

# Mixture Multigroup Bayesian SEM With Approximate Measurement Invariance for Comparing Structural Relations Across Many Groups

Hongwei Zhao<sup>1</sup> , Jeroen K. Vermunt<sup>2</sup> , Kim De Roover<sup>1,2</sup> 

[1] *Quantitative Psychology and Individual Differences, KU Leuven, Leuven, Belgium.* [2] *Department of Methodology, Tilburg School of Social and Behavioral Sciences, Tilburg University, Tilburg, the Netherlands.*

---

Methodology, 2025, Vol. 21(4), 286–312, <https://doi.org/10.5964/meth.16411>

**Received:** 2024-12-15 • **Accepted:** 2025-09-16 • **Published (VoR):** 2025-12-18

**Handling Editor:** Jone Aliri, University of the Basque Country, Leioa, Spain.

**Corresponding Author:** Hongwei Zhao, Quantitative Psychology and Individual Differences, KU Leuven, Tiensestraat 102, 3000 Leuven, Belgium. E-mail: [hongwei.zhao@kuleuven.be](mailto:hongwei.zhao@kuleuven.be)

**Supplementary Materials:** Code, Materials [see [Index of Supplementary Materials](#)]



## Abstract

In social sciences, researchers often compare relations between constructs, referred to as “structural relations”, across a large number of groups. This paper proposes Mixture Multigroup Bayesian SEM (MixMG-BSEM), a novel method for comparing structural relations across many groups while accounting for approximate measurement invariance in factor loadings. Traditional methods often assume exact measurement invariance, which may not reflect real-world data where small differences in measurement parameters commonly occur across many groups. MixMG-BSEM addresses this by using Multigroup Bayesian CFA with small-variance priors to allow for these small differences, and groups are then clustered based on their structural relations using Mixture Modeling. This is done in a stepwise estimation procedure built on the structural-after-measurement approach. By combining cluster-specific structural relations with small between-group differences in measurement parameters, MixMG-BSEM obtains a clustering that is driven only by the structural relations. The robustness and effectiveness of MixMG-BSEM are demonstrated through a simulation study.

## Keywords

Multigroup Bayesian SEM, Approximate Measurement Invariance, Mixture Modeling, Structural Relation



This is an open access article distributed under the terms of the [Creative Commons Attribution 4.0 International License, CC BY 4.0](#), which permits unrestricted use, distribution, and reproduction, provided the original work is properly cited.

In social sciences, Structural Equation Modelling (SEM; [Bollen, 1989](#); [Hoyle, 2012](#)) is widely used to investigate relations between constructs (e.g., emotions, motivation), referred to as “structural relations” within SEM. Researchers are often interested in how these structural relations vary across groups. For instance, [Michael and Kyriakides \(2023\)](#) examined how academic motivation mediated the effect of socioeconomic status on reading achievement among 15-year-old students and how this differed across 38 countries.

To study differences in structural relations, Multigroup SEM (MG-SEM) and Multilevel SEM (ML-SEM) can be used. MG-SEM estimates the structural relations for each group and allows testing whether they are equal across groups. ML-SEM captures variations in structural relations by normally distributed random effects around the overall mean estimate for each relation. Even though group-specific estimates of relations can be derived from random effects, only the mean and variance of each random effect are part of the model parameters, which makes ML-SEM more parsimonious, allowing for accurate parameter estimates in case of very small sample sizes per group. To pinpoint which groups have the same relations and for which groups they differ, MG-SEM and ML-SEM require pairwise comparisons of group-specific relations. As the number of groups increases, performing pairwise comparisons quickly becomes infeasible. For example, for 38 groups, this requires 703 pairwise comparisons per structural relation. To reduce the number of comparisons, mixture modeling ([McLachlan & Peel, 2000](#)) can be used to cluster groups based on similarity of the structural relations. Before performing such a clustering, it is essential to ensure that the structural relations are validly comparable across groups and that they are the only source of differences driving the clustering.

In social sciences, the constructs of interest are typically unobserved or latent variables, also known as “factors” in SEM. SEM addresses their latent nature by including a measurement model (MM), which specifies how latent variables are measured by observed indicators (often questionnaire items), whereas the relations of interest among the latent variables are part of the structural model (SM). For valid comparisons of constructs and their relations, measurement invariance (MI) must hold across the groups. MI implies that the MM is equal across groups, meaning that the constructs are measured in the same way, so that observed differences reflect differences in the constructs rather than differences in measurement.

MI is examined at different levels by assessing the equality of different subsets of MM parameters. Configural invariance evaluates whether the factor structure is the same across groups, meaning that, in each group, the same set of indicators relates to a factor. The strength and direction of the relations between factors and indicators are quantified by factor loadings. Whereas configural invariance only deals with which factor loadings are non-zero, weak or metric invariance requires the loadings to be equal across groups. Next, strong and strict invariance impose equality of the items’ intercepts and residual or

‘unique’ variances, respectively. Metric invariance is a prerequisite for validly comparing structural relations, whereas strong and strict invariance are not required. When full metric invariance (i.e., invariance of all loadings) does not hold, partial metric invariance (i.e., invariance of some loadings) still enables valid comparisons of structural relations (Byrne et al., 1989), as long as the loading differences are captured in the model (e.g., by group specific loadings). The same holds for differences in item intercepts and unique variances.

When combining SEM with mixture modeling, groups can be clustered on their structural relations by making the structural relations cluster-specific (i.e., the same for all groups assigned to a cluster). In traditional mixture SEM methods (Arminger & Stein, 1997; Dolan & van der Maas, 1998; Jedidi et al., 1997), MM parameters can be specified as invariant or cluster-specific, implying that MM differences can either be ignored or captured by the same clustering. To cluster groups only on the structural relations rather than also on differences in measurement, a framework of novel mixture SEM methods emerged recently. Perez Alonso and colleagues (2024) introduced Mixture Multigroup SEM (MixMG-SEM), which combines MG-SEM with mixture modeling. Zhao and colleagues (2025a) proposed Mixture Multilevel SEM (MixML-SEM), which builds the mixture clustering onto the more parsimonious ML-SEM. The aim of both methods is to cluster groups specifically on the structural relations while accounting for measurement non-invariance, but the difference is that MixML-SEM uses Multilevel Confirmatory Factor Analysis (ML-CFA) with random effects to deal with MM differences, whereas MixMG-SEM uses Multigroup Confirmatory Factor Analysis (MG-CFA) with group-specific MM parameters. Their estimation builds on the stepwise “Structural-After-Measurement” (Rossee & Loh, 2024) approach, where the MM is estimated first, using either MG-CFA or ML-CFA, followed by the SM, which includes clustering the groups on their structural relations. For comparability of structural relations, both methods require at least partial metric invariance and impose exact equality for the invariant factor loadings (i.e., exact MI). However, with a large number of groups, achieving exact MI is often unrealistic. To address this, Multigroup Bayesian SEM (MG-BSEM; Muthén & Asparouhov, 2012, 2013a, 2013b) with Approximate MI (AMI) uses priors with small variances for the MM parameters to allow for small differences across groups while keeping them approximately equal. In this paper, we present Mixture Multigroup BSEM (MixMG-BSEM), which extends MG-BSEM with mixture modeling to cluster groups on the structural relations while capturing approximate invariance of factor loadings.

MixMG-BSEM, MixMG-SEM and MixML-SEM differ in their first estimation step only, that is, in their MM and the corresponding MI assumptions. MixMG-SEM and MixML-SEM require exact invariance for (at least) some loadings, whereas the first step of MixMG-BSEM is a MG-CFA with Bayesian estimation (MG-BCFA) that assumes approximately invariant loadings. Approximate invariance lies between exact invariance (where parameters are exactly equal across groups) and non-invariance (where param-

eters can differ substantially across groups), where exact invariance is more closely approximated as the variances of the priors become smaller.

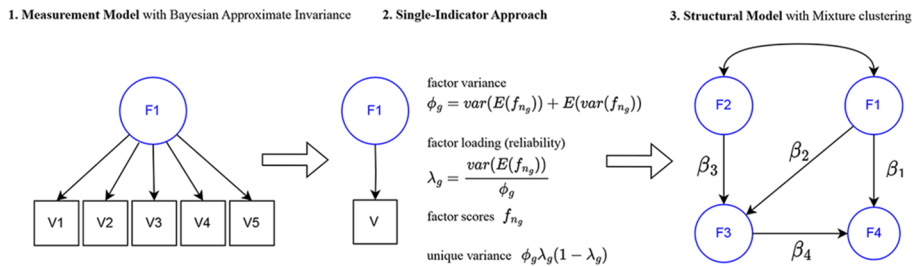
The paper is structured as follows: We begin with a description of MixMG-BSEM in the Method section. Next, we evaluate its performance through a Simulation Study. Finally, the Discussion section summarizes the main findings and addresses limitations and future directions.

## Method

As mentioned above, MixMG-BSEM is estimated in a stepwise manner, building on the SAM approach (see Figure 1). In Step 1, MG-BCFA with small-variance priors is performed for each factor, and factor scores are extracted. In Step 2, these factor scores are used as single indicators to obtain group-specific factor covariances with Croon’s correction (Croon, 2002). In Step 3, the SM is estimated, including the clustering and the cluster-specific structural relations, using an Expectation-Maximization (Dempster et al., 1977) algorithm for maximum likelihood estimation. Note that Steps 2 and 3 are the same as for MixML-SEM and are therefore only briefly described below (for details, see Zhao et al., 2025a).

Figure 1

Mixture Multigroup Bayesian SEM with Approximate Measurement Invariance



### Step 1: Measurement Model With Bayesian Approximate Measurement Invariance

The MM defines how the factors are measured by the items and MG-CFA is used to compare MMs across groups. Note that we estimate the MM per factor, which implies that we assume the factors to be independent in Step 1. Indicating an individual in Group  $g$  ( $g = 1, \dots, G$ ) by  $n_g$  and gathering the responses on the  $J_q$  items measuring Factor  $q$  ( $q = 1, \dots, Q$ ) in the Vector  $\mathbf{x}_{n_g}$ , the MG-CFA model for Factor  $q$  is expressed as:

$$\mathbf{x}_{n_g} = \boldsymbol{\tau}_g + \boldsymbol{\lambda}_g \eta_{n_g} + \boldsymbol{\epsilon}_{n_g} \text{ with } \boldsymbol{\epsilon}_{n_g} \sim MVN(0, \boldsymbol{\Theta}_g) \quad (1)$$

where  $\boldsymbol{\tau}_g$  is a  $J_q$ -dimensional vector of intercepts for Group  $g$ ,  $\boldsymbol{\lambda}_g$  is a  $J_q$ -dimensional vector of factor loadings (i.e., item-factor relations) for Group  $g$ ,  $\eta_{n_g}$  denotes the latent variable score for the individual, and  $\boldsymbol{\epsilon}_{n_g}$  is a  $J_q$ -dimensional vector of residuals, with the diagonal of  $\boldsymbol{\Theta}_g$  containing the group-specific unique variances of the items. To set the scale of each factor, one can either set its variance to one or use the marker variable approach by fixing one loading (ideally, a strong and invariant loading) to one, for each group. In this paper, we adopt the marker variable approach to ensure that a one-unit change in the underlying factor has the same meaning across groups.

Since small differences in MM parameters are common across many groups and still allow for latent variable comparisons, we apply the assumption of approximate metric invariance (i.e.,  $\boldsymbol{\lambda}_g \approx \boldsymbol{\lambda}$  for all Groups  $g$ ) instead of exact invariance (i.e.,  $\boldsymbol{\lambda}_g = \boldsymbol{\lambda}$ ) in Step 1 of MixMG-BSEM. This is accomplished by using MG-CFA with Bayesian estimation<sup>1</sup> (MG-BCFA) and applying small-variance, normally distributed priors to the corresponding factor loadings, which constrain the group-specific loadings to be approximately equal (Muthén & Asparouhov, 2012). For this, both *Mplus* (Muthén & Muthén, 1998) and the R-package *blavaan* (Merkle et al., 2021) are available, but we use *blavaan* by default because it is free and open-source. In *blavaan*, AMI is achieved by applying small-variance priors in every group except for the reference group, which is the first group (by default). A non-informative prior is used for the parameter in the first group and the parameter estimate for this group is used as the mean of the small-variance priors for that same parameter in the other groups.

Since Bayesian estimation can be computationally challenging, two measures are taken to lower the computation time of Step 1 of MixMG-BSEM: (1) the data are centered per group to remove the mean structure (i.e.,  $\boldsymbol{\tau}_g = 0$  and  $\boldsymbol{\alpha}_g = 0$ ), which is irrelevant to the comparison of structural relations, and (2) MG-BCFA is performed for each factor separately, which is in line with the “measurement blocks” approach in SAM (Rosseel & Loh, 2024) with one factor per measurement block. This approach lowers the number of parameters to be estimated and also enhances the model’s robustness against MM misspecifications, such as a few unmodeled crossloadings.

In this paper, we assume that all factor loadings, except for the marker variable loadings, are approximately invariant, while the unique variances and factor variances are estimated as group-specific parameters (i.e., with the default non-informative priors per group). Note that it is harmless to specify exactly invariant loadings as approximately invariant since they will then be estimated as nearly identical across groups. Of course, in practice, combinations of exactly and approximately invariant loadings can be applied

1) The possibility to use a different estimator in each step, such as Bayesian estimation for the MM and maximum likelihood for the SM (see also Zhao et al., 2025a) is an important advantage of the SAM approach.

in Step 1 of MixMG-BSEM. Moreover, in theory, all combinations of exact invariance, approximate invariance and non-invariance can be used, but complex combinations may cause convergence problems.

To determine which parameters are (approximately) invariant or non-invariant, MI testing should be performed prior to using MixMG-BSEM. Note that, if exact invariance does not hold for a parameter, standard MG-CFA requires a tedious process of comparing group-specific parameter estimates to determine whether differences reflect non-invariance or approximate invariance. Instead, MG-BCFA allows to test the tenability of AMI directly by imposing small-variance priors on MM parameters and assessing model fit. [Muthén and Asparouhov \(2012\)](#) recommend starting with a very small variance (e.g., 0.001) and, if needed, the priors' variances can be increased to reach a good model fit. In this way, MG-BCFA provides information on how large the parameter differences are (i.e., on the level of AMI). Model fit can be assessed using the posterior predictive  $p$  value ([Gelman et al., 1996](#)), but it is not very sensitive to the prior variances in case of large samples. Other fit measures include the Bayesian RMSEA (BRMSEA; [Hoofs et al., 2018](#)) and the Deviance Information Criterion (DIC; [Spiegelhalter et al., 2002](#)). The DIC balances model fit (i.e., the posterior mean deviance) and complexity (i.e., the effective number of parameters) in Bayesian models, with smaller values indicating a better balance. Regarding the prior selection in MG-BSEM, [Kim et al. \(2017\)](#) found that the DIC often selected models with smaller prior variances when the sample size was small and [Pokropek et al. \(2020\)](#) found that the DIC performed better as sample size increased, and recommended using the DIC with thresholds tailored to different sample sizes.

Once the marker variables and the approximately invariant loadings are confirmed by the MI testing, we obtain the specification of the MG-BCFA model that corresponds to the first step of MixMG-BSEM. In the next step, we need estimates of the factor scores and their uncertainty. To this end, the means and standard deviations of the posterior distributions of the individuals' factor scores (i.e., estimated latent variable scores) are appended to the data file.

## Step 2: Single-Indicator Approach to Obtain Group-Specific Factor Covariances

In a single-indicator approach, the factor scores are used as the "observed" proxy (or a single indicator) for the latent variable ([Vermunt, 2025](#)). Since factor scores are only estimates of the true latent variable scores, we apply Croon's correction (2002) to the factor score covariances ( $\text{cov}(\mathbf{f}_g)$ ) to obtain unbiased estimates of the true latent variable covariances ( $\text{cov}(\boldsymbol{\eta}_g)$ ), here denoted as  $\Phi_g^{s2}$ :

$$\Phi_g^{s2} = \widehat{\Lambda}_g^{-1} (\text{cov}(\mathbf{F}_g) - \widehat{\Theta}_g) (\widehat{\Lambda}_g')^{-1} \quad (2)$$

where  $\widehat{\Lambda}_g$  corresponds to the  $Q \times Q$  diagonal matrix of group-specific factor loadings (reflecting the reliability of the factor scores) and  $\widehat{\Theta}_g$  is the  $Q \times Q$  diagonal matrix of group-specific unique variances. These estimates correspond to the MM parameters of the single indicators of the factors (i.e., the factor scores) rather than the original, observed indicators. These MM parameters are derived from the posterior mean and standard deviation estimates for the factor scores, obtained from Step 1. For details, please refer to Equations (7–8) in the MixML-SEM paper (Zhao et al., 2025a).

### Step 3: Structural Model With Mixture Clustering of the Groups

In Step 3, MixMG-BSEM clusters the groups and estimates cluster-specific structural relations. The SM is thus conditional on the cluster membership,  $z_{gk}$ , which denotes whether Group  $g$  belongs to Cluster  $k$ . Whereas the true cluster membership is assumed to be either 1 or 0, its estimation,  $\widehat{z}_{gk}$ , ranges from 0 to 1 and represents the probability of Group  $g$  belonging to Cluster  $k$ . The model-implied factor covariance matrix  $\Phi_{gk}$ , given that  $z_{gk} = 1$ , is defined as:

$$\Phi_{gk} = (\mathbf{I} - \mathbf{B}_k)^{-1} \Psi_{gk} (\mathbf{I} - \mathbf{B}_k)^{-1'} \quad (3)$$

where  $\mathbf{B}_k$  contains the cluster-specific regression coefficients between latent variables, and  $\Psi_{gk}$  is the residual factor covariance matrix, which is specified as group-and-cluster-specific to ensure that clustering is driven only by the regressions  $\mathbf{B}_k$  (for details, see Perez Alonso et al., 2024). The SM is estimated with maximum likelihood estimation using  $\Phi_g^s$  as input.

For the mixture clustering in MixMG-BSEM, it is assumed that the (true) latent variable scores  $\boldsymbol{\eta}_{n_g}$  are sampled from a mixture of  $K$  multivariate normal distributions. Specifically, all latent variable scores of Group  $g$ ,  $\mathbf{H}_g$ , are assumed to be sampled from the same distribution:

$$f(\mathbf{H}_g; \mathbf{v}) = \sum_{k=1}^K \pi_k \prod_{n_g=1}^{N_g} MVN(\boldsymbol{\eta}_{n_g}; \boldsymbol{\alpha}_g, \Phi_{gk}) \text{ with } \sum_{k=1}^K \pi_k = 1 \quad (4)$$

where  $f$  is the population density function,  $\mathbf{v}$  represents the set of population parameters, and  $\pi_k$  is the prior probability that Group  $g$  belongs to Cluster  $k$ . The scores in  $\mathbf{H}_g$  are assumed to follow a normal distribution with  $\boldsymbol{\alpha}_g$  as the factor mean (which is zero due to centering) and  $\Phi_{gk}$  as the factor covariance matrix. The unknown parameters  $\mathbf{v}$  are estimated by maximizing the following log-likelihood function:

$$\begin{aligned} \log L_{\eta} &= \log \left( \prod_{g=1}^G \sum_{k=1}^K \pi_k \frac{1}{(2\pi)^{Q/2} |\Phi_{gk}|^{1/2}} \exp \left( -\frac{1}{2} \text{tr}(\Phi_g^{s2} \Phi_{gk}^{-1}) \right) \right)^{N_g} \\ &= \sum_{g=1}^G \log \left( \sum_{k=1}^K \pi_k \frac{1}{(2\pi)^{Q/2} |\Phi_{gk}|^{1/2}} \exp \left( -\frac{1}{2} \text{tr}(\Phi_g^{s2} \Phi_{gk}^{-1}) \right) \right)^{N_g} \end{aligned} \quad (5)$$

where  $\Phi_g^{s2}$  is the group-specific factor covariance matrix from Step 2 (Equation 2), and  $\Phi_{gk}$  is the group-and-cluster-specific factor covariance matrix from Step 3 (Equation 3). The maximum likelihood estimation is performed using the EM algorithm (Dempster et al., 1977). Specifically, in the E-step, the algorithm estimates the classification probabilities  $\hat{z}_{gk}$  given the current parameter estimates. In the M-step, the algorithm estimates the unknown parameters  $\nu$  given the classification probabilities obtained from the E-step. The E- and M-steps are iterated until convergence. A multi-start procedure is applied to mitigate convergence to local maxima, where the converged solution with the highest loglikelihood across the different starts is selected as the final result. For an in-depth explanation of the technical details of Step 3, readers are referred to Appendix A of the paper by Perez Alonso et al. (2024).

## Simulation

In the simulation study, we evaluated the performance of MixMG-BSEM, assuming the true number of clusters was known. Firstly, we aimed to examine how MixMG-BSEM's performance was affected by the within-group sample size, the number of groups, the number of clusters, the item reliability, the AMI of the loadings, and the size of (differences in) regression parameters. On top of that, since the first step of MixMG-BSEM estimates the MM per factor, we evaluated the consequences of ignoring crossloadings in this step. Literature on traditional SEM has shown that factor correlations tend to be overestimated when crossloadings are constrained to zero (e.g., Asparouhov et al., 2015; Marsh et al., 2009, 2010, 2014), which may affect the comparison of structural relations. However, given its stepwise estimation and measurement block approach, MixMG-BSEM may be relatively robust to overlooked crossloadings (Rosseeel & Loh, 2024), but the recovery of clusters and regression parameters may still decline in case of multiple crossloadings. Secondly, in terms of the analysis, we examined the impact of a key aspect of the Bayesian estimation; that is, the impact of different prior variances for the loadings on the recovery of clusters and cluster-specific regressions. We expected that using too narrow priors might fail to capture the loading differences across groups, which may affect the estimation of and clustering on the structural relations. Additionally, we also evaluated which prior was selected by the Deviance Information Criterion (DIC), since selecting this prior is an important step in empirical practice.

In a complete factorial design, the following factors were manipulated:

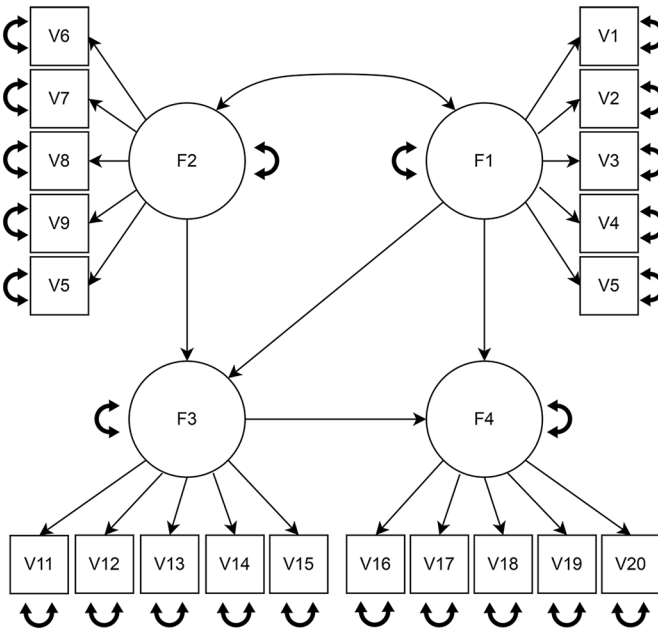
1. Total number of groups  $G$  (2 levels): 24, 48.
2. Within-group sample size  $N_g$  (3 levels): 50, 100, 200.
3. Number of clusters  $K$  (2 levels): 2, 4.
4. Size of regression parameters  $\beta$  (3 levels): 0.2, 0.3, 0.4.
5. Item reliability (2 levels): high, low.
6. Level of AMI for loadings (5 levels): 0.001, 0.005, 0.01, 0.05, 0.1.
7. Size of crossloadings (3 levels): 0, 0.2, 0.4.

We chose a minimum of 24 groups with group sizes  $N_g$  ranging from 50 to 200, which partially correspond to the group sizes in other simulation studies on Bayesian AMI (Kim et al., 2017; Lek et al., 2018). The number of groups in each cluster depended on the number of groups  $G$ , and the number of Clusters  $K$ , where each cluster contained an equal number of groups. Note that larger  $G$ , larger  $N_g$ , and smaller  $K$  result in larger within-cluster sample sizes, which were expected to improve the performance of MixMG-BSEM.

The data were generated from a SEM model with four latent variables, each measured by five items (see Figure 2), as in Perez Alonso et al. (2024) and Zhao et al. (2025a). Specifically, the data were generated from a multivariate normal distribution (MVN) with covariance matrix  $\Sigma_{gk}$ , determined by the parameters  $\mathbf{B}_k$ ,  $\Psi_{gk}$ ,  $\Lambda_g$  and  $\Theta_g$  (see Equation (6) in Perez Alonso et al., 2024).

**Figure 2**

*The Data-Generating Model With Exogenous Factors F1 and F2 and Endogenous Factors F3 and F4*

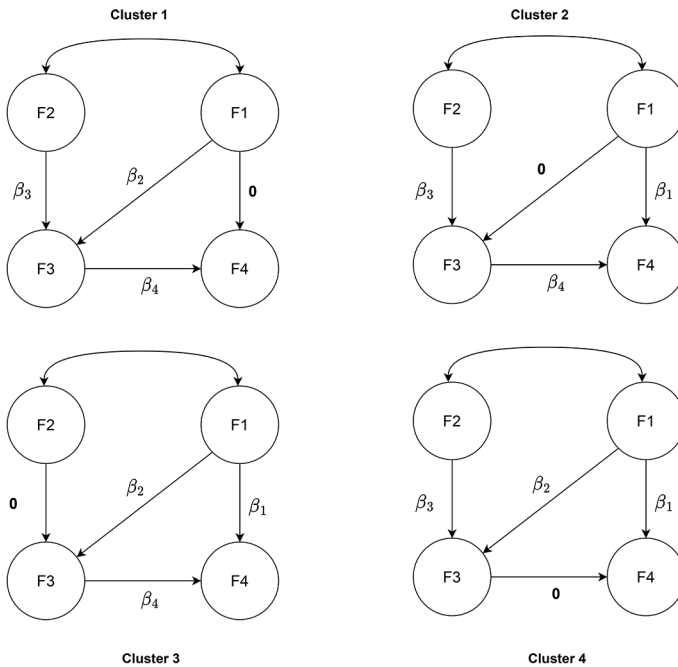


The size of the regression parameters was set to  $\beta$  and, as shown in Figure 3, the differences between clusters were introduced by setting one regression parameter to zero in each cluster. Hence, larger values of  $\beta$  resulted in larger differences and thus in greater separation between clusters, which should make the clusters easier to recover.

For the group-and-cluster-specific residual factor covariances  $\Psi_{gk}$ , we sampled the variances of the exogenous factors  $F1$  and  $F2$  from a uniform distribution  $U(0.75, 1.25)$  and their covariance from  $U(-0.3, 0.3)$ . The total variances of the endogenous factors  $F3$  and  $F4$  were also sampled from  $U(0.75, 1.25)$  and their residual variances are determined as follows: For  $F3$  and  $F4$ , it was computed as  $\text{Var}(F3)_g - (\beta_{2,k}^2 \text{Var}(F1)_g + \beta_{3,k}^2 \text{Var}(F2)_g + 2\beta_{2,k}\beta_{3,k} \text{Cov}(F1, F2)_g)$  and  $\text{Var}(F4)_g - (\beta_{1,k}^2 \text{Var}(F1)_g + \beta_{4,k}^2 \text{Var}(F3)_g + 2\beta_{1,k}\beta_{4,k}(\beta_{2,k} \text{Var}(F1)_g + \beta_{3,k} \text{Cov}(F1, F2)_g))$ , respectively.

Figure 3

The Cluster-Specific Structural Relations



In Loading Matrix  $\Lambda_g$ , the first loading of each factor was fixed to one. The other loadings (except for crossloadings) were approximately invariant across groups and were sampled from a normal distribution with a mean of  $\sqrt{0.6}$  in case of a high item reliability and  $\sqrt{0.4}$  in case of a low item reliability, and a variance that depended on the level of the AML. According to a summary by Arts et al. (2021), prior variances used in simulation and empirical studies typically range from 0.001 to 0.1. Following this, our simulations included five levels of AMI (i.e., 0.001, 0.005, 0.01, 0.05, 0.1). For instance, to obtain an AMI level of 0.01, which implies a variance of 0.01 for differences in loadings, we sampled loadings from a normal distribution with a variance of 0.005 for all groups.<sup>2</sup> Per factor, one crossloading was added to the third item measuring the next factor (i.e., Item

2) In *blavaan*, the estimate of a parameter in the first group is used as the mean of the prior for that same parameter in the other groups. Consequently, the prior reflects the differences of the other groups to the reference group. The variance of the difference between two factor loadings equals the sum of their individual variances, assuming there is no covariance between them. For all groups, including the reference group, we sampled loadings from a normal distribution with a variance that is half the targeted MI level for all groups, so that the variance of the loading differences toward the reference group equals the targeted MI level.

8 crossloaded on Factor 1, Item 13 on Factor 2, Item 18 on Factor 3, and Item 3 on Factor 4). A value of 0 corresponded to no crossloading, 0.2 to a moderate crossloading, and 0.4 to a large crossloading. The unique variances on the diagonal of  $\Theta_g$  were sampled from  $U(0.3, 0.5)$  under high reliability and from  $U(0.50, 0.70)$  under low reliability.

Finally, the data were sampled from  $MVN(0, \Sigma_{gk})$  for each group. In total, we generated  $2$  (number of groups)  $\times 3$  (within-group sample size)  $\times 2$  (number of clusters)  $\times 3$  (size of regression parameters)  $\times 2$  (reliability)  $\times 5$  (size of AMI)  $\times 3$  (size of crossloadings)  $\times 50$  (replications) = 54,000 data sets according to the described procedure, using R Version 4.4 (R Core Team, 2022). All data sets were analyzed with MixMG-BSEM with 50 random starts and the true number of clusters. For each data set, we performed the analysis five times, with different prior variances for the loadings (i.e., 0.001, 0.005, 0.01, 0.05, 0.1) in Step 1, to examine the performance of MixMG-BSEM across different prior variances. The analyses were performed on a supercomputer consisting of 2 Intel Xeon Platinum 8468 CPUs (Sapphire Rapids). The average computation time was 50.5 minutes ( $SD = 24.8$ ) for Step 1 (mainly influenced by  $G$  and  $N_g$ ) with the four factors estimated sequentially, 0.03 minutes ( $SD = 0.02$ ) for the intermediate Step 2, and 5.7 minutes ( $SD = 5.7$ ) for Step 3 (mainly influenced by  $N_g$ ,  $K$ , and  $\beta$ ).<sup>3</sup> We assessed MCMC convergence of the Bayesian estimation using the potential scale reduction factor (PSRF). On average, the PSRF was 1.000 ( $SD = 0.001$ ), indicating sufficient convergence, with no notable differences across conditions. These results suggest that the MCMC chains converged properly for all simulated data sets.

## Results

### Recovery of Factor Loadings

We evaluated the recovery of the group-specific factor loading estimates for each item  $j$ , using the Root Mean Squared Error (RMSE) across groups as follows:

---

3) The first step of MixMG-BSEM (i.e., estimating the MM using *blavaan*) can be computationally demanding, especially for larger sample sizes. Luckily, the stepwise estimation of MixMG-BSEM implies that the MM needs to be estimated only once, even when estimating the SM with different numbers of clusters for model selection. To illustrate, we report computation times for one of the largest data sets in our simulation study, which included 48 groups and 200 observations per group. Using *blavaan*, estimating the four factors sequentially without parallel computing for the MCMC chains (per factor) took around 62 minutes in total. With parallel computing applied to the three MCMC chains (per factor), the computation time decreased to 42 minutes. Note that additional speed gains could be achieved by further parallelizing across the four factors, depending on hardware capabilities and user preferences. Alternatively, *Mplus* offers a more time-efficient estimation of the MM with AMI, though it is commercial software. For the same data set, *Mplus* completed the estimation of the four factors sequentially in only 2 minutes. It is worth noting that eliminating the mean structure by centering per group (see [Method](#) section) helped as well, since these computation times of *blavaan* and *Mplus* increased to 66 and 91 minutes, respectively, when including the mean structure in the model.

$$RMSE_{\lambda_j} = \sqrt{\frac{\sum_{g=1}^G (\hat{\lambda}_{gj} - \lambda_{gj})^2}{G}} \quad (3)$$

where  $\lambda_{gj}$  is the true group-specific loading of the  $j$ -th item on the factor, and  $\hat{\lambda}_{gj}$  is the corresponding estimate. Note that  $RMSE_{\lambda_j}$  is also affected by the variance of the estimates, rather than just the bias.

When using MixMG-BSEM with the true prior variances for the loadings, the average  $RMSE_{\lambda_j}$  across the four factors and all simulated data sets was 0.066, 0.091, 0.066, and 0.066, respectively, for the loadings of the second to the fifth item of each factor (Table 1, last row). Note that  $RMSE_{\lambda_3}$  was larger due to the disregarded crossloadings on that item. This was also the only  $RMSE_{\lambda_j}$  value that differed across the four factors. Specifically, the  $RMSE_{\lambda_3}$  values were 0.093, 0.072, 0.098, and 0.101, for  $F1$  to  $F4$ , respectively. It seems that the third loading for  $F2$  is affected by the ignored crossloading the least, which may be explained by the fact that, unlike the other factors,  $F2$  is involved in only one direct regression relation with the other factors<sup>4</sup> and is thus less correlated with the other factors. When the crossloadings were zero (i.e., without crossloadings),  $RMSE_{\lambda_3}$  took on the same values as for the other loadings ( $RMSE_{\lambda_3} = 0.066$ ), whereas they increased with larger crossloadings (see Table 1).  $RMSE_{\lambda_3}$  was also higher in case of larger  $\beta$ , which implies stronger correlations between factors (see Table 1). Note that larger  $N_g$ , higher item reliability, and lower levels of AMI – thus applying lower prior variances – resulted in lower  $RMSE_{\lambda}$  for all items. The latter is explained by the fact that a lower prior variance more strongly approximates an equality constraint, which lowers the sample size requirements and thus the estimates' variability for a given sample size.

**Table 1**

*The Average  $RMSE_{\lambda_j}$  for Factor Loading Estimates When Using the True Prior Variances for the Loadings*

Factor	Level	$RMSE_{\lambda_2}$ (SD)	$RMSE_{\lambda_3}$ (SD)	$RMSE_{\lambda_4}$ (SD)	$RMSE_{\lambda_5}$ (SD)
G	24	0.066 (0.031)	0.091 (0.041)	0.066 (0.031)	0.066 (0.031)
	48	0.065 (0.031)	0.090 (0.041)	0.065 (0.031)	0.065 (0.031)
$N_g$	50	0.079 (0.038)	0.102 (0.046)	0.079 (0.038)	0.079 (0.038)
	100	0.065 (0.028)	0.091 (0.039)	0.065 (0.028)	0.065 (0.028)
	200	0.052 (0.019)	0.081 (0.035)	0.052 (0.019)	0.052 (0.019)
K	2	0.066 (0.031)	0.091 (0.041)	0.066 (0.031)	0.066 (0.031)
	4	0.066 (0.031)	0.091 (0.041)	0.066 (0.031)	0.066 (0.031)

4) It has indirect relations with the other factors via the correlation between  $F1$  and  $F2$ , but the expected value of this correlation is zero.

Factor	Level	$RMSE_{\lambda_2}$ (SD)	$RMSE_{\lambda_3}$ (SD)	$RMSE_{\lambda_4}$ (SD)	$RMSE_{\lambda_5}$ (SD)
$\beta$	0.2	0.066 (0.031)	0.080 (0.035)	0.066 (0.031)	0.066 (0.031)
	0.3	0.066 (0.031)	0.090 (0.039)	0.066 (0.031)	0.066 (0.031)
	0.4	0.065 (0.031)	0.103 (0.045)	0.065 (0.031)	0.065 (0.031)
Reliability	high	0.061 (0.028)	0.087 (0.039)	0.061 (0.027)	0.061 (0.027)
	low	0.070 (0.034)	0.095 (0.043)	0.070 (0.034)	0.070 (0.034)
AMI	0.001	0.028 (0.006)	0.055 (0.029)	0.028 (0.006)	0.028 (0.006)
	0.005	0.048 (0.006)	0.072 (0.027)	0.048 (0.006)	0.048 (0.006)
	0.01	0.060 (0.009)	0.084 (0.027)	0.060 (0.009)	0.060 (0.009)
	0.05	0.090 (0.020)	0.116 (0.032)	0.090 (0.020)	0.090 (0.020)
	0.1	0.102 (0.027)	0.128 (0.037)	0.102 (0.026)	0.102 (0.027)
Crossloadings	0	0.066 (0.031)	0.066 (0.031)	0.066 (0.031)	0.066 (0.031)
	0.2	0.066 (0.031)	0.086 (0.032)	0.066 (0.031)	0.065 (0.031)
	0.4	0.066 (0.031)	0.122 (0.037)	0.066 (0.031)	0.066 (0.031)
Total		0.066 (0.031)	0.091 (0.041)	0.066 (0.031)	0.066 (0.031)

To illustrate the effect of the prior variances for the loadings,  $RMSE_{\lambda_2}$  across different prior variances is shown in Figure 4. The diagonal of the plot represents cases where the prior variances were correctly specified, while the lower part shows cases where the priors were narrower than the true level of AMI. In general, applying too narrow or too wide priors resulted in larger  $RMSE_{\lambda_2}$  values. Since the prior variance affected the loading recovery, we also evaluated prior selection using the Deviance Information Criterion (DIC), which balances model fit and complexity. When looking at the prior selection per loading, the correct selection rate was 59.2% across all loadings and all simulated data sets.<sup>5</sup> For 28.2% of the data sets, the prior selection was flawless in the sense that true priors were selected for *all* loadings. Generally, the DIC tended to select either the true or slightly smaller prior variances. Specifically, for an AMI level of 0.001, DIC correctly selected the prior variance of 0.001 for 86.6% of the loadings, with smaller proportions selecting 0.005 (12.5%) and 0.01 (0.9%). For an AMI level of 0.005, the correct selection rate was 59.4%, followed by 0.001 (36.6%) and 0.01 (4.0%). For an AMI level of 0.01, DIC most frequently selected 0.005 (71.1%), followed by 0.01 (23.5%) and 0.001 (5.4%). For an AMI level of 0.05, DIC primarily selected the true prior variance (81.5%), followed by prior variances of 0.1 (10.0%) and 0.01 (8.5%). For an AMI level of 0.1, DIC mostly selected prior variances of 0.05 (55.1%) and 0.1 (44.9%). Although prior selection based on the DIC varied across different levels of AMI and was not always optimal, the

5) Similar results were found with the widely applicable information criterion (WAIC; Watanabe, 2010) and leave-one-out information criterion (LOOIC; Vehtari et al., 2017): WAIC: 58.2%; LOOIC: 57.6%.

resulting  $RMSE_{\lambda}$  values remained relatively stable given the selected priors, suggesting stable factor loading recovery despite prior variance misspecification.

**Figure 4**

$RMSE_{\lambda_2}$  Across Different Prior Variances, Indicated by the Columns, Whereas the Rows Represent the True Levels of AMI

**Recovery of factor loadings (RMSE)**

		Prior variance				
		0.001	0.005	0.01	0.05	0.1
Approximate MI	0.001	0.028	0.039	0.050	0.082	0.096
	0.005	0.048	0.048	0.055	0.083	0.096
	0.01	0.065	0.058	0.060	0.084	0.096
	0.05	0.139	0.108	0.094	0.090	0.098
	0.1	0.198	0.151	0.126	0.098	0.102

*Note.* The diagonal (in white) contains cases where the prior variances were correctly specified, while the lower part represents cases where the priors were too narrow. For each row, the cells are colored red if the  $RMSE_{\lambda_2}$  is larger than the  $RMSE_{\lambda_2}$  on the diagonal, and blue if it is smaller.

## Sensitivity to Local Maxima

To evaluate how often (Step 3 of) MixMG-BSEM converged to a local maximum, we compared the log-likelihood of the final best solution (out of 50 random starts) to the one obtained when starting from the true clustering, which is a proxy for the global maximum. If  $\log L_{\eta}$  was more than 0.001 lower than the proxy, the solution was considered a local maximum. Overall, when applying the true priors, MixMG-BSEM ended up in a local maximum for 0.01% of the data sets, with all local maxima occurring in case of  $N_g = 50$  with  $K = 4$ .

## Recovery of Clusters

The Adjusted Rand Index (ARI; Hubert & Arabie, 1985) measures the similarity between two partitions while correcting for chance, with a value of one indicating perfect agreement and zero indicating the level of agreement between two random partitions. To compute the ARI, the modal clustering (i.e., assigning each group to the cluster with the highest classification probability) was compared to the true clustering. Additionally, the correct clustering rate (%CC) was computed as the percentage of correctly clustered groups.

When using the true priors, the average ARI across all simulated data was 0.644 and the %CC was 84.8%.<sup>6</sup> To check which main effects and two-way interactions among the manipulated factors significantly influenced the ARI, we conducted an analysis of variance (ANOVA) using the *aov* function in R. The ANOVA results table is provided in Supplementary Material (S1; Zhao et al., 2025b). Firstly, all main effects were significant at the  $\alpha = 0.01$  level. From Table 2, we see that larger  $G$ , larger  $N_g$ , smaller  $K$ , larger  $\beta$ , higher reliability, all led to better recovery of clusters, with  $N_g$  and  $\beta$  having the strongest effect, as reflected in higher ARI and %CC. The differences in ARI across levels of AMI and crossloadings were very small, suggesting that cluster recovery was relatively unaffected by these factors. Secondly,  $N_g$ ,  $K$ ,  $\beta$  and reliability all interacted significantly with one another (i.e., all six of their two-way interactions were significant), which is why we further illustrate the four-way interaction effect of these four manipulated factors on ARI in Figure 5. We see that the ARI was very sensitive to  $N_g$ , and more so in case of  $K = 4$ , smaller  $\beta$ , or low reliability. According to Steinley (2004), ARI values above 0.80 indicate good cluster recovery. Using this rule-of-thumb, the cluster recovery was good when  $\beta = 0.4$  with  $N_g \geq 100$ , or when  $\beta = 0.3$  with  $N_g = 200$ . When  $\beta = 0.2$ , the ARI only exceeded 0.80 with  $N_g = 200$  for two clusters and a high reliability.

Across different prior variances, the ARI remained relatively stable. When the applied prior was too narrow or too large, the ARI slightly dropped. For example, for an AMI level of 0.1, it decreased from 0.635 when using the true prior variance to 0.608 when using a prior variance of 0.001. For an AMI level of 0.001, it decreased from 0.651 when using the true prior variance to 0.638 when using a prior variance of 0.1.

---

6) We further examined the distributions of the ARI using boxplots across relevant conditions (see Figure 1 in the Supplementary Material, S3; Zhao et al., 2025b). Under more challenging conditions (e.g., small sample sizes, low reliability, small regression effects), we observed larger variability. For example, with  $\beta = 0.2$  and  $N_g = 50$ , the distributions were more right-skewed, indicating a majority of low ARI values. In contrast, under easier conditions, such as  $\beta = 0.4$  and  $N_g = 200$ , the distributions became symmetric or slightly left-skewed, with less variability, indicating more consistent high ARI values. Therefore, we have added a table reporting the medians and MADs for ARI and %CC in the Supplementary Material (S4; Zhao et al., 2025b), to complement the means and SDs reported in the main text.

**Table 2**

The Average ARI and Correct Clustering Rate (%CC) When Using the True Prior Variances for the Loadings

Factor	Level	ARI (SD)	%CC (SD)
G	24	0.628 (0.331)	0.840 (0.170)
	48	0.660 (0.301)	0.856 (0.159)
$N_g$	50	0.402 (0.280)	0.733 (0.174)
	100	0.663 (0.284)	0.861 (0.145)
	200	0.865 (0.187)	0.949 (0.082)
K	2	0.719 (0.295)	0.913 (0.106)
	4	0.569 (0.320)	0.782 (0.187)
$\beta$	0.2	0.400 (0.282)	0.733 (0.175)
	0.3	0.685 (0.280)	0.871 (0.142)
	0.4	0.849 (0.201)	0.940 (0.095)
Reliability	high	0.682 (0.309)	0.865 (0.159)
	low	0.606 (0.320)	0.830 (0.170)
AMI	0.001	0.651 (0.314)	0.851 (0.163)
	0.005	0.647 (0.317)	0.849 (0.164)
	0.01	0.649 (0.316)	0.850 (0.164)
	0.05	0.638 (0.319)	0.845 (0.167)
	0.1	0.635 (0.318)	0.843 (0.166)
Crossloadings	0	0.654 (0.316)	0.852 (0.164)
	0.2	0.644 (0.318)	0.848 (0.166)
	0.4	0.634 (0.318)	0.844 (0.166)
Total		0.644 (0.317)	0.848 (0.165)

## Recovery of Regression Parameters

To evaluate the recovery of the regression parameters, we computed the  $ME_\beta$  and the  $RMSE_\beta$  per regression parameter (i.e.,  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$ , and  $\beta_4$ ):

$$ME_\beta = \frac{\sum_{k=1}^K (\hat{\beta}_k - \beta_k)}{K} \# (7)$$

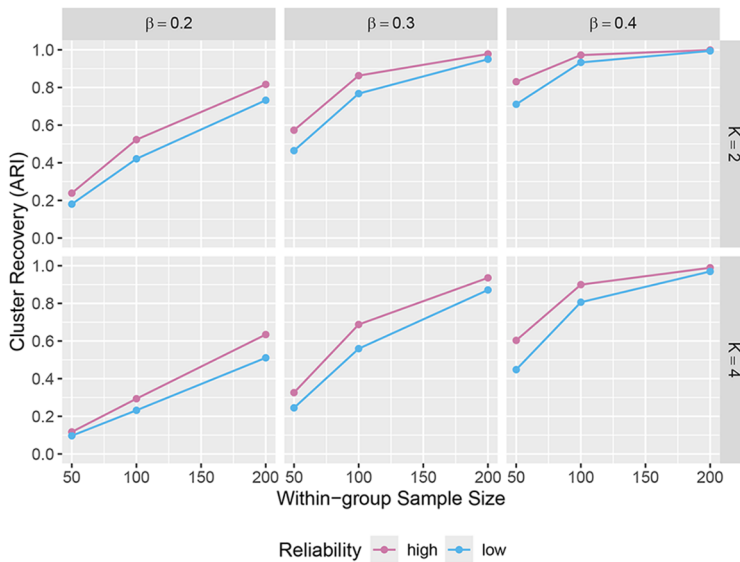
$$RMSE_\beta = \sqrt{\frac{\sum_{k=1}^K (\hat{\beta}_k - \beta_k)^2}{K}} \# (8)$$

where  $\hat{\beta}_k$  is the estimated regression coefficient in cluster  $k$  and  $\beta_k$  is the corresponding true value. To get a better idea of the separate  $\beta$  coefficients in different clusters (i.e., without averaging across clusters), we also reported the bias for  $\beta_1$  in the Supplementary

Material as an example (S2; Zhao et al., 2025b), with positive values indicating overestimation in most conditions. Similar trends were observed for  $\beta_3$  and  $\beta_4$ . For  $\beta_2$ , the bias values were the smallest and close to zero (slightly negative).

**Figure 5**

*The ARI for MixMG-BSEM in Function of the Within-Group Sample Sizes, Number of Clusters, Size of Regression Parameters, and Reliability*



On average,  $ME_{\beta}$  was 0.037, -0.004, 0.039, and 0.031 (Table 3), and  $RMSE_{\beta}$  was 0.073, 0.054, 0.072 and 0.069 (Table 4) for  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$ , and  $\beta_4$ , respectively.<sup>7</sup> We performed an ANOVA on  $ME_{\beta_1}$  as an example (S1). All main effects were significant at the  $\alpha = 0.01$  level, except for AMI. Similar to the trends observed for the cluster recovery, larger  $G$ , larger  $N_g$ , smaller  $K$ , larger  $\beta$ , higher reliability, and smaller crossloadings resulted in smaller  $ME_{\beta}$ , with crossloadings having the strongest effect (Table 3). To get a better idea of how  $ME_{\beta}$  got influenced by the manipulated factors, we depicted the interaction effects of  $N_g$ ,  $K$ , and  $\beta$  across different levels of crossloadings and reliability (Figure 6). Specifically, larger crossloadings resulted in larger  $ME_{\beta}$  values, especially in case of lower item reliability. Note that the recovery of the regression parameters was barely affected

<sup>7</sup> We also checked the distributions of  $ME_{\beta_i}$ , which were mostly symmetric (see Figure 2 in the Supplementary Material, S3; Zhao et al., 2025b). Thus, we retained the use of means and SDs and did not include a separate table for the medians and MADs for these outcomes.

by using different prior variances, even more narrow ones, likely due to the fact that the cluster recovery was hardly affected as well.

**Table 3**

*The Average  $ME_{\beta}$  for Each of the Four Estimated Regression Parameters When Using the True Prior Variances for the Loadings*

Factor	Level	$ME_{\beta_1}$ (SD)	$ME_{\beta_2}$ (SD)	$ME_{\beta_3}$ (SD)	$ME_{\beta_4}$ (SD)
G	24	0.037 (0.057)	-0.004 (0.045)	0.039 (0.053)	0.032 (0.051)
	48	0.036 (0.046)	-0.004 (0.033)	0.039 (0.043)	0.031 (0.041)
$N_g$	50	0.040 (0.071)	-0.002 (0.059)	0.040 (0.063)	0.034 (0.063)
	100	0.035 (0.044)	-0.005 (0.031)	0.038 (0.043)	0.030 (0.040)
	200	0.034 (0.033)	-0.006 (0.018)	0.038 (0.036)	0.029 (0.031)
K	2	0.035 (0.046)	-0.005 (0.036)	0.038 (0.042)	0.031 (0.039)
	4	0.038 (0.056)	-0.004 (0.043)	0.040 (0.054)	0.032 (0.053)
$\beta$	0.2	0.040 (0.059)	-0.002 (0.048)	0.041 (0.055)	0.036 (0.052)
	0.3	0.037 (0.050)	-0.004 (0.038)	0.039 (0.047)	0.032 (0.046)
	0.4	0.033 (0.044)	-0.007 (0.031)	0.037 (0.042)	0.026 (0.041)
Reliability	high	0.033 (0.045)	-0.005 (0.035)	0.035 (0.043)	0.028 (0.042)
	low	0.040 (0.058)	-0.004 (0.044)	0.042 (0.053)	0.034 (0.051)
AMI	0.001	0.037 (0.050)	-0.004 (0.039)	0.040 (0.048)	0.032 (0.046)
	0.005	0.036 (0.050)	-0.004 (0.039)	0.040 (0.048)	0.033 (0.046)
	0.01	0.037 (0.051)	-0.004 (0.039)	0.039 (0.048)	0.032 (0.046)
	0.05	0.037 (0.053)	-0.005 (0.041)	0.039 (0.050)	0.030 (0.047)
	0.1	0.036 (0.054)	-0.005 (0.041)	0.037 (0.048)	0.029 (0.048)
Crossloadings	0	0.001 (0.043)	0.001 (0.040)	-0.000 (0.036)	0.002 (0.040)
	0.2	0.040 (0.043)	-0.004 (0.040)	0.042 (0.037)	0.033 (0.040)
	0.4	0.069 (0.045)	-0.010 (0.039)	0.075 (0.039)	0.059 (0.041)
Total		0.037 (0.052)	-0.004 (0.040)	0.039 (0.048)	0.031 (0.047)

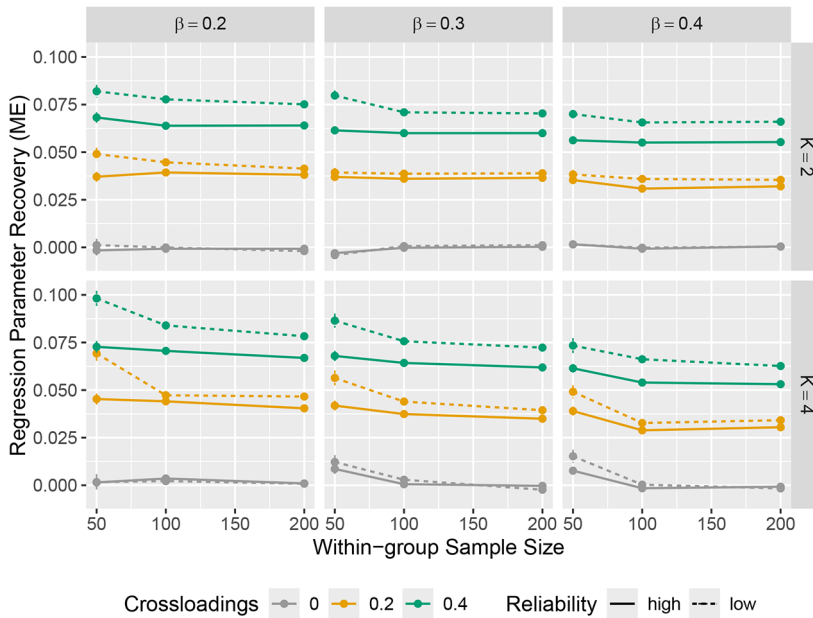
**Table 4**

*The Average RMSE<sub>β</sub> for Each of the Four Estimated Regression Parameters When Using the True Prior Variances for the Loadings*

Factor	Level	RMSE <sub>β<sub>1</sub></sub> (SD)	RMSE <sub>β<sub>2</sub></sub> (SD)	RMSE <sub>β<sub>3</sub></sub> (SD)	RMSE <sub>β<sub>4</sub></sub> (SD)
G	24	0.081 (0.065)	0.062 (0.057)	0.079 (0.058)	0.077 (0.061)
	48	0.064 (0.051)	0.045 (0.044)	0.065 (0.046)	0.061 (0.048)
N <sub>g</sub>	50	0.105 (0.079)	0.087 (0.069)	0.100 (0.070)	0.101 (0.075)
	100	0.064 (0.042)	0.047 (0.035)	0.065 (0.039)	0.061 (0.039)
	200	0.048 (0.027)	0.027 (0.017)	0.051 (0.029)	0.045 (0.024)
K	2	0.055 (0.039)	0.036 (0.032)	0.056 (0.037)	0.052 (0.035)
	4	0.090 (0.070)	0.071 (0.061)	0.088 (0.061)	0.087 (0.066)
β	0.2	0.086 (0.067)	0.067 (0.061)	0.085 (0.062)	0.083 (0.062)
	0.3	0.070 (0.057)	0.051 (0.050)	0.070 (0.051)	0.067 (0.054)
	0.4	0.061 (0.049)	0.043 (0.038)	0.060 (0.041)	0.058 (0.048)
Reliability	high	0.064 (0.048)	0.048 (0.045)	0.065 (0.047)	0.062 (0.048)
	low	0.081 (0.068)	0.059 (0.057)	0.079 (0.058)	0.076 (0.062)
AMI	0.001	0.072 (0.055)	0.052 (0.050)	0.071 (0.050)	0.068 (0.053)
	0.005	0.072 (0.056)	0.053 (0.051)	0.071 (0.052)	0.069 (0.054)
	0.01	0.072 (0.057)	0.053 (0.051)	0.071 (0.051)	0.069 (0.054)
	0.05	0.074 (0.062)	0.055 (0.052)	0.073 (0.057)	0.070 (0.057)
	0.1	0.074 (0.066)	0.055 (0.054)	0.072 (0.054)	0.070 (0.061)
Crossloadings	0	0.055 (0.058)	0.052 (0.052)	0.051 (0.052)	0.056 (0.055)
	0.2	0.070 (0.058)	0.053 (0.053)	0.069 (0.050)	0.067 (0.055)
	0.4	0.092 (0.056)	0.055 (0.050)	0.096 (0.048)	0.085 (0.053)
Total		0.073 (0.059)	0.054 (0.052)	0.072 (0.053)	0.069 (0.056)

**Figure 6**

*The  $ME_{\beta}$  for MixMG-BSEM in Function of the Within-Group Sample Sizes, Number of Clusters, Size of Regression Parameters, Crossloadings and Reliability*



**Section Conclusion**

We assessed the performance of MixMG-BSEM when the true number of clusters was known. We found that performing 50 random starts in Step 3 largely prevented local maxima. The recovery of clusters and regression parameters was good when the within-group sample size was at least 200 and/or in case of a larger cluster separation (i.e.,  $\beta = 0.4$ ). Ignoring larger crossloadings (by estimating the MM per factor) resulted in more biased estimates for factor loadings and regression parameters, while cluster recovery was less affected. Although prior selection based on the DIC was not always optimal, the recovery of clusters and regression parameters was relatively stable across the priors.

## Discussion

We presented MixMG-BSEM as a new addition to the novel mixture SEM framework for comparing structural relations across many groups. Unlike the existing approaches that rely on the exact MI assumption, MixMG-BSEM adopts the more realistic assumption of AMI, which accommodates small differences in MM parameters across groups. Specifically, after estimating the MM using MG-BCFA with small-variance priors, MixMG-BSEM clusters groups with the same structural relations, thereby eliminating the need for pairwise comparisons of group-specific structural relations. Since its results may depend on the specified prior variances, results obtained with different prior variances should be compared and the best prior selected based on model fit measures such as the Deviance Information Criterion (DIC). Our simulation study results showed that both the clustering and regression parameter estimates were relatively insensitive to the choice of prior variances, however.

Currently, MixMG-BSEM estimates the MM per factor (i.e., with one factor per measurement block). In the simulation study, the cluster recovery was unaffected by ignoring crossloadings, but the recovery of the factor loadings and regression estimates was affected. Therefore, it would be valuable to investigate the performance of MixMG-BSEM when including factors with crossloadings in the same measurement block, at the cost of a (much) longer computation time, where small-variance priors could also be applied to the crossloadings to allow for small differences (Muthén & Asparouhov, 2012). However, it is important to note that the default prior mean for crossloadings is zero, whereas applying a prior mean of zero to a sizeable crossloading can negatively impact the regression parameter estimates (Wei et al., 2022). Therefore, researchers should gather prior information about crossloadings before choosing an appropriate prior (Wei et al., 2022). Note that, in cases with several crossloadings connecting more than two factors, the estimation may fail when all these factors are included in one measurement block, especially in case of many groups. Thus, in such cases, model estimation is only possible when partitioning the factors into smaller measurement blocks, highlighting the benefits of the measurement block approach even more.

While the simulation study evaluated the performance of MixMG-BSEM with approximate metric invariance for all loadings, except for the invariant marker variable loading, MixMG-BSEM can theoretically accommodate all combinations of exact, approximate and non-invariance for the loadings. The stepwise estimation of MixMG-BSEM conveniently allows to tweak the MG-BCFA model, for instance, by specifying certain loadings as non-invariant, before moving onto the next steps. Similarly, if group-specific loading estimates are virtually identical across groups, one may consider specifying the loading as exactly invariant. Specifying an invariant parameter as approximately invariant is rather harmless, whereas specifying a non-invariant parameter as approximately invariant may introduce bias in parameter estimation and affect the clustering. Note that MG-BCFA allows to evaluate non-invariances for all parameters, which is achieved by

comparing group-specific estimates to the credible intervals of the average posterior estimates across all groups (e.g., Winter & Depaoli, 2020). In future research, it would be interesting to evaluate the performance of MixMG-BSEM when non-invariant loadings are specified as approximately invariant.

The simulation study assumed the number of clusters to be known, whereas this is typically unknown for empirical data. To determine the number of clusters, different methods are available, such as the Bayesian Information Criterion (BIC; Schwarz, 1978), Akaike Information Criterion (AIC; Akaike, 1974), and convex hull procedure (CHull; Ceulemans & Kiers, 2006). In brief, all these methods balance model fit (i.e., the log-likelihood) and model complexity (i.e., the number of parameters). BIC and AIC do so by combining model fit and a penalty for model complexity into a single criterion, whereas CHull uses a generalized scree test. Previous studies on model selection for MixMG-SEM (Perez Alonso et al., 2025) and MixML-SEM (Zhao et al., 2025a) have shown that combining AIC, BIC, and CHull – with visual inspection of the scree plot – is an effective way to determine the number of clusters. Since MixMG-BSEM performs the same mixture clustering on group-specific factor covariances as these methods, we expect these recommendations to generalize to MixMG-BSEM. However, in the future, it would still be useful to evaluate model selection for MixMG-BSEM specifically.

Currently, MixMG-BSEM combines Bayesian and maximum likelihood estimation, assuming continuous items. In empirical practice, we often work with ordinal items with a few response categories (e.g., Likert scale items). To accommodate ordinal data in MixMG-BSEM, only the first step (i.e., MG-BCFA) would need to be adjusted to deal with ordinal data (Muthén & Asparouhov, 2013a), whereas the subsequent steps would remain unchanged. In future studies, it will be valuable to evaluate the performance of MixMG-BSEM adapted to ordinal data.

In conclusion, MixMG-BSEM is an effective method for accommodating AMI while clustering structural relations of interest. By relaxing the strict assumption of exact MI, it extends the framework of novel mixture SEM methods in an important way, making it more suited for empirical applications where small differences in parameters across groups are expected.

---

**Funding:** Funded by the European Union (ERC, PROCESSHETEROGENEITY, 101040754, awarded to Kim De Roover). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them. The resources and services used in this work were provided by the VSC (Flemish Supercomputer Center), funded by the Research Foundation – Flanders (FWO) and the Flemish Government.

---

**Acknowledgments:** The authors have no additional (i.e., non-financial) support to report.

---

**Competing Interests:** The authors have declared that no competing interests exist.

---

## Supplementary Materials

Type of supplementary materials	Availability/Access
<b>Data</b>	
No study data is available.	—
<b>Code</b>	
a. R Code - 2MixMG-SEM.	Zhao et al. (2024)
b. R Code - MixMG-BSEM.	Zhao et al. (2024)
c. R Code - MixML-SEM.	Zhao et al. (2024)
d. R Code - Empirical Application.	Zhao et al. (2024)
<b>Material</b>	
a. S1. ANOVA Tables.	Zhao et al. (2025b)
b. S2. Bias for $\beta_1$ as an example.	Zhao et al. (2025b)
c. S3. Boxplots for ARI and $ME_{\beta_1}$ .	Zhao et al. (2025b)
d. S4. Table for ARI and %CC with Median and MAD.	Zhao et al. (2025b)
e. S5. Empirical Application of MixMG-BSEM.	Zhao et al. (2025b)
<b>Study/Analysis preregistration</b>	
The study was not preregistered.	—
<b>Other</b>	
No other material to report.	—

## References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*(6), 716–723. <https://doi.org/10.1109/TAC.1974.1100705>
- Arminger, G., & Stein, P. (1997). Finite mixtures of covariance structure models with regressors: Loglikelihood function, minimum distance estimation, fit indices, and a complex example. *Sociological Methods & Research*, *26*(2), 148–182. <https://doi.org/10.1177/0049124197026002002>
- Arts, I., Fang, Q., van de Schoot, R., & Meitinger, K. (2021). Approximate measurement invariance of willingness to sacrifice for the environment across 30 countries: The importance of prior distributions and their visualization. *Frontiers in Psychology*, *12*, Article 624032. <https://doi.org/10.3389/fpsyg.2021.624032>
- Asparouhov, T., Muthén, B. O., & Morin, A. J. S. (2015). Bayesian structural equation modeling with cross-loadings and residual covariances: Comments on Stromeier et al. *Journal of Management*, *41*(6), 1561–1577. <https://doi.org/10.1177/0149206315591075>
- Bollen, K. A. (1989). *Structural equations with latent variables*. John Wiley & Sons. <https://doi.org/10.1002/9781118619179>
- Byrne, B. M., Shavelson, R. J., & Muthén, B. O. (1989). Testing for the equivalence of factor covariance and mean structures: The issue of partial measurement invariance. *Psychological Bulletin*, *105*(3), 456–466. <https://doi.org/10.1037/0033-2909.105.3.456>

- Ceulemans, E., & Kiers, H. A. L. (2006). Selecting among three-mode principal component models of different types and complexities: A numerical convex hull based method. *British Journal of Mathematical & Statistical Psychology*, 59(1), 133–150. <https://doi.org/10.1348/000711005X64817>
- Croon, M. (2002). Using predicted latent scores in general latent structure models. In G. A. Marcoulides & I. Moustaki (Eds.), *Latent variable and latent structure models* (1<sup>st</sup> ed., pp. 195–223). Lawrence Erlbaum.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1), 1–22. <https://doi.org/10.1111/j.2517-6161.1977.tb01600.x>
- Dolan, C. V., & van der Maas, H. L. J. (1998). Fitting multivariate normal finite mixtures subject to structural equation modeling. *Psychometrika*, 63(3), 227–253. <https://doi.org/10.1007/BF02294853>
- Gelman, A., Meng, X.-L., & Stern, H. (1996). Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica*, 6(4), 733–760.
- Hoofs, H., Van De Schoot, R., Jansen, N. W. H., & Kant, I. (2018). Evaluating model fit in Bayesian confirmatory factor analysis with large samples: Simulation study introducing the BRMSEA. *Educational and Psychological Measurement*, 78(4), 537–568. <https://doi.org/10.1177/0013164417709314>
- Hoyle, R. H. (2012). *Handbook of structural equation modeling*. Guilford Press.
- Hubert, L., & Arabie, P. (1985). Comparing partitions. *Journal of Classification*, 2(1), 193–218. <https://doi.org/10.1007/BF01908075>
- Jedidi, K., Jagpal, H. S., & Desarbo, W. S. (1997). Finite-mixture structural equation models for response-based segmentation and unobserved heterogeneity. *Marketing Science*, 16(1), 39–59. <https://doi.org/10.1287/mksc.16.1.39>
- Kim, E. S., Cao, C., Wang, Y., & Nguyen, D. T. (2017). Measurement invariance testing with many groups: A comparison of five approaches. *Structural Equation Modeling*, 24(4), 524–544. <https://doi.org/10.1080/10705511.2017.1304822>
- Lek, K., Oberski, D., Davidov, E., Ciecuch, J., Seddig, D., & Schmidt, P. (2018). Approximate measurement invariance. In T. P. Johnson, B. Pennell, I. A. L. Stoop & B. Dorer (Eds.), *Advances in comparative survey methods* (pp. 911–929). John Wiley & Sons. <https://doi.org/10.1002/9781118884997.ch41>
- Marsh, H. W., Lüdtke, O., Muthén, B. O., Asparouhov, T., Morin, A. J. S., Trautwein, U., & Nagengast, B. (2010). A new look at the big five factor structure through exploratory structural equation modeling. *Psychological Assessment*, 22(3), 471–491. <https://doi.org/10.1037/a0019227>
- Marsh, H. W., Morin, A. J. S., Parker, P. D., & Kaur, G. (2014). Exploratory structural equation modeling: An integration of the best features of exploratory and confirmatory factor analysis. *Annual Review of Clinical Psychology*, 10, 85–110. <https://doi.org/10.1146/annurev-clinpsy-032813-153700>
- Marsh, H. W., Muthén, B. O., Asparouhov, T., Lüdtke, O., Robitzsch, A., Morin, A. J. S., & Trautwein, U. (2009). Exploratory structural equation modeling, integrating CFA and EFA: Application to

- students' evaluations of university teaching. *Structural Equation Modeling*, 16(3), 439–476.  
<https://doi.org/10.1080/10705510903008220>
- McLachlan, G., & Peel, D. (2000). *Finite mixture models*. Wiley. <https://doi.org/10.1002/0471721182>
- Merkle, E. C., Fitzsimmons, E., Uanhoro, J., & Goodrich, B. (2021). Efficient Bayesian structural equation modeling in Stan. *Journal of Statistical Software*, 100(6), 1–22.  
<https://doi.org/10.18637/jss.v100.i06>
- Michael, D., & Kyriakides, L. (2023). Mediating effects of motivation and socioeconomic status on reading achievement: A secondary analysis of PISA 2018. *Large-Scale Assessments in Education*, 11, Article 31. <https://doi.org/10.1186/s40536-023-00181-9>
- Muthén, B., & Asparouhov, T. (2012). Bayesian structural equation modeling: A more flexible representation of substantive theory. *Psychological Methods*, 17(3), 313–335.  
<https://doi.org/10.1037/a0026802>
- Muthén, B., & Asparouhov, T. (2013a). *BSEM measurement invariance analysis* (Mplus Web Notes: No. 17). Mplus. <https://www.statmodel.com/examples/webnotes/webnote17.pdf>
- Muthén, B., & Asparouhov, T. (2013b). *New methods for the study of measurement invariance with many groups*. Mplus. <https://www.statmodel.com/download/PolAn.pdf>
- Muthén, L. K., & Muthén, B. O. (1998). *Mplus user's guide: Statistical analysis with latent variables* (8<sup>th</sup> ed.). Muthén & Muthén.  
[https://www.statmodel.com/download/usersguide/MplusUserGuideVer\\_8.pdf#:~:text=Following%20is%20the%20correct%20citation%20for%20this%20document%3A,Muth%C3%A9n%20%26%20Muth%C3%A9n%20Version%208%20April%202017%20](https://www.statmodel.com/download/usersguide/MplusUserGuideVer_8.pdf#:~:text=Following%20is%20the%20correct%20citation%20for%20this%20document%3A,Muth%C3%A9n%20%26%20Muth%C3%A9n%20Version%208%20April%202017%20)
- Perez Alonso, A. F., Rosseel, Y., Vermunt, J. K., & De Roover, K. (2024). Mixture multigroup structural equation modeling: A novel method for comparing structural relations across many groups. *Psychological Methods*. Advance online publication. <https://doi.org/10.1037/met0000667>
- Perez Alonso, A. F., Vermunt, J. K., Rosseel, Y., & Roover, K. D. (2025). Selecting the number of clusters in mixture multigroup structural equation modeling. *Methodology: European Journal of Research Methods for the Behavioral and Social Sciences*, 21(1), Article e14931.  
<https://doi.org/10.5964/meth.14931>
- Pokropek, A., Schmidt, P., & Davidov, E. (2020). Choosing priors in Bayesian measurement invariance modeling: A Monte Carlo simulation study. *Structural Equation Modeling*, 27(5), 750–764. <https://doi.org/10.1080/10705511.2019.1703708>
- R Core Team. (2022). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <https://www.R-project.org/>
- Rosseel, Y., & Loh, W. W. (2024). A structural after measurement approach to structural equation modeling. *Psychological Methods*, 29(3), 561–588. <https://doi.org/10.1037/met0000503>
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, 6(2), 461–464.  
<https://doi.org/10.1214/aos/1176344136>

- Spiegelhalter, D. J., Best, N. G., Carlin, B. P., & Van Der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(4), 583–639. <https://doi.org/10.1111/1467-9868.00353>
- Steinley, D. (2004). Properties of the Hubert-Arable Adjusted Rand Index. *Psychological Methods*, 9(3), 386–396. <https://doi.org/10.1037/1082-989X.9.3.386>
- Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, 27, 1413–1432. <https://doi.org/10.1007/s11222-016-9696-4>
- Vermunt, J. K. (2025). Stepwise estimation of latent variable models: An overview of approaches. *Statistical Modelling*, 25(6), 530–551. <https://doi.org/10.1177/1471082X251355693>
- Watanabe, S. (2010). Asymptotic equivalence of Bayes cross validation and widely applicable information criterion in singular learning theory. *Journal of Machine Learning Research*, 11(116), 3571–3594.
- Wei, X., Huang, J., Zhang, L., Pan, D., & Pan, J. (2022). Evaluation and comparison of SEM, ESEM, and BSEM in estimating structural models with potentially unknown cross-loadings. *Structural Equation Modeling*, 29(3), 327–338. <https://doi.org/10.1080/10705511.2021.2006664>
- Winter, S. D., & Depaoli, S. (2020). An illustration of Bayesian approximate measurement invariance with longitudinal data and a small sample size. *International Journal of Behavioral Development*, 44(4), 371–382. <https://doi.org/10.1177/0165025419880610>
- Zhao, H., Vermunt, J. K., & De Roover, K. (2024). *Novel mixture SEM methods for comparing structural relations among many groups* [OSF project page containing R study codes]. Open Science Framework. <https://osf.io/rtp78/>
- Zhao, H., Vermunt, J. K., & De Roover, K. (2025a). Mixture Multilevel SEM vs. Multilevel SEM for comparing structural relations across groups in presence of measurement non-invariance. *Frontiers in Psychology*, 16, Article 1463790. <https://doi.org/10.3389/fpsyg.2025.1463790>
- Zhao, H., Vermunt, J. K., & De Roover, K. (2025b). *Supplementary materials to “Mixture multigroup Bayesian SEM with approximate measurement invariance for comparing structural relations across many groups”* [Supplementary materials with ANOVA, bias, and ARI tables; boxplots for ARI and ME, and empirical application of MixMG-BSEM]. PsychOpen GOLD. <https://doi.org/10.23668/psycharchives.21454>



*Methodology* (METH) is the official journal of the European Association of Methodology (EAM).



PsychOpen GOLD is a publishing service provided by the Leibniz Institute for Psychology (ZPID), Germany.