

# Optimal forecast reconciliation with time series selection

Xiaoqian Wang\*

and

Rob J Hyndman

and

Shanika L Wickramasuriya

Department of Econometrics & Business Statistics, Monash University

December 11, 2023

## Abstract

Forecast reconciliation ensures forecasts of time series in a hierarchy adhere to aggregation constraints, enabling aligned decision making. While forecast reconciliation can improve overall accuracy in hierarchical or grouped structures, the most substantial improvements occur in series with initially poor-performing base forecasts, and some series may still experience a deterioration in reconciled forecasts. However, in practice, some series in a structure often have poor base forecasts due to model misspecification or low forecastability. To address this, we propose two categories of forecast reconciliation methods that incorporate time series selection based on out-of-sample and in-sample information, respectively. Our methods keep “poor” base forecasts of some series unused in forming reconciled forecasts, preventing their negative impact on reconciliation. This process adjusts weights allocated to the remaining series accordingly when generating bottom-level reconciled forecasts. Additionally, our methods mitigate disparities arising from using different estimates of the base forecast error covariance matrix, thus alleviating the challenge of estimator selection. We evaluate the proposed methods through two simulation studies and empirical applications using Australian labour force data and domestic tourism data, showing improved accuracy compared to alternative methods, especially for higher aggregation levels, longer forecast horizons, and model misspecification.

*Keywords:* Coherent, Hierarchical time series, Grouped time series, Linear forecast reconciliation, Optimization problem

---

\*Corresponding author.

# 1 Introduction

Forecast reconciliation is a post-processing method that ensures forecasts of multivariate time series adhere to known linear constraints (Hyndman et al. 2011, Wickramasuriya et al. 2019). For example, the sum of regional unemployment forecasts should be equal to the national unemployment forecast.

Hyndman et al. (2011) introduced optimal forecast reconciliation, whereby “base” forecasts of all series are generated independently, and then adjusted to satisfy the constraints, leading to a set of coherent reconciled forecasts. Subsequent research has extended and developed the idea in the context of cross-sectional data (Hyndman et al. 2016, Wickramasuriya et al. 2019, Panagiotelis et al. 2021), temporal data (Athanasopoulos et al. 2017), and cross-temporal data (Di Fonzo & Girolimetto 2023). Athanasopoulos et al. (2024) provided a comprehensive introduction to the forecast reconciliation literature.

Reconciliation is known to improve overall forecast accuracy in collections of time series with aggregation constraints. On average, when the base forecasts are unbiased, the mean squared reconciled forecast error from the minimum trace reconciliation method (Wickramasuriya et al. 2019) is lower than that from the base forecasts (Wickramasuriya 2021). Most of the improvements attributed to reconciliation are observed in series with initially poor-performing base forecasts (Athanasopoulos et al. 2017). In practice, it is not uncommon for some series to have poor base forecasts due to challenges such as model misspecification or low signal-to-noise ratio (SNR). In such cases, it may be advantageous to exclude the worst base forecasts when performing reconciliation. This is the motivation for our proposed methods.

First, we propose forecast reconciliation methods that incorporate time series selection based on out-of-sample information, assuming unbiased base forecasts. We formulate this as an optimization problem, using diverse penalty functions to control the number of nonzero column entries in the weighting matrix for linear forecast reconciliation. We show that the number of selected time series is at least equal to the number of series at the bottom level, and we can reconstruct the entire structure by aggregating/disaggregating the selected series. Second, we relax the unbiasedness assumption and introduce an additional reconciliation method with selection, utilizing in-sample observations and their fitted values. This enables us to use the

in-sample reconciliation performance for selection purposes. In this case, it may happen that fewer than the number of series at the bottom level are used for reconciliation. Through simulation experiments and two empirical applications, we demonstrate that our proposed methods guarantee coherent forecasts that outperform or match their respective benchmark methods. The improvements are particularly pronounced when focusing on higher aggregation levels, longer forecast horizons, and cases of model misspecification. A remarkable feature of the proposed methods is their ability to diminish disparities arising from using different estimates of the base forecast error covariance matrix, thereby mitigating challenges associated with estimator selection, which is a prominent concern in the field of forecast reconciliation research.

