauer, 2011, for a discussion on between-person and within-person inferences).

## Model Description

As we adopt Bayesian inference from a pragmatic perspective, our exposition of it is kept to a minimum. Readers interested in further methodological guidance on Bayesian analyses are referred to Kruschke, 2014; McElreath, 2018 for introductions and Gelman et al., 2013; Bürkner, 2018 for more advanced usage. In the following section, multilevel models with compositional predictors and their associated post-hoc substitution analyses are first described, followed by multilevel models with compositional responses.

### Multilevel Models with Compositional Predictors

To express a linear model for the time-varying $D$-part multilevel compositional predictor, we first denote the outcome variable observed at time point $i$ for individual $j$ as $y_{ij}$. The prediction of a continuous, normally distributed outcome $y_{ij}$ is the linear combination of the between-person and within-person effects of a $D$-part composition (expressed as a set of $(D-1)$-dimension ilr coordinates). A linear multilevel model of $y_{ij}$ can be written as

$$y_{ij} = \beta_{0j} + \overbrace{\sum_{k=1}^{D-1} \beta_k z_{k \cdot j}^{(b)}}^{\text{between}} + \underbrace{\sum_{k=1}^{D-1} \beta_{(k+D-1),j}\, z_{kij}^{(w)}}_{\text{within}} + \varepsilon_{ij} \tag{8}$$

where

$$
\begin{bmatrix} \beta_1 \\ \vdots \\ \beta_{(D-1)} \end{bmatrix} = \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_{(D-1)} \end{bmatrix}
$$

$$
\begin{bmatrix} \beta_{0j} \\ \beta_{Dj} \\ \vdots \\ \beta_{2(D-1)j} \end{bmatrix} = \begin{bmatrix} \gamma_0 \\ \gamma_D \\ \vdots \\ \gamma_{2(D-1)} \end{bmatrix} + \begin{bmatrix} u_{0j} \\ u_{1j} \\ \vdots \\ u_{(D-1)j} \end{bmatrix}
$$

$$
\begin{bmatrix} u_{0j} \\ u_{1j} \\ \vdots \\ u_{(D-1)j} \end{bmatrix} \sim \text{MVNormal}(\mathbf{0}, \mathbf{\Sigma_u})
$$

$$
\varepsilon_{ij} \sim \text{Normal}(0, \sigma_\varepsilon^2)
$$

The $\gamma$s are the population-level effects, $u$s are group-level effects, and $\mathbf{\Sigma_u}$ is a variance-covariance matrix for the group-level effects. The between- and within-person components of the composition (expressed as a set of ilr coordinates) are $z_{\cdot j}^{(b)}$ and $z_{ij}^{(w)}$, with the subscripts denoting that the between component is unique to individual $j$ and the within component is unique to time $i$ for individual $j$. Thus, all $z_{\cdot j}^{(b)}$ and $z_{ij}^{(w)}$ can be included as population-level effects ($\gamma$), and only $z_{ij}^{(w)}$ can be included as group-level effects ($u$). The between- and within-person effects of the $k^{\text{th}}$ ($k = 1, 2, \ldots, D-1$) ilr coordinates are $\beta_k$ and $\beta_{k+D-1}$, respectively. Because each ilr coordinate is decomposed into its between- and within-person components, the total number of $\beta$ parameters for the ilr coordinates is twice the number of the ilr coordinates. Further (time varying) population- and/or group-level covariates are not included here but can easily be incorporated.

**Multilevel Compositional Substitution Analysis**

When examining the relationships between compositional predictors and an outcome, we often are interested in the expected difference in the outcome when a fixed amount of the composition is reallocated from one compositional component to another, while the other

components remain constant. These changes can be examined using isotemporal compositional substitution analysis. In the following, we describe this model in the multilevel framework.

### *Prediction for A Given (Reference) Composition*

For a $D$-part composition for person $j$ at time $i$, $\boldsymbol{x}_{ij} = \boldsymbol{x}_{\cdot j}^{(b)} \oplus \boldsymbol{x}_{ij}^{(w)}$, and the corresponding set of ilr coordinates $\boldsymbol{z}_{ij} = \boldsymbol{z}_{\cdot j}^{(b)} + \boldsymbol{z}_{ij}^{(w)}$, the predicted $y_{ij}$ is

$$\hat{y}_{ij} = \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k \cdot j}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{kij}^{(w)} \tag{9}$$

Now consider the reallocation of a given amount from one part of the composition, denoted $d$, to another part, denoted $d'$, where $d' \neq d \in \{1, \dots, D\}$. This is only possible with reference to a starting composition. The starting composition where compositional components are reallocated from/to is referred to as the reference composition (commonly the *compositional mean*, although any reference composition could be used). The decomposition of a reference composition $\boldsymbol{x}_0$ is

$$\begin{aligned}
\boldsymbol{x}_0 &= \boldsymbol{x}_0^{(b)} \oplus \boldsymbol{x}_0^{(w)} \\
&= \mathscr{C}\left(x_{10}^{(b)} \cdot x_{10}^{(w)}, \dots, x_{d0}^{(b)} \cdot x_{d0}^{(w)}, \dots, x_{d'0}^{(b)} \cdot x_{d'0}^{(w)}, \dots, x_{D0}^{(b)} \cdot x_{D0}^{(w)}\right)
\end{aligned} \tag{10}$$

Note when the reference composition is a compositional mean value at the between-person level, the within-person subcomposition $\boldsymbol{x}_0^{(w)}$ becomes the neutral element of the simplex, $\boldsymbol{1}_D = \mathscr{C}(1, 1, \dots, 1) = (\kappa/D, \kappa/D, \dots, \kappa/D)$ as there is no within-person deviation. In such cases, the reference composition and its corresponding ilr transformation can be simplified to

$$\boldsymbol{x}_0 = \boldsymbol{x}_0^{(b)} \oplus \boldsymbol{1}_D = \boldsymbol{x}_0^{(b)}$$

The predicted outcome at a reference composition $\boldsymbol{x}_0$ is

$$\hat{y}_0 = \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)} \tag{11}$$

where $z_{k0}^{(b)}$ and $z_{k0}^{(w)}$ are the between- and within-person ilr coordinates at the reference composition, respectively.

Given that multilevel composition contains both between- and within-person variability, we can investigate the changes in the outcome associated with the reallocation of compositional parts at between- and within-person levels. There are important distinctions between the two approaches. A between-person substitution examines the differences in the outcome between individuals with different mean compositions, whereas a within-person substitution examines the differences in the outcome associated with the changes in the composition within an individual (i.e., the deviations from their own mean composition).

### *Between-person Substitution*

We denote the two compositional parts involved in a given between-person pairwise substitution as $x_{d0}^{(b)}$ and $x_{d'0}^{(b)}$. The reallocation of a fixed amount $t$ from $x_{d0}^{(b)}$ to $x_{d'0}^{(b)}$ (that is, adding $t$ to $x_{d'}^{(b)}$ and subtracting $t$ from $x_{d0}^{(b)}$ simultaneously) around a reference composition $\boldsymbol{x}_0$ at the between-person level is

$$
\begin{aligned}
x_d^{(b)'} &= x_{d0}^{(b)} - t \\
x_{d'}^{(b)'} &= x_{d'0}^{(b)} + t
\end{aligned}
\tag{12}
$$

where $d' \neq d \in \{1, \ldots, D\}$, $t$ is the reallocated change (e.g., minutes/1440 if $\kappa = 1440$), and $0 < t < \min\left\{x_d^{(b)}, \kappa - x_{d'}^{(b)}\right\}$. Keeping the remaining parts of the composition constant, the new $D$-part composition $\boldsymbol{x}_{(d-d')}^{(b)'}$ can be expressed as

$$
\begin{aligned}
\boldsymbol{x}_{(d-d')}^{(b)'} &= \mathscr{C}(x_{10}^{(b)} \cdot x_{10}^{(w)}, \ldots, x_d^{(b)'} \cdot x_{d0}^{(w)}, \ldots, x_{d'}^{(b)'} \cdot x_{d'0}^{(w)}, \ldots, x_{D0}^{(b)} \cdot x_{D0}^{(w)}) \\
&= \mathscr{C}(x_{10}^{(b)} \cdot x_{10}^{(w)}, \ldots, (x_{d0}^{(b)} - t) \cdot x_{d0}^{(w)}, \ldots, (x_{d'0}^{(b)} + t) \cdot x_{d'0}^{(w)}, \ldots, x_{D0}^{(b)} \cdot x_{D0}^{(w)})
\end{aligned}
\tag{13}
$$

The predicted outcome at the between-person reallocation is given as

$$\hat{y}_{(d-d')}^{(b)'} = \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)'} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)} \tag{14}$$

where $z_{k0}^{(b)'}$ indicates the new between-person ilr coordinates resulted from the between-person reallocation in the composition $z_{k0}^{(w)}$ (within-person ilr coordinates) is the same as the reference ilr coordinates. The predicted difference in the outcome, $\Delta\hat{y}_{(d-d')}^{(b)}$, for the between-person changes in compositional parts (i.e., between the reference composition and the reallocated composition at between-person level) is therefore

$$\begin{aligned} \Delta\hat{y}_{(d-d')}^{(b)} =\ & \hat{y}_{(d-d')}^{(b)'} - \hat{y}_0 \\ =\ & \left( \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)'} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)} \right) \\ & - \left( \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)} \right) \\ =\ & \sum_{k=1}^{D-1} \hat{\beta}_k \left( z_{k0}^{(b)'} - z_{k0}^{(b)} \right) \end{aligned} \tag{15}$$

### *Within-person Substitution*

The reallocation of a fixed amount $t$ between two compositional parts at the within-person level ($x_{d0}^{(w)}$ and $x_{d'0}^{(w)}$) around a reference composition can be expressed as

$$\begin{aligned} x_d^{(w)'} &= x_{d0}^{(w)} - t \\ x_{d'}^{(w)'} &= x_{d'0}^{(w)} + t. \end{aligned} \tag{16}$$

The new $D$-part composition for within-person level reallocation of $t$ becomes

$$\begin{aligned} \boldsymbol{x}_{(d-d')}^{(w)'} &= \mathscr{C}(x_{10}^{(b)} \cdot x_{10}^{(w)}, \ldots, x_{d0}^{(b)} \cdot x_d^{(w)'}, \ldots, x_{d'0}^{(b)} \cdot x_{d'}^{(w)'}, \ldots, x_{D0}^{(b)} \cdot x_{D0}^{(w)}) \\ &= \mathscr{C}(x_{10}^{(b)} \cdot x_{10}^{(w)}, \ldots, x_{d0}^{(b)} \cdot (x_{d0}^{(w)} - t), \ldots, x_{d'0}^{(b)} \cdot (x_{d'}^{(w)} + t), \ldots, x_{D0}^{(b)} \cdot x_{D0}^{(w)}). \end{aligned} \tag{17}$$

The predicted outcome for the within-person reallocation is

$$\hat{y}^{(w)'}_{(d-d')} = \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)'} \tag{18}$$

where the $z_{k0}^{(b)}$ remains the same as the reference between-person ilr coordinates, whereas the $z_{k0}^{(w)'}$ is the new within-person ilr coordinates, showing the change in within-person ilr coordinates relative to the reference point. Thus, the predicted changes in the outcome due to the changes across the compositional parts at the within-person level, $\Delta\hat{y}^{(w)}_{(d-d')}$, is

$$
\begin{aligned}
\Delta\hat{y}^{(w)}_{(d-d')} =\ & \hat{y}^{(w)'}_{(d-d')} - \hat{y}_0 \\
=\ & \left( \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)'} \right) \\
& - \left( \hat{\beta}_{0j} + \sum_{k=1}^{D-1} \hat{\beta}_k z_{k0}^{(b)} + \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} z_{k0}^{(w)} \right) \\
=\ & \sum_{k=1}^{D-1} \hat{\beta}_{(k+D-1),j} \left( z_{k0}^{(w)'} - z_{k0}^{(w)} \right).
\end{aligned}
\tag{19}
$$

### Substitution Analysis Framework

We propose two frameworks for the substitution analysis (Table 2), with noteworthy distinctions. The **Simple substitution analysis** provides simple effects of the change in a composition on an outcome, where the reference composition could be grand compositional mean or any (constant) hypothetical set of values. The *Average substitution analysis* is motivated by average marginal effects (Norton et al., 2019; Mize et al., 2019). That is, using the cluster (e.g., person) compositional mean as the reference composition to estimate the predicted changes in the outcome for each cluster, then averaging across the prediction to obtain the average change of the sample. This estimate reflects the change in outcome when every cluster (e.g., person) in the sample reallocates a *t* unit from one compositional part to another, which demonstrates the change for the full distribution of the predictor(s) rather than an arbitrary prediction (Leeper, 2017).

For linear outcomes, the results produced by the two models are expected to be

comparable. *Average substitution analysis* provides better estimates than *Simple substitution analysis* particularly in the cases of non-linear outcomes, models with covariates, large reallocation across compositional parts resulting in one part approaching zero, or imbalanced data, such as the unequal balance of time spent in sleep-wake behaviours across individuals (e.g., shift workers vs non-shift workers, male vs females). In addition, the credible intervals estimated by the *Simple substitution analysis* only reflect the population level effects, whereas the credible intervals estimated by the average substitution analysis incorporate the variability at the group-level by including all group-level effects.

Average substitution analysis* generally require more time and computational resources than *Simple substitution analysis*, as the estimation takes place at the cluster-level. However, all `substitution()` analyses can be executed in parallel using available **R** packages, such as ***doFuture*** (Bengtsson, 2023), to optimise computational time and performance.

## Package Overview

Package ***multilevelcoda*** provides functions for fitting multivariate multilevel models with compositional data using full Bayesian inference. The package is open-source software for the **R** programming platform. The latest release version of ***multilevelcoda*** from the Comprehensive **R** Archive Network (CRAN) can be installed via `install.packages("multilevelcoda")`. Alternatively, the current developmental version can be downloaded from GitHub via

```
R> devtools::install_github(
+ "florale/multilevelcoda")
```

***multilevelcoda*** uses the **R** package ***brms*** to build models, which in turn uses the probabilistic programming language **Stan** as the backend that dynamically generates and compiles **C++** code for specific, Bayesian models. Thus, a **C++** compiler is required, beyond just having a functional **R** installation. For Windows, the program ***Rtools*** (R Core Team, 2022) comes with a **C++** compiler. On Mac, Installation of **Xcode** (Apple Inc, 2022) for Mac, is required. Linux requires **g++** or **Clang**. Detailed instructions on how to get the compilers and running can be found in the prerequisites section on the RStan package's website. Note that the *rstan* package (the **R** interface of **Stan**; Stan Development Team, 2020) also depends heavily on several other R

packages; these dependencies are automatically installed if the *rstan* package (R interface to

Stan) is installed via one of the conventional mechanisms. Users will find further assistance

through **R** documentation and vignettes to guide them through the functionality of the package, as

well as examples of its use. Contributions are welcome both in terms of bug reports and feature

enhancements, via the standard mechanism of GitHub issues and pull requests.

Analysis in *multilevelcoda* follows the procedure in Figure 1. First, the user calculates the

composition and the corresponding log-ratio transforms (e.g., ilr coordinates) using `complr()`

function. Next, this information is used to fit Bayesian (multivariate) multilevel models using the

`brmcoda()` function. During this step, the model is passed to `brm()` to generate **Stan** model

code, which is then passed to either the *rstan* package (Stan Development Team, 2020) or the

*cmdstanr* package (the **R** interface of *CmdStan*; Stan Development Team, 2022). Models are

compiled in **C++**, fitted by **Stan**, and post-processed in *brms* before being saved in

*multilevelcoda*'s `brmcoda()` in **R**. The results from `brmcoda` can be used to estimate the pivot

coordinates of the composition using the `pivot_coord()` function, and substitution analysis

using the `substitution()` function. Finally, results from all functions can be investigated in **R**

using various methods such as `summary()`, `plot()`, or `predict()` (for example, for a complete

list of methods defined on the `brmcoda` object, type `methods(class = "brmcoda")`).

### Example Application

This section presents examples to implement Bayesian multilevel compositional data

analysis following workflow in Figure 1. Reproducible code can be found at

https://github.com/florale/multilevelcoda-overview. The example data set is a simulated, built-in

data set in a long format, with repeated measurements of stress and 24h time use separated into

five behaviours: total sleep time, time awake in bed, moderate-to-vigorous physical activity

(MVPA), light physical activity (LPA), and sedentary behaviour (SB). Daily stress was measured

on a 0-10 scale. The five behaviours make up a 5-part composition.

```
R> library(multilevelcoda)
R> data(mcompd)
R> data(psub)
```

The example data set `mcompd` is in long format and consists of 3540 entries of 10 variables

```
R> head(mcompd)
ID Stress   TST WAKE MVPA  LPA   SB Time Age Female
185      4   542   99  297  460   41    1  30      0
185      7   458   49  117  653  162    2  30      0
185      3   271   41  489  625   15    3  30      0
185      2   525   76  259  398  182    4  30      0
185      8   651   86  112  436  155    5  30      0
185      8   431   84  264  476  185    6  30      0
```

Variable `ID` is participant id. `Stres` is self-reported stress measured on a 0-10 scale. Sleep duration (i.e., total sleep time, `TST`), time awake in bed (`WAKE`), moderate-to-vigorous physical activity (`MVPA`), light physical activity (`LPA`), and sedentary behaviour (`SB`) make up a 5-part composition. `Time` is time point id at which stress and the 5-part composition were repeatedly measured. Finally, variables `Age` and `Female` are individual baseline factors.

**Transforming Multilevel Compositional Data**

The `complr()` function processes compositional data and performs log-ratio transformation. First, to build a set of ilr coordinates, a SBP is required. The construction and interpretation of *ilr* coordinates may depend on specific application. Alternatively, multilevelcoda uses pivot coordinates as the default SBP, where the first pivot coordinate represent the ratio of the first compositional part relative to the remaining parts. We may use the following code to process and transform compositional data:

```
R> cilr <- complr(
+   data  = mcompd,
+   parts = c("TST", "WAKE", "MVPA", "LPA", "SB"),
+   idvar = "ID",
+   total = 1440
+ )
```

We specify the variable that identifies how units are clustered `idvar` as `ID` based on our data, and specify value of `total` to be 1440, which is the total minutes of a 24-hour day, to which the 5 behaviours must sum. In this example (when a SBP is not specified), the default SBP is

```
R> sbp
       TST WAKE MVPA LPA SB
[1,]     1   -1   -1  -1 -1
[2,]     0    1   -1  -1 -1
[3,]     0    0    1  -1 -1
[4,]     0    0    0   1 -1
```

where each column represents one of the 5 parts of the composition and each row represents one of the transformed 4 (5-1) ilr coordinates. Here, the first coordinate represent the ratio of sleep relative to the remaining behaviours. Intepretation of the coordinates based on this SBP is in Table 3. A summary the transformed data can be obtained by running

```
R> summary(cilr)
composition_parts   TST, WAKE, MVPA, LPA, SB
logratios              ilr1, ilr2, ilr3, ilr4
idvar                                     NULL
nobs                                      3540
ngrps                                      266
transform_type                             ilr
total                                     1440
composition_geometry                     acomp
logratio_class                           rmult
```

The output provides information about the data and transformation. Some general information on the data include the names of compositional parts (`composition_parts`) and log-ratio variables (`logratios`), the ID variable (`idvar`, for multilevel dataset), number of observations (`nobs`) and number of groups (`ngrps`). Other information to perform transformation on compositional data includes transformation methods (either `ilr`, `alr`, or `clr`), the closure value (for ilr transformation), the geometries and classes of the composition and the ilr coordinates. The class `acomp` indicates the composition class that aligns with the philosophical framework of the Aitchison Simplex, whereas ilr are real multivariate vectors. Within a `complr` object (not shown), data sets of composition and ilr coordinates are stored alongside the original data set, which are used for subsequent analyses.

**Fitting Bayesian Multilevel Models with Compositional Predictors**

Multilevel models for compositional data are estimated using `brmcoda()`. To examine how between-person and within-person 24h behaviours predict stress, we may fit the following model:

```
R> m <- brmcoda(
+    complr  = cilr,
+    formula = Stress ~
+              bilr1 + bilr2 + bilr3 + bilr4 +
+              wilr1 + wilr2 + wilr3 + wilr4 +
+              (1 | ID),
+    warmup = 1000, iter = 2000, seed = 123,
```

```
+    chains = 4, cores = 4, backend = "cmdstanr"
+ )
```

The structure of `brmcoda()` has two core arguments: the `complr` object, which replaces the standard `data`, argument for models, and a model `formula`. The `formula` argument takes information on the outcomes and predictors of the model, separated by the $\sim$. Models are fitted using ***brms***, therefore, the syntax follows the form of a `brmsfit` object and is similar to ***lme4***'s. The left side of the `formula` is the outcome, `Stress` in this example. The right side of the `formula` specifies the predictors, including both population-level and group-level terms, separated by the $+$. In the present example, the population-level terms are 4 `bilrs` representing the set of ilr coordinates at between-person level and 4 `wilrs` representing the ilr coordinates at within-person level. The group-level terms follow the form `(coef | group)`, allowing the intercept to vary by `ID`. Note that a `data` argument is not required, as the data set is supplied by the `complr` object. Additional arguments used for `brm` model function are specified in the . . . argument, such as prior specifications and distribution of the response variable. If not otherwise specified, default `brm` priors and link functions are applied. This example model `m` is fitted using 4 chains, each with 2000 iterations including 1000-warmup iterations for the sampler, running on 4 cores. The model produces 4000 posterior draws using the HMC sampler. Weakly-informative priors were used (see Supplementary materials for prior information), which play a minimal role in the computation of the posterior distribution, and maximise the influence of the data. Student's t distribution was used for the fixed intercept, and flat priors (improper priors over the reals) were used for the fixed parameters of the predictors. The standard deviation parameters of the random intercept and residual were specified using student's t distributions.

### *Model Summary*

The output of `brmcoda()` is a **R** `brmcoda` object with 2 elements: an fitted `brm()` model with class `brmsfit` and the input data from `complr()`. A model summary is available via

```
R> summary(m)
 Family: gaussian
  Links: mu = identity; sigma = identity
Formula: Stress ~ bilr1 + bilr2 + bilr3 + bilr4 +
         wilr1 + wilr2 + wilr3 + wilr4 + (1 | ID)
```

```
   Data: tmp (Number of observations: 3540)
  Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
         total post-warmup draws = 4000

Multilevel Hyperparameters:
~ID (Number of levels: 266)
              Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
sd(Intercept)     1.00      0.06     0.88     1.13 1.00     1573     2899

Regression Coefficients:
          Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
Intercept     2.59      0.48     1.59     3.51 1.00     1463     2242
bilr1         0.39      0.43    -0.43     1.26 1.00     1530     2041
bilr2        -0.10      0.17    -0.46     0.23 1.00     1344     1829
bilr3         0.11      0.21    -0.31     0.54 1.00     1469     2333
bilr4        -0.01      0.28    -0.56     0.55 1.00     1463     1907
wilr1        -0.16      0.16    -0.48     0.15 1.00     4646     2869
wilr2        -0.30      0.08    -0.47    -0.14 1.00     6135     3094
wilr3        -0.10      0.08    -0.25     0.06 1.00     3741     2891
wilr4         0.24      0.10     0.04     0.43 1.00     3951     3377

Further Distributional Parameters:
      Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
sigma     2.38      0.03     2.33     2.44 1.00     6030     2794

Draws were sampled using sample(hmc). For each parameter, Bulk_ESS
and Tail_ESS are effective sample size measures, and Rhat is the potential
scale reduction factor on split chains (at convergence, Rhat = 1).
```

The model output follows the standard output from `brm()`. The top of the output shows the general information of the model, followed by the group-level (random) effects and population-level (fixed) effects. At the bottom of the output, family specific parameters and autocorrelation (if incorporated) are also provided. Every parameter is summarised using the mean (`Estimate`), standard deviation (`Est.Error`) of the posterior distribution (the standard error of the estimate), and two-sided 95% Credible intervals based on the quantiles (`l-95% CI` and `u-95% CI`) (Bürkner, 2018). Additional information about the model were also provided, including `Rhat` for information on the convergence of the algorithm and `Bulk_ESS` and `Tail_ESS` for effective sample size (ESS).

Model convergence can be evaluated using diagnostic statistic $\hat{R} < 1.05$ (Vehtari et al., 2021) and ESS $> 400$ (Vehtari et al., 2021). Examining the population-level effects of the ilr coordinates, only `wilr2` and `wilr4` have the two-sided 95% Credible Intervals not containing

zero. Therefore, we have evidence for the association between `wilr2` and `wilr4` and `Stress`, repsectively. Recall that the interpretation of the ilr depends on the SBP. For example, the significant coefficient for `wilr2` shows that the increase time awake in bed while proportionally decreasing waking behaviours (MVPA, LPA, and SB) on a given day, predicted lower stress (-0.30 [95%CI -0.47, -0.14]).

### *Estimating and Interpreting Pivot Coordinate Coefficients*

Pivot coordinates represent the relative importance of each behaviour in the 24h composition (with respect to the geometric average of the remaining behaviours). Pivot coordinates may be an easier interpretation as they are always contrasting one part of the composition to all remaining parts. Using them also makes the results independent of any one SBP specified, because pivot coordinates can be calculated as a rotation of any given SBP. Pivot coordinates can be obtained by running

```
R> m_coordinates <- pivot_coord(m, method = "rotate")
+   summary(m_coordinates)
```

We supplied the `pivot_coord()` function with a `brmcoda` object, and specified `method = "rotate"` to indicate we want to rotate all possible ilr basis matrix to estimate the pivot coordinates representing each 24h behaviour. The results showing the association between the pivot coordinates representing the 24h behaviours and stress are in Table 3. Results showed that higher time spent in awake in bed relative to the remaining behaviours was associated with -0.25 [95%CI-0.42, -0.08] lower stress, whereas higher time spent in LPA relative to the remaining behaviours predicted 0.37 [0.12, 0.63]) higher stress.

### Running Multilevel Compositional Substitution Analysis

Beyond understanding the independent and compositional association between behaviours and stress, the changes in stress for different pairwise reallocation of behaviours (e.g., reallocation between MVPA and SB while keeping the remaining fixed) can be estimated using compositional substitution analysis. The below example shows how to conduct a *Simple substitution analysis* (by automating the steps described in Table 2) to examine the changes in stress associated with

behaviour reallocation for 1 to 10 minutes, at between-person and within-person levels using the

using the `substitution()` function. We use the below code

```
R> sub_simple <- substitution(
+    object = m,
+    delta  = 1:10,
+    ref    = "grandmean",
+    level  = c("between", "within")
+ )
```

`substitution()` requires a `brmcoda` object. `delta = 1:10` indicates the estimation of the

changes in the outcome `Stress` for the reallocation from 1 to 10 minutes across behaviours. We

also specify `ref = "grandmean"` to indicate *simple substitution analysis*. If desired, `ref` can

also take a reference grid that contains the combination of predictors (i.e., reference composition

and other covariates) over which predictions are made. If an user's specified reference grid is not

supplied, the default reference grid (imported from `emmeans` package; Lenth, 2023) consisting of

average value of numeric predictors and the levels of the categorical predictors is used. The

default 95% credible interval are used here, however, can be any desired intervals, such as `ci =`

`0.99`. As we are interested in both between- and within-person changes, we specified `level =`

`c("between", "within")`.

Results of the Bayesian compositional substitution analyses are in Table 2. At

between-person level, none of the results are significant, as the 95% credible intervals contain 0,

showing that reallocation between behaviours was not associated in changes in stress. At

within-person level, there were significant results for the substitution of time awake in bed and

total sleep time (TST), respectively, with other behaviours. More 10 minutes in time awake in bed

at the expense of any other behaviours predicted lower stress (estimates range from -0.03 [95%CI

-0.06, -0.00] to -0.04 [95%CI -0.06, -0.02]. The opposite reallocations were also supported, with

reallocation of 10 minutes from time awake in bed to other behaviours was associated with higher

stress (estimates range from 0.04 [95%CI 0.01, 0.07] to 0.05 [95%CI 0.02, 0.07]. Additionally,

less time in TST predicted -0.03 lower stress [95% CI -0.06, -0.00] when compensated by time

awake in bed, but 0.01 higher stress [95% CI 0.00, 0.02] when compensated by LPA.

**Presenting Substitution Analysis Results**

      *multilevelcoda* offers a streamlined way of visualising the results from the `substitution` models, using the `plot()` method that is built on the well-known `ggplot2` package (Wickham, 2016). For example, we can graph the between substitution results of sleep by running

```
R> plot(sub_simple, to = "TST", ref = "grandmean", level = "between")
```

      The estimated differences in stress associated with both the between- and within substitution results of sleep are in Figure 2 (some additional parameters had to be set for the figure to be in the format shown, see supplementary code for details). Figure 2 showed that reallocation from other behaviours to time awake in bed was associated with lower stress level. These associations were only significant at the within, but not between-person levels.

<p align="center"><strong>Comparison Between Packages</strong></p>

      Many existing **R** packages implement general functions for compositional data, however they do not accommodate multilevel data or provide functions to fit multilevel models with compositional variables and perform post hoc analyses. The **R** package *multilevelcoda* stands out by enabling a streamlined workflow for analysing compositional data in a multilevel framework. Features unique to *multilevelcoda* are the calculation of multilevel composition and log-ratios at between- and within-person levels and the capacity to estimate wide range of (multivariate) multilevel models. Other features currently exclusive to *multilevelcoda* is the streamlined estimation of pivot coordinates and the multilevel compositional substitution analysis for different types of variability (between and within-person), as well as types of reference composition (grand mean, cluster mean, and user's specified). Beyond features, another important focus of *multilevelcoda* is on speed. Models using `brmcoda()` are fitted using package *brms*, which generally require more time and computational resources than package *lme4*. Given the complexity of the substitution analysis, `substitution()` supports parallel execution via package *foreach* (Daniel et al., 2022) and *doFuture* (Bengtsson, 2023), which enables models to run faster in shorter walltime. A comparison across packages for working with compositional data is provided in Table 5.

**Future extension**

In this article, we introduced the analysis of multilevel compositional data in a Bayesian framework using the our **R** package *multilevelcoda*. In the current limited landscape for modelling multilevel compositional data, *multilevelcoda* is a contribution that integrates three methods: compositional data analysis, multilevel modelling, and Bayesian inference into one, open-source program. The implementation of *multilevelcoda* enables a streamlined and efficient workflow from dealing with raw multilevel compositional data, estimating models, and presenting final results, making the analysis of multilevel compositional data faster and more accessible. To the best of our knowledge, this is the first statistical package provides tools for estimating multilevel isotemporal compositional substitution analysis at between and within-person levels. As this method provides a unique opportunity to gain novel insights in the fields of epidemiology and psychology, such as the integrated and interactive effects of sleep-wake behaviours on health outcomes, *multilevelcoda* may be particularly useful for intensive longitudinal studies in this landscape. The support for optional parallel execution further promotes an efficient and powerful performance, enabling the estimation of complex models with less walltime.

*multilevelcoda* is under active development. A current priority is to support various outcome families (e.g., multivariate) in the substitution analysis. We also plan on integrating features to deal with missing data, zeros, and outliers from existing packages to enable more streamlined workflow. Functions to estimate the marginal means for Bayesian multivariate multilevel models with compositional outcomes when integrating out group-level effects if desired will also be added. These extensions will be made available on the developmental version on GitHub before releasing on CRAN, along with vignettes to demonstrate new functionality. Interested users are welcome to follow the GitHub version of *multilevelcoda* and provide feedback for future development of the package.

# References

Aitchison, J. (1982). The statistical analysis of compositional data. *Journal of the Royal Statistical Society: Series B (Methodological)*, *44*(2), 139–160.

Aitchison, J. (1986). *The statistical analysis of compositional data*. Chapman; Hall.

Aitchison, J. (1994). Principles of compositional data analysis. *Lecture Notes-Monograph Series*, 73–81.

Apple Inc. (2022). *Xcode software* [Version 14]. Cupertino, USA. https://developer.apple.com/xcode/

Atchison, J., & Shen, S. M. (1980). Logistic-normal distributions: Some properties and uses. *Biometrika*, *67*(2), 261–272.

Bengtsson, H. (2023). *doFuture: Use foreach to parallelize via the future framework* [R package version 1.0.0]. https://CRAN.R-project.org/package=doFuture

Betancourt, M. (2017). A conceptual introduction to hamiltonian monte carlo. *arXiv preprint arXiv:1701.02434*. https://arxiv.org/abs/1701.02434

Betancourt, M., Byrne, S., Livingstone, S., & Girolami, M. (2014). The geometric foundations of hamiltonian monte carlo. *arXiv preprint arXiv:1410.5110*. https://arxiv.org/abs/1410.5110

Bürkner, P.-C. (2018). Advanced Bayesian Multilevel Modeling with the R Package brms. *The R Journal*, *10*(1), 395–411. https://doi.org/10.32614/RJ-2018-017

Curran, P. J., & Bauer, D. J. (2011). The disaggregation of within-person and between-person effects in longitudinal models of change. *Annual Review of Psychology*, *62*, 583. https://doi.org/10.1146/annurev.psych.093008.100356

Daniel, F., Ooi, H., Calaway, R., Microsoft, & Weston, S. (2022). *foreach: Provides foreach looping construct* [R package 1.5.2]. https://CRAN.R-project.org/package=foreach

Dumuid, D., Pedisic, Z., Stanford, T. E., Martín-Fernández, J.-A., Hron, K., Maher, C. A., Lewis, L. K., & Olds, T. (2019). The compositional isotemporal substitution model: A method for estimating changes in a health outcome for reallocation of time between sleep,

physical activity and sedentary behaviour. *Statistical Methods in Medical Research*, *28*(3), 846–857. https://doi.org/10.1177/0962280217737805

Dumuid, D., Stanford, T. E., Martin-Fernández, J.-A., Pedišić, Ž., Maher, C. A., Lewis, L. K., Hron, K., Katzmarzyk, P. T., Chaput, J.-P., Fogelholm, M., et al. (2018). Compositional data analysis for physical activity, sedentary time and sleep research. *Statistical Methods in Medical Research*, *27*(12), 3726–3738. https://doi.org/10.1177/0962280217710835

Egozcue, J. J., & Pawlowsky-Glahn, V. (2005). Groups of parts and their balances in compositional data analysis. *Mathematical Geology*, *37*(7), 795–828. https://doi.org/10.1007/s11004-005-7381-9

Egozcue, J. J., Pawlowsky-Glahn, V., Mateu-Figueras, G., & Barcelo-Vidal, C. (2003). Isometric logratio transformations for compositional data analysis. *Mathematical Geology*, *35*(3), 279–300. https://doi.org/10.1023/A:1023818214614

Fišerová, E., & Hron, K. (2011). On the interpretation of orthonormal coordinates for compositional data. *Mathematical Geosciences*, *43*, 455–468.

Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian data analysis*. CRC press.

Grgic, J., Dumuid, D., Bengoechea, E. G., Shrestha, N., Bauman, A., Olds, T., & Pedisic, Z. (2018). Health outcomes associated with reallocations of time between sleep, sedentary behaviour, and physical activity: A systematic scoping review of isotemporal substitution studies. *International Journal of Behavioral Nutrition and Physical Activity*, *15*(1), 1–68.

Gueorguieva, R., Rosenheck, R., & Zelterman, D. (2008). Dirichlet component regression and its applications to psychiatric data. *Computational Statistics & Data Analysis*, *52*(12), 5344–5355.

Hoffman, M. D., & Gelman, A. (2014). The no-u-turn sampler: Adaptively setting path lengths in hamiltonian monte carlo. *Journal of Machine Learning Research*, *15*(1), 1593–1623.

Hron, K., Filzmoser, P., & Thompson, K. (2012). Linear regression with compositional explanatory variables. *Journal of Applied Statistics*, *39*(5), 1115–1128. https://doi.org/10.1080/02664763.2011.644268

Janssen, I., Clarke, A. E., Carson, V., Chaput, J.-P., Giangregorio, L. M., Kho, M. E., Poitras, V. J., Ross, R., Saunders, T. J., Ross-White, A., et al. (2020). A systematic review of compositional data analysis studies examining associations between sleep, sedentary behaviour, and physical activity with health outcomes in adults. *Applied Physiology, Nutrition, and Metabolism*, *45*(10), S248–S257.

Kruschke, J. (2014). *Doing bayesian data analysis: A tutorial with r, jags, and stan.*

Le, F., Stanford, T. E., Dumuid, D., & Wiley, J. F. (2024). Bayesian multilevel compositional data analysis: Introduction, evaluation, and application. *arXiv preprint arXiv:2405.03985*.

Le, F., & Wiley, J. F. (2024). *multilevelcoda: Estimate Bayesian multilevel models for compositional data* [R package version 1.3.0]. https://CRAN.R-project.org/package=multilevelcoda

Leeper, T. J. (2017). Interpreting regression results using average marginal effects with R's margins. *Available at the comprehensive R Archive Network (CRAN)*, 1–32.

Lenth, R. V. (2023). *emmeans: Estimated marginal means, aka least-squares means* [R package version 1.8.5]. https://CRAN.R-project.org/package=emmeans

Levy, R., & McNeish, D. (2023). Perspectives on bayesian inference and their implications for data analysis. *Psychological Methods*, *28*(3), 719.

Mateu-Figueras, G., Pawlowsky-Glahn, V., & Egozcue, J. J. (2011). The principle of working on coordinates, 29–42. https://doi.org/10.1002/9781119976462.ch3

McElreath, R. (2018). *Statistical rethinking: A bayesian course with examples in r and stan.* Chapman; Hall/CRC.

Miatke, A., Olds, T., Maher, C., Fraysse, F., Mellow, M. L., Smith, A. E., Pedisic, Z., Grgic, J., & Dumuid, D. (2023). The association between reallocations of time and health using compositional data analysis: A systematic scoping review with an interactive data

exploration interface. *International Journal of Behavioral Nutrition and Physical Activity*, *20*(1), 127.

Mize, T. D., Doan, L., & Long, J. S. (2019). A general framework for comparing predictions and marginal effects across models. *Sociological Methodology*, *49*(1), 152–189. https://doi.org/10.1177/0081175019852763

Norton, E. C., Dowd, B. E., & Maciejewski, M. L. (2019). Marginal effects—quantifying the effect of changes in risk factors in logistic regression models. *JAMA*, *321*(13), 1304–1305. https://doi.org/10.1001/jama.2019.1954

Palarea-Albaladejo, J., & Martín-Fernández, J. A. (2015). zCompositions — R package for multivariate imputation of left-censored data under a compositional approach. *Chemometrics and Intelligent Laboratory Systems*, *143*, 85–96. https://doi.org/10.1016/j.chemolab.2015.02.019

Pawlowsky-Glahn, V., & Buccianti, A. (2011). *Compositional data analysis: Theory and applications*. John Wiley & Sons. https://doi.org/10.1002/9781119976462

R Core Team. (2022). *RTools: Toolchains for building R and R packages from source on windows* [R package version 4.2]. R Foundation for Statistical Computing. Vienna, Austria. https://cran.r-project.org/bin/windows/Rtools/

Stan Development Team. (2020). *rstan: The R interface to Stan* [R package version 2.21.2]. http://mc-stan.org/

Stan Development Team. (2022). *cmdstanr: The R interface to CmdStan* [R package version 2.30.1]. https://mc-stan.org/cmdstanr/

Stanford, T. E., Rasmussen, C. L., & Dumuid, D. (2022). *codaredistlm: Compositional data linear models with composition redistribution* [R package version 0.1.0]. https://CRAN.R-project.org/package=codaredistlm

Stegmueller, D. (2013). How many countries for multilevel modeling? a comparison of frequentist and bayesian approaches. *American Journal of Political Science*, *57*(3), 748–761.

Templ, M., Hron, K., & Filzmoser, P. (2011). robCompositions: An r-package for robust
       statistical analysis of compositional data. *Compositional data analysis: Theory and
       applications*, 341–355.

Tsagris, M., Athineou, G., Alenazi, A., & Adam, C. (2023). *Compositional: Compositional data
       analysis* [R package version 6.4]. https://CRAN.R-project.org/package=Compositional

Van den Boogaart, K. G., & Tolosana-Delgado, R. (2008). compositions: A unified R package to
       analyze compositional data. *Computers and Geosciences*, *34*(4), 320–338.
       https://doi.org/10.1016/j.cageo.2006.11.017

Van den Boogaart, K. G., & Tolosana-Delgado, R. (2013). *Analyzing compositional data with R*
       (Vol. 122). Springer-Verlag. https://doi.org/10.1007/978-3-642-36809-7

van de Schoot, R., Sijbrandij, M., Depaoli, S., Winter, S. D., Olff, M., & Van Loey, N. E. (2018).
       Bayesian PTSD-trajectory analysis with informed priors based on a systematic literature
       search and expert elicitation. *Multivariate Behavioral Research*, *53*(2), 267–291.

Vehtari, A., Gelman, A., Simpson, D., Carpenter, B., & Bürkner, P.-C. (2021).
       Rank-normalization, folding, and localization: An improved r$\hat{R}$ for assessing convergence
       of mcmc (with discussion). *Bayesian Analysis*, *16*(2), 667–718.
       https://doi.org/https://doi.org/10.48550/arXiv.1903.08008

Wagenmakers, E.-J., Morey, R. D., & Lee, M. D. (2016). Bayesian benefits for the pragmatic
       researcher. *Current Directions in Psychological Science*, *25*(3), 169–176.

Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
       https://ggplot2.tidyverse.org

**Table 1**

*Example Sequential Binary Partition of A D = 5 Compositional Parts to Construct (D -1) = 4Pivot Balance Coordinates.*

| Coordinate | $x_{1ij}$ | $x_{2ij}$ | $x_{3ij}$ | $x_{4ij}$ | $x_{5ij}$ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | +1 | -1 | -1 | -1 | -1 |
| 2 | 0 | +1 | -1 | -1 | -1 |
| 3 | 0 | 0 | +1 | -1 | -1 |
| 4 | 0 | 0 | 0 | +1 | -1 |

**Table 2**

*Substitution analysis framework.*

---

*Simple substitution analysis*

---

Examines the change in the outcome for the *between-person* and *within-person* reallocation of compositional parts, using the **grand compositional mean** or a **user's specified composition** as the reference composition. Estimated change in outcome is the simple effect of the compositional reallocation, and if incorporated, at different levels of the other (categorical) predictors or is an unweighted average of them.

The fitting procedure of the *Simple substitution analysis* consists of **6** main steps:

1. Calculate the reference composition ($x_0$) (e.g., grand compositional mean) and its corresponding ilr coordinates ($z_0$)

2. Estimate the outcome at the reference composition, $\hat{y}_0$.

3. Using the reference composition, generate new composition(s), $x'_0$, and the corresponding ilr coordinates, $z'_0$, for the reallocation(s) of $t$ from one part of the composition to another.

4. Estimate the outcome at the reallocated compositions, $\hat{y}'_0$.

5. Calculate the changes in the outcome for the $t$ reallocation(s), $(\Delta\hat{y})$.

6. Repeat this procedure for all compositional parts (end after $D$ steps).

---

*Average substitution analysis*

---

Examines the average change in the outcome for the *between-person* and *within-person* reallocation of compositional parts, using the **cluster (e.g., individual) compositional mean** as the reference composition. Change in the outcome is estimated for each cluster then averaged over the data to obtain an average change of the compositional reallocation, which gives weighted prediction over the empirical (sample) distribution.

The estimation of *Average substitution analysis* follows **7** main steps:

1. Calculate the reference composition ($x_0$) (i.e., cluster compositional mean) and its corresponding ilr coordinates ($z_0$).

2. Estimate the outcome at the cluster compositional mean, $\hat{y}_0$.

3. Using the cluster compositional mean, generate new composition(s) $x'_0$, and the corresponding ilr coordinates, $z'_0$, for the reallocation(s) of $t$ from one part of the composition to another for each cluster.

4. Estimate the outcome at the reallocated compositions for each cluster, $\hat{y}'_0$.

5. Calculate the changes in the outcome for the $t$ reallocation(s), $(\Delta\hat{y})$, at the cluster level.

6. Average the changes in $\Delta\hat{y}$, that is $\overline{\Delta\hat{y}}$ over the clusters/data.

7. Repeat this procedure for all compositional parts (end after $D$ steps).

---

**Table 3**

*Bayesian Multilevel Model with Compositional Predictor Examining the Associations of the 24-hour Sleep-Wake Behaviours and Stress.*

| Pivot Coordinate | Posterior mean and 95% credible intervals |
|---|---|
| **Between-person level** | |
| Sleep vs remaining | 0.39 $[-0.43, 1.26]$ |
| Awake in bed vs remaining | $-0.20$ $[-0.57, 0.16]$ |
| MVPA vs remaining | 0.04 $[-0.34, 0.41]$ |
| LPA vs remaining | $-0.13$ $[-0.81, 0.58]$ |
| SB vs remaining | $-0.11$ $[-0.52, 0.30]$ |
| **Within-person level** | |
| Sleep vs remaining | $-0.16$ $[-0.48, 0.15]$ |
| Awake in bed vs remaining | $-0.25^*$ $[-0.42, -0.08]$ |
| MVPA vs remaining | 0.05 $[-0.09, 0.19]$ |
| LPA vs remaining | $0.37^*$ $[0.12, 0.63]$ |
| SB vs remaining | $-0.01$ $[-0.15, 0.14]$ |

*Notes.* MVPA = moderate-to-vigorous physical activity, LPA = light physical activity, SB = sedentary behaviour. $^*$95% credible intervals not containing 0.

**Table 4**

*Bayesian Multilevel Compositional Substitution Analysis Estimating the Difference in Stress Associated with Reallocation of 30 minutes across 24-hour Sleep-Wake Behaviours.*

| | ↓ Sleep | ↓ Awake in bed | ↓ MVPA | ↓ LPA | ↓ SB |
|---|---|---|---|---|---|
| **Between-person level** | | | | | |
| ↑ Sleep | - | 0.04 [−0.03, 0.11] | 0.01 [−0.03, 0.04] | 0.01 [−0.01, 0.03] | 0.01 [−0.02, 0.04] |
| ↑ Awake in bed | −0.03 [−0.10, 0.02] | - | −0.03 [−0.08, 0.03] | −0.02 [−0.08, 0.03] | −0.02 [−0.07, 0.02] |
| ↑ MVPA | −0.01 [−0.04, 0.03] | 0.03 [−0.02, 0.09] | - | 0.00 [−0.02, 0.03] | 0.01 [−0.02, 0.03] |
| ↑ LPA | −0.01 [−0.03, 0.01] | 0.03 [−0.03, 0.09] | 0.00 [−0.03, 0.03] | - | 0.00 [−0.02, 0.02] |
| ↑ SB | −0.01 [−0.04, 0.02] | 0.03 [−0.03, 0.09] | −0.01 [−0.03, 0.02] | 0.00 [−0.02, 0.02] | - |
| **Within-person level** | | | | | |
| ↑ Sleep | - | 0.04* [0.01, 0.07] | −0.01 [−0.02, 0.01] | −0.01 [−0.02, 0.00] | 0.00 [−0.01, 0.01] |
| ↑ Awake in bed | −0.03* [−0.06, −0.00] | - | −0.04* [−0.06, −0.00] | −0.04* [−0.06, −0.02] | −0.03* [−0.06, −0.01] |
| ↑ MVPA | 0.01 [−0.01, 0.02] | 0.04* [0.01, 0.07] | - | 0.00 [−0.01, 0.01] | 0.00 [−0.01, 0.01] |
| ↑ LPA | 0.01* [0.00, 0.02] | 0.05* [0.02, 0.07] | 0.00 [−0.01, 0.01] | - | 0.00 [−0.00, 0.01] |
| ↑ SB | 0.00 [−0.01, 0.01] | 0.04* [0.01, 0.07] | 0.00 [−0.01, 0.01] | −0.01 [−0.01, 0.00] | - |

*Notes.* MVPA = moderate-to-vigorous physical activity, LPA = light physical activity, SB = sedentary behaviour. Values are posterior means and 95% credible intervals. *95% credible intervals not containing 0.

**Table 5**

*Comparison across packages for compositional data. Notes. ilr = isometric log-ratio, alr = additive log-ratio, clr = centered log-ratio. *Models with compositional outcomes can include compositional predictors. †Only available for Bayesian models.*

| | multilevelcoda | compositions | Compositional | compositions | robCompositions |
|---|---|---|---|---|---|
| **Basic functions for composition** | | | | | |
| Composition | Yes | Yes | Yes | Yes | No |
| Multilevel composition | Yes | No | No | No | No |
| Logratio transformation | ilr, alr, clr | ilr, alr, clr | ilr, alr | No | ilr, alr, clr |
| Multilevel logratios | Yes | No | No | No | No |
| Missing value and zero imputation | No | Yes | Yes | Yes | Yes |
| Outlier detection | No | Yes | No | No | Yes |
| **Frequentist models** | | | | | |
| Single-level with compositional predictors | No | No | Yes | No | No |
| Single-level with compositional outcomes | No | No | Yes | No | No |
| Multilevel with compositional predictors | No | No | No | No | No |
| Multilevel with compositional outcomes | No | No | No | No | No |
| **Bayesian models** | | | | | |
| Single-level with compositional predictors | Yes | No | No | No | No |
| Single-level with compositional outcomes* | Yes | No | No | No | No |
| Multilevel with compositional predictors | Yes | No | No | No | No |
| Multilevel with compositional outcomes* | Yes | No | No | No | No |
| **Pivot coordinate estimation** | | | | | |
| Single-level with compositional predictors | Yes† | No | No | No | No |
| Single-level with compositional outcomes* | No | No | No | No | No |
| Multilevel with compositional predictors | Yes† | No | No | No | No |
| Multilevel with compositional outcomes* | No | No | No | No | No |
| **Compositional substitution analysis** | | | | | |
| Simple single-level | Yes† | No | No | No | No |
| Simple multilevel | Yes† | No | No | No | No |
| Average single-level | Yes† | No | No | No | No |
| Average multilevel | Yes† | No | No | No | No |

**Figure 1**

*Implementing the Bayesian multilevel compositional data analysis using R package*
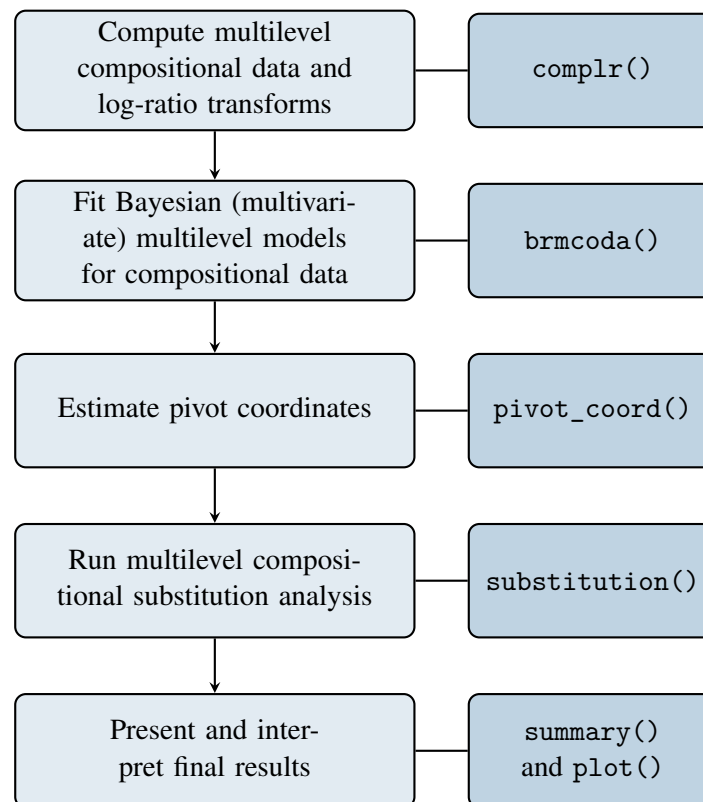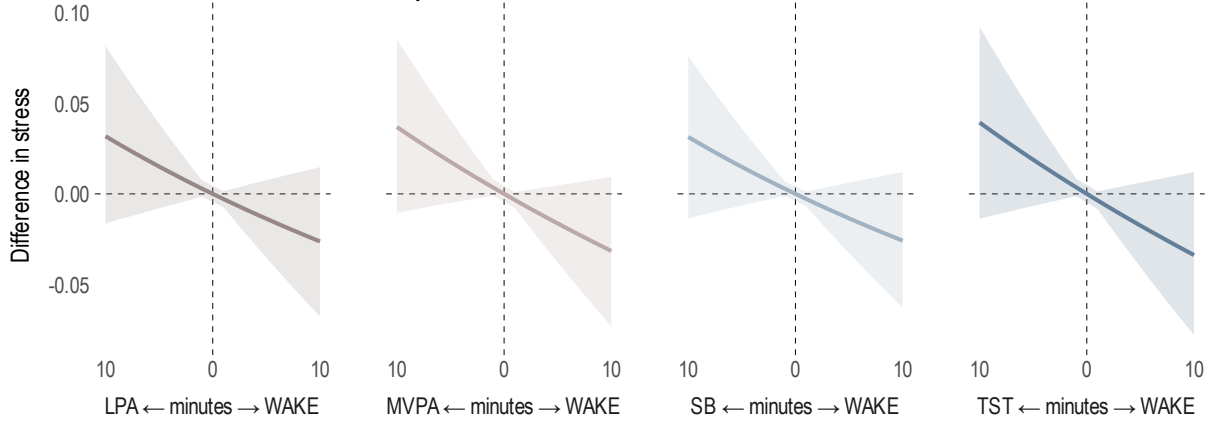***multilevelcoda****.*

**Figure 2**

*Difference in stress for 1-10 minutes of reallocation between time awake in bed and other behaviours. TST = total sleep time, WAKE = time awake in bed, MVPA = moderate-to-vigorous physical activity, LPA = light physical activity, SB = sedentary behaviour.*



**A. Awake in Bed Reallocations at Between-person level**



**B. Awake in Bed Reallocations at Within-person level**