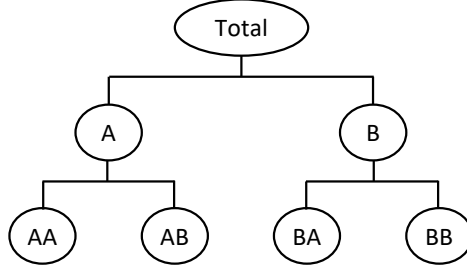
The remainder of the paper is structured as follows. Section 2 presents the notations and a review of linear forecast reconciliation methods. Section 3 introduces our proposed methods to achieve time series selection in reconciliation, and provides some theoretical insights. Section 4 and Section 5 show the results from simulations and two real-world datasets, respectively, followed by concluding remarks in Section 6. The R code for reproducing the results is available at <https://github.com/xqnwang/hfs>.

## 2 Preliminaries

### 2.1 Notation

A *hierarchical time series* is an  $n$ -dimensional multivariate time series that adheres to known linear constraints. Let  $\mathbf{y}_t \in \mathbb{R}^n$  be a vector comprising observations from all time series in the hierarchy at time  $t$ , and  $\mathbf{b}_t \in \mathbb{R}^{n_b}$  be a vector comprising observations of only the most disaggregated (“bottom-level”) time series at time  $t$ . The full hierarchy at time  $t$  can be written as

$$\mathbf{y}_t = \mathbf{S}\mathbf{b}_t,$$



**Figure 1:** An example of a two-level hierarchical time series.

for  $t = 1, 2, \dots, T$ , where  $T$  is the length of the time series, and  $\mathbf{S}$  is an  $n \times n_b$  *summing matrix* that defines the aggregation constraints. We can write the summing matrix as  $\mathbf{S} = \begin{bmatrix} \mathbf{A} \\ \mathbf{I}_{n_b} \end{bmatrix}$ , where  $\mathbf{A}$  is an  $n_a \times n_b$  *aggregation matrix* with  $n = n_a + n_b$ , and  $\mathbf{I}_{n_b}$  is an  $n_b$ -dimensional identity matrix.

For example, Figure 1 shows a simple hierarchy with  $n = 7$ ,  $n_b = 4$ ,  $n_a = 3$ ,  $\mathbf{y}_t = [y_{\text{Total},t}, y_{A,t}, y_{B,t}, y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}]'$ ,  $\mathbf{b}_t = [y_{AA,t}, y_{AB,t}, y_{BA,t}, y_{BB,t}]'$ , and

$$\mathbf{S} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ & & \mathbf{I}_4 \end{bmatrix}.$$

The notation is general enough to include aggregation constraints that are non-hierarchical. Please refer to [Hyndman & Athanasopoulos \(2021\)](#) for further details.

## 2.2 Linear forecast reconciliation

Let  $\hat{\mathbf{y}}_{T+h|T} \in \mathbb{R}^n$  be a vector of  $h$ -step-ahead *base forecasts* for all time series in the structure, given observations up to and including time  $T$ , and stacked in the same order as  $\mathbf{y}_t$ . We can use any method to generate these forecasts, but in general they will not be coherent (i.e., they won't satisfy the constraints).

Let  $\tilde{\mathbf{y}}_{T+h|T} \in \mathbb{R}^n$  denote a vector of  $h$ -step-ahead *reconciled forecasts* given by

$$\tilde{\mathbf{y}}_{T+h|T} = \mathbf{S} \mathbf{G}_h \hat{\mathbf{y}}_{T+h|T}, \quad (1)$$

where  $\mathbf{G}_h$  is an  $n_b \times n$  weighting matrix and  $\mathbf{S}$  is an  $n \times n_b$  summing matrix ([Wickramasuriya et al. 2019](#)).

## Minimum trace reconciliation

Let  $\hat{e}_{t+h|t} = \mathbf{y}_{t+h} - \hat{\mathbf{y}}_{t+h|t}$  denote the  $h$ -step-ahead in-sample *base forecast errors*, and  $\tilde{e}_{t+h|t} = \mathbf{y}_{t+h} - \tilde{\mathbf{y}}_{t+h|t}$  denote the  $h$ -step-ahead *reconciled forecast errors*, for  $t = 1, 2, \dots, T - h$ . Wickramasuriya et al. (2019) formulated the reconciliation problem as minimizing the trace (MinT) of the  $h$ -step-ahead covariance matrix of the reconciled forecast errors,  $\text{Var}(\tilde{e}_{t+h|t})$ . Assuming the base forecasts are unbiased, and that we want the reconciled forecasts to be unbiased, the unique solution to the minimization problem is

$$\mathbf{G}_h = (\mathbf{S}' \mathbf{W}_h^{-1} \mathbf{S})^{-1} \mathbf{S}' \mathbf{W}_h^{-1}, \quad (2)$$

where  $\mathbf{W}_h$  is the positive definite covariance matrix of the  $h$ -step-ahead base forecast errors.

The trace minimization problem can be reformulated as a least squares problem with linear constraints:

$$\min_{\tilde{\mathbf{y}}_{T+h|T}} \frac{1}{2} (\hat{\mathbf{y}}_{T+h|T} - \tilde{\mathbf{y}}_{T+h|T})' \mathbf{W}_h^{-1} (\hat{\mathbf{y}}_{T+h|T} - \tilde{\mathbf{y}}_{T+h|T}) \quad \text{s.t.} \quad \tilde{\mathbf{y}}_{T+h|T} = \mathbf{S} \tilde{\mathbf{b}}_{T+h|T}, \quad (3)$$

where  $\tilde{\mathbf{b}}_{T+h|T} \in \mathbb{R}^{n_b}$  comprises the  $h$ -step-ahead bottom-level reconciled forecasts, made at time  $T$ . The intuition behind MinT reconciliation is that **the larger the estimated variance of the base forecast errors, the larger the range of adjustments permitted for forecast reconciliation**.

It is challenging to estimate  $\mathbf{W}_h$ , especially for  $h > 1$ . It is common to assume  $\mathbf{W}_h = k_h \mathbf{W}_1, \forall h$ , where  $k_h > 0$ ; then the MinT solution of  $\mathbf{G}$  does not change with the forecast horizon,  $h$ . Hence, we will drop the subscript  $h$  for ease of exposition. Table 1 lists the most popularly used candidate estimators for  $\mathbf{W}_h$ .

## Relaxation of the unbiasedness assumptions

Both Hyndman et al. (2011) and Wickramasuriya et al. (2019) imposed unbiasedness conditions on the base forecasts and the reconciled forecasts. Ben Taieb & Koo (2019) proposed a reconciliation method relaxing the assumption of unbiasedness. Specifically, by expanding the training window incrementally, one observation at a time, they formulated the reconciliation problem as a regularized empirical risk

**Table 1:** Forecast reconciliation methods for which different estimators of  $\mathbf{W}$  are used.

Reconciliation method	$\mathbf{W}_h \propto$
OLS (Hyndman et al. 2011)	$\mathbf{I}$
WLSs (Athanasopoulos et al. 2017)	$\text{diag}(\mathbf{S}\mathbf{1})$
WLSv (Hyndman et al. 2016)	$\text{Diag}(\hat{\mathbf{W}}_1)$
MinT (Wickramasuriya et al. 2019)	$\hat{\mathbf{W}}_1$
MinTs (Wickramasuriya et al. 2019)	$\lambda \text{Diag}(\hat{\mathbf{W}}_1) + (1 - \lambda)\hat{\mathbf{W}}_1$

Note:  $\mathbf{1}$  is a vector of 1s of size  $n_b$ ,  $\text{diag}(\cdot)$  constructs a diagonal matrix using a given vector,  $\hat{\mathbf{W}}_1$  denotes the unbiased covariance estimator based on the in-sample one-step-ahead base forecast errors (i.e., residuals), and  $\text{Diag}(\cdot)$  forms a diagonal matrix using the diagonal elements of the input matrix.

minimization (RERM) problem given by

$$\min_{\mathbf{G}_h} \frac{1}{(T - T_1 - h + 1)n} \|\mathbf{Y}_h^* - \hat{\mathbf{Y}}_h^* \mathbf{G}_h' \mathbf{S}'\|_F^2 + \lambda \|\text{vec}(\mathbf{G}_h)\|_1,$$

where  $T_1$  denotes the minimum number of observations used for model training,  $\|\cdot\|_F$  is the Frobenius norm,  $\|\cdot\|_1$  is the  $L_1$  norm,  $\text{vec}(\cdot)$  denotes the vectorization of a matrix (stacking the columns of the matrix),  $\mathbf{Y}_h^* = [\mathbf{y}_{T_1+h}, \dots, \mathbf{y}_T]'$ ,  $\hat{\mathbf{Y}}_h^* = [\hat{\mathbf{y}}_{T_1+h|T_1}, \dots, \hat{\mathbf{y}}_{T|T-h}]'$ , and  $\lambda \geq 0$  is a regularization parameter.

When  $\lambda = 0$ , the problem reduces to an empirical risk minimization (ERM) problem without regularization.

Assuming that the series in the structure are jointly weakly stationary and  $\hat{\mathbf{Y}}_h^{*'} \hat{\mathbf{Y}}_h^*$  is invertible, it has a closed-form solution given by

$$\hat{\mathbf{G}}_h = \mathbf{B}_h^{*'} \hat{\mathbf{Y}}_h^* (\hat{\mathbf{Y}}_h^{*'} \hat{\mathbf{Y}}_h^*)^{-1},$$

where  $\mathbf{B}_h^* = [\mathbf{b}_{T_1+h}, \dots, \mathbf{b}_T]'$ . If  $\hat{\mathbf{Y}}_h^{*'} \hat{\mathbf{Y}}_h^*$  is not invertible, Ben Taieb & Koo (2019) suggested using a generalized inverse. When  $\lambda > 0$ , imposing the  $L_1$  penalty on  $\mathbf{G}_h$  will introduce sparsity and reduce estimation variance, albeit at the cost of introducing some bias.

Wickramasuriya (2021) proposed an empirical MinT (**EMinT**) solution without the unbiasedness constraint by minimizing the trace of the covariance matrix of the reconciled forecast errors,  $\text{Var}(\tilde{\mathbf{e}}_{T+h|T})$ . Assuming the series are jointly weakly stationary, she derived the solution

$$\hat{\mathbf{G}}_h = \mathbf{B}_h' \hat{\mathbf{Y}}_h (\hat{\mathbf{Y}}_h' \hat{\mathbf{Y}}_h)^{-1},$$

where  $\mathbf{B}_h = [\mathbf{b}_h, \dots, \mathbf{b}_T]'$ , and  $\hat{\mathbf{Y}}_h = [\hat{\mathbf{y}}_{h|0}, \dots, \hat{\mathbf{y}}_{T|T-h}]'$ .

The difference between EMinT and ERM lies in the data sources, as EMinT uses in-sample observations and base forecasts, while ERM relies on observations and base forecasts from a holdout validation set. Both ERM and EMinT consider an estimate of  $\mathbf{G}$  that changes over the forecast horizon, which is why we keep the subscript  $h$  here.

A challenge in forecast reconciliation arises when some base forecasts perform poorly, as the role of the weighting matrix  $\mathbf{G}$  is to assimilate *all* base forecasts and map them into bottom-level disaggregated forecasts, which are subsequently summed by  $\mathbf{S}$ . While the RERM method proposed by [Ben Taieb & Koo \(2019\)](#) introduces sparsity by shrinking some elements of  $\mathbf{G}$  towards zero, it remains incapable of mitigating the adverse impact of underperforming base forecasts. Moreover, the method is time-consuming because it uses expanding windows to recursively generate out-of-sample base forecasts.

We therefore propose two new forecast reconciliation methods involving time series selection: constrained out-of-sample (under the unbiasedness assumption) and unconstrained in-sample (without the unbiasedness assumption). These methods aim to address the negative effect of some poor base forecasts on the overall performance of the reconciled forecasts. Additionally, through the incorporation of regularization in the objective function, our method improves reconciliation outcomes produced with a “poor” choice of  $\mathbf{W}$ .

### 3 Forecast reconciliation with time series selection

In this section, we introduce our methods for forecast reconciliation while automatically avoiding the use of poor base forecasts. Section 3.1 introduces constrained reconciliation methods in a formulation of the problem based on out-of-sample base forecasts, while Section 3.2 presents an unconstrained reconciliation method, where we formulate the problem based on in-sample observations and base forecasts.

### 3.1 Series selection with unbiasedness constraint

As  $\mathbf{S}$  is fixed and  $\hat{\mathbf{y}}_{T+h|T}$  is given,  $\mathbf{G}_h$  determines the linear reconciliation performance, as shown in Equation 1. We drop the subscript  $h$  here as we assume  $\mathbf{W}$  and  $\mathbf{G}$  do not vary with the forecast horizon. A natural way to remove forecasts of some series is by controlling the number of nonzero column entries in  $\mathbf{G}$ . This leads to a generalization of the MinT optimization problem with an additional penalty term:

$$\min_{\mathbf{G}} \quad \frac{1}{2} (\hat{\mathbf{y}} - \mathbf{S}\mathbf{G}\hat{\mathbf{y}})' \mathbf{W}^{-1} (\hat{\mathbf{y}} - \mathbf{S}\mathbf{G}\hat{\mathbf{y}}) + \lambda \mathbf{g}(\mathbf{G}) \quad \text{s.t.} \quad \mathbf{G}\mathbf{S} = \mathbf{I}, \quad (4)$$

where  $\hat{\mathbf{y}} := \hat{\mathbf{y}}_{T+1|T}$ ,  $\mathbf{g}(\cdot)$  penalizes the columns of  $\mathbf{G}$  towards zero, and  $\lambda$  is a penalty parameter. This can be considered a *grouped variable selection problem*, with each group corresponding to a column of  $\mathbf{G}$ . The constraint,  $\mathbf{G}\mathbf{S} = \mathbf{I}$ , ensures that the reconciled forecasts are unbiased if the base forecasts are unbiased (Wickramasuriya et al. 2019). When  $\lambda = 0$ , the problem reduces to the MinT optimization problem in Equation 3 with a closed-form solution given by Equation 2.

**Proposition 1.** *Under the assumption of unbiasedness, the count of nonzero column entries of  $\mathbf{G}$  (i.e., the number of time series selected for reconciliation), derived through solving Equation 4, is at least equal to the number of time series at the bottom level. In addition, we can restore the full hierarchical structure by aggregating/disaggregating the selected time series.*

*Proof.* Under unbiasedness,  $\mathbf{G}\mathbf{S} = \mathbf{I}$ , and

$$\min(\text{rank}(\mathbf{G}), \text{rank}(\mathbf{S})) \geq \text{rank}(\mathbf{I}_{n_b}) = n_b,$$

which indicates that the count of nonzero column entries of  $\mathbf{G}$  is at least equal to  $n_b$ .

Let  $\mathbf{X}_{\mathbb{S}} \in \mathbb{R}^{r \times |\mathbb{S}|}$  denote the submatrix of the  $r \times c$  matrix  $\mathbf{X}$  with column indices forming a set  $\mathbb{S}$  (and when  $\mathbb{S} = \{j\}$ , we simply use  $\mathbf{X}_{\cdot j}$ ), where  $|\mathbb{S}|$  denotes the size of the set  $\mathbb{S}$ . Similarly, let  $\mathbf{X}_{\mathbb{S}} \in \mathbb{R}^{|\mathbb{S}| \times c}$  denote the submatrix of  $\mathbf{X}$  whose rows are indexed by a set  $\mathbb{S}$  (and when  $\mathbb{S} = \{i\}$ , we simply use  $\mathbf{X}_{i \cdot}$ ). Assuming that the set  $\mathbb{S}$  consists of the indices of nonzero columns in the solution to Equation 4,  $\hat{\mathbf{G}}$ , the



following equations hold:

$$GS = \hat{G}_{\cdot\mathbb{S}} S_{\mathbb{S}} = I, \quad \text{and} \quad \min(\text{rank}(\hat{G}_{\cdot\mathbb{S}}), \text{rank}(S_{\mathbb{S}})) \geq \text{rank}(I_{n_b}) = n_b.$$

Additionally, we have  $\text{rank}(S_{\mathbb{S}}) \leq n_b$  as  $S$  has  $n_b$  columns. Therefore, we can conclude that  $\text{rank}(S_{\mathbb{S}}) = n_b$ . Moreover, we have

$$y_t = Sb_t = SGSb_t = S\hat{G}_{\cdot\mathbb{S}} S_{\mathbb{S}} b_t = S\hat{G}_{\cdot\mathbb{S}}(y_t)_{\mathbb{S}},$$

which implies that the hierarchical structure can be fully restored by aggregating/disaggregating the selected time series denoted by  $(y_t)_{\mathbb{S}}$ .

For example, consider the simple hierarchy shown in Figure 1, it is not possible for our constrained reconciliation methods with selection to simultaneously zero out columns of  $G$  associated with series AA and AB. However, it is possible to zero out columns related to series AA and BA simultaneously.

**Proposition 2.** *The optimization problem in Equation 4 can be reformulated as a least squares problem with regularization and linear equality constraint as follows:*

$$\begin{aligned} \min_{\text{vec}(G)} \quad & \frac{1}{2} (\hat{y} - (\hat{y}' \otimes S) \text{vec}(G))' W^{-1} (\hat{y} - (\hat{y}' \otimes S) \text{vec}(G)) + \lambda g(\text{vec}(G)) \\ \text{s.t.} \quad & (S' \otimes I_{n_b}) \text{vec}(G) = \text{vec}(I_{n_b}), \end{aligned} \tag{5}$$

which is characterized as a high-dimensional problem in which the number of features, denoted as  $p = n_b \times n$ , is much larger than the number of observations,  $n$ .

*Proof.* We have

$$\text{vec}(\hat{y}) = \hat{y},$$

$$\text{vec}(SG\hat{y}) = (\hat{y}' \otimes S) \text{vec}(G),$$

$$\text{vec}(GS) = \text{vec}(I_{n_b} GS) = (S' \otimes I_{n_b}) \text{vec}(G).$$

Substituting these into Equation 4, the previous problem now takes the form of a regression problem with an additional regularization term and an equality constraint on the coefficients, as shown in Equation 5.

Next, we present three classes of regularizations that allow forecast reconciliation with series selection, resulting in three optimization problems: (i) group best-subset selection with ridge regularization, (ii) an intuitive method with  $L_0$  regularization, and (iii) a group lasso method.

### Group best-subset selection with ridge regularization

In a high-dimensional context with  $p \gg n$ , it is common to assume that the true regression coefficient (i.e.,  $\text{vec}(\mathbf{G})$  in our problem) is sparse. We apply a combination of  $L_0$  and  $L_2$  regularization to control the nonzero column entries in  $\mathbf{G}$ :

$$\begin{aligned} \min_{\text{vec}(\mathbf{G})} \quad & \frac{1}{2} (\hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G}))' \mathbf{W}^{-1} (\hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G})) + \lambda_0 \sum_{j=1}^n 1(\mathbf{G}_{\cdot j} \neq \mathbf{0}) + \lambda_2 \|\text{vec}(\mathbf{G})\|_2^2 \\ \text{s.t.} \quad & (\mathbf{S}' \otimes \mathbf{I}_{n_b}) \text{vec}(\mathbf{G}) = \text{vec}(\mathbf{I}_{n_b}), \end{aligned} \quad (6)$$

where  $1(\cdot)$  is the indicator function,  $\lambda_0 \geq 0$  controls the number of nonzero columns of  $\mathbf{G}$ ,  $\lambda_2 \geq 0$  controls the strength of the ridge regularization, and  $\|\cdot\|_2$  is the  $L_2$  norm. In a hierarchical or grouped time series context,  $\text{vec}(\mathbf{G})$  has an inherent non-overlapping grouping structure, wherein each group corresponds to a single column of  $\mathbf{G}$ , each of size  $n_b$ . Therefore, we call this reconciliation method *group best-subset selection with ridge regularization*. In the results that follow, we label the **Subset** method differently based on various estimators for  $\mathbf{W}$ , referring to them as **OLS-subset**, **WLSs-subset**, **WLSv-subset**, **MinT-subset**, and **MinTs-subset**, respectively.

The inclusion of the ridge term in Equation 6 is motivated by earlier work on best-subset selection (e.g., [Hazimeh & Mazumder 2020](#), [Mazumder et al. 2022](#)), which suggests that additional ridge regularization can mitigate the poor predictive performance of best-subset selection method in low SNR regimes.

We present a Big-M based mixed integer programming (MIP) formulation for the problem in Equation 6:

$$\min_{\text{vec}(\mathbf{G}), \mathbf{z}, \check{\mathbf{e}}, \mathbf{g}^+} \quad \frac{1}{2} \check{\mathbf{e}}' \mathbf{W}^{-1} \check{\mathbf{e}} + \lambda_0 \sum_{j=1}^n z_j + \lambda_2 \mathbf{g}^{+'} \mathbf{g}^+ \quad (7)$$

$$\begin{aligned}
\text{s.t.} \quad & (\mathbf{S}' \otimes \mathbf{I}_{n_b}) \text{vec}(\mathbf{G}) = \text{vec}(\mathbf{I}_{n_b}) \\
& \hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G}) = \check{\mathbf{e}} \\
& \sum_{i=1}^{n_b} g_{i+(j-1)n_b}^+ \leq \mathcal{M} z_j, \quad j \in [n] \\
& \mathbf{g}^+ \geq \text{vec}(\mathbf{G}) \\
& \mathbf{g}^+ \geq -\text{vec}(\mathbf{G}) \\
& z_j \in \{0, 1\}, \quad j \in [n],
\end{aligned}$$

where  $\mathcal{M}$  is a Big-M parameter (specified a-priori) that is sufficiently large that the optimal solution to Equation 7,  $\mathbf{g}^{+*}$ , satisfies  $\max_{j \in [n]} \sum_{i=1}^{n_b} g_{i+(j-1)n_b}^+ \leq \mathcal{M}$ . The binary variable  $z_j = 0$  implies that  $\mathbf{G}_{\cdot j} = \mathbf{0}$ , and  $z_j = 1$  implies that  $\sum_{i=1}^{n_b} g_{i+(j-1)n_b}^+ \leq \mathcal{M}$ . Such Big-M formulations are commonly used in MIP problems to model relations between discrete and continuous variables, and have been recently explored in regression with  $L_0$  regularization (Bertsimas et al. 2016). The problem is a mixed integer quadratic program (MIQP) that can be solved using commercial MIP solvers, e.g., Gurobi and CPLEX.

**Parameter tuning.** To avoid computationally expensive cross-validation, we tune the parameters to minimize the sum of squared reconciled forecast errors on the truncated training set, comprising only the  $\max\{h, s\}$  observations closest to the forecast origin, where  $s$  is the seasonal period for seasonal data and  $s = T$  for non-seasonal data. Let  $\lambda_0^1 = \frac{1}{2} (\hat{\mathbf{y}} - \tilde{\mathbf{y}}^{\text{bench}})' \mathbf{W}^{-1} (\hat{\mathbf{y}} - \tilde{\mathbf{y}}^{\text{bench}})$ , which captures the scale of the first term in the objective function, where  $\tilde{\mathbf{y}}^{\text{bench}}$  is a vector of reconciled forecasts obtained using Equation 2 with the same estimator of  $\mathbf{W}$ , and define  $\lambda_0^k = 0.0001 \lambda_0^1$ . For the parameter  $\lambda_0$ , we consider a grid of  $k + 1$  values,  $\{\lambda_0^1, \dots, \lambda_0^k, 0\}$ , where  $\lambda_0^j = \lambda_0^1 (\lambda_0^k / \lambda_0^1)^{(j-1)/(k-1)}$  for  $j \in [k]$ . So  $\lambda_0^1, \dots, \lambda_0^k$  is a sequence decreasing on the log scale. We use a grid of six values for the parameter  $\lambda_2$ ,  $\{0, 10^{-2}, 10^{-1}, 10^0, 10^1, 10^2\}$ . Therefore, we tune over a two-dimensional grid of  $(k + 1) \times 6$  values to find the optimal combination of  $\lambda_0$  and  $\lambda_2$ .

**Computation details.** The MIQP problem in Equation 7 is NP-hard and computationally intensive. Bertsimas et al. (2016) showed that commercial MIP solvers are capable of tackling problem instances for

$p$  up to a thousand. To address larger instances, there has been impressive work on developing MIP-based approaches for solving  $L_0$ -regularized regression problem; e.g., [Bertsimas et al. \(2016\)](#), [Hazimeh & Mazumder \(2020\)](#), and [Hazimeh et al. \(2022\)](#). However, it is challenging to extend these approaches to accommodate additional constraints within the optimization problem. Despite the potential sluggishness of handling large instances with commercial MIP solvers, in our experiments, we use Gurobi to solve our problem in Equation 7 by configuring parameters such as MIPGap = 0.001 and TimeLimit = 600 seconds for cases with  $p > 1000$ . This enables us to terminate the solver before reaching the global optimum and return a suboptimal solution instead. This strategy is motivated by our need to consider numerous parameter candidates, and the final solution will be validated against the training set, which prevents the utilization of a poor estimate of  $G$ .

### Intuitive method with $L_0$ regularization

Instead of estimating the entire matrix  $G$  as above, we leverage the MinT solution in Equation 2 to streamline the optimization problem under consideration. Specifically, we define  $\tilde{S} = AS$ , where  $A = \text{diag}(z)$  is an  $n \times n$  diagonal matrix, and  $z$  is an  $n$ -dimensional vector with elements either equal to 0 or 1. Taking the MinT solution in Equation 2, we have  $\tilde{G} = (S'A'W^{-1}AS)^{-1}S'A'W^{-1}$ . Given fixed  $S$  and estimation of  $W$ ,  $\tilde{G}$  is entirely determined by  $A$ . By this way, when the  $j$ th diagonal element of  $A$  equals zero, the  $j$ th column of  $\tilde{G}$  becomes entirely composed of zeros. Therefore, the optimization problem can be reduced to an integer quadratic programming (IQP) problem in which all of the variables are restricted to being integers:

$$\begin{aligned} \min_A \quad & \frac{1}{2} (\hat{y} - S\tilde{G}\hat{y})' W^{-1} (\hat{y} - S\tilde{G}\hat{y}) + \lambda_0 \sum_{j=1}^n A_{jj} \\ \text{s.t.} \quad & \tilde{G} = (S'A'W^{-1}AS)^{-1}S'A'W^{-1} \quad \text{and} \quad \tilde{G}S = I, \end{aligned}$$

where  $\lambda_0 \geq 0$  controls the number of nonzero diagonal elements in  $A$ , consequently affecting the number of nonzero columns (i.e., selected time series) in  $G$ . We call this reconciliation method the *intuitive*

method with  $L_0$  regularization. In the results that follow, we label the **Intuitive** method differently based on various estimators for  $W$ , referring to them as **OLS-intuitive**, **WLSs-intuitive**, **WLSv-intuitive**, **MinT-intuitive**, and **MinTs-intuitive**, respectively.

We should note that implementing grouped variable selection with this optimization problem can be challenging because it imposes restrictions  $\bar{G}$  to ensure it adheres rigorously to the analytical solution of MinT while making the selection. Therefore, the resulting solution tends to be dense and may not have zero columns.

To ensure the invertibility of  $S' A' W^{-1} A S$ , and make the problem compatible with Gurobi, we reformulate the problem as

$$\begin{aligned}
& \min_{A, \bar{G}, C, \check{e}, z} \quad \frac{1}{2} \check{e}' W^{-1} \check{e} + \lambda_0 \sum_{j=1}^n z_j \\
& \text{s.t.} \quad \bar{G} S = I \\
& \quad \hat{y} - (\hat{y}' \otimes S) \text{vec}(\bar{G}) = \check{e} \\
& \quad \bar{G} A S = I \\
& \quad \bar{G} = C S' A' W^{-1} \\
& \quad z_j \in \{0, 1\}, \quad j \in [n].
\end{aligned} \tag{8}$$

**Parameter tuning.** Similar to the setup in the group best-subset selection, we select the tuning parameter,  $\lambda_0$ , by minimizing the sum of squared reconciled forecast errors on a truncated training set, comprising only the  $\max\{h, s\}$  observations that occurred prior to the forecast origin. Let  $\lambda_0^1 = \frac{1}{2} (\hat{y} - \tilde{y}^{\text{bench}})' W^{-1} (\hat{y} - \tilde{y}^{\text{bench}})$ , and  $\lambda_0^k = 0.0001 \lambda_0^1$ , the collection of candidate values for  $\lambda_0$  we consider is  $\{\lambda_0^1, \dots, \lambda_0^k, 0\}$ , where  $\lambda_0^j = \lambda_0^1 (\lambda_0^k / \lambda_0^1)^{(j-1)/(k-1)}$  for  $j \in [k]$ .

**Computation details.** Following a setup akin to that in the group best-subset selection, we employ Gurobi to solve Equation 8 by configuring parameters such as MIPGap = 0.001 and TimeLimit = 600 seconds for problems with  $p > 1000$ .

## Group lasso method

Lasso is another popular method for the selection and estimation of parameters in the context of linear regression. [Yuan & Lin \(2006\)](#) introduced the group lasso method that can be used when there is a grouped structure among the variables. Here, we consider *a group lasso problem under the unbiasedness assumption* given by

$$\begin{aligned} \min_{\mathbf{G}} \quad & \frac{1}{2} (\hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G}))' \mathbf{W}^{-1} (\hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G})) + \lambda \sum_{j=1}^n w_j \|\mathbf{G}_{\cdot j}\|_2 \\ \text{s.t.} \quad & (\mathbf{S}' \otimes \mathbf{I}_{n_b}) \text{vec}(\mathbf{G}) = \text{vec}(\mathbf{I}_{n_b}), \end{aligned} \quad (9)$$

where  $\lambda \geq 0$  is a tuning parameter,  $w_j \neq 0$  is the penalty weight assigned in  $\mathbf{G}_{\cdot j}$  to make the model more flexible, and the second term in the objective is the penalty function that is intermediate between the  $L_1$ -penalty that is used in the lasso and the  $L_2$ -penalty that is used in ridge regression. In the results that follow, we label the **Lasso** method based on various estimators for  $\mathbf{W}$ , referring to them as **OLS-lasso**, **WLSs-lasso**, **WLSv-lasso**, **MinT-lasso**, and **MinTs-lasso**, respectively.

Next, we present the second order cone programming (SOCP) formulation for the group lasso based estimators given by

$$\begin{aligned} \min_{\text{vec}(\mathbf{G}), \check{\mathbf{e}}, \mathbf{g}^+} \quad & \frac{1}{2} \check{\mathbf{e}}' \mathbf{W}_h^{-1} \check{\mathbf{e}} + \lambda \sum_{j=1}^n w_j c_j \\ \text{s.t.} \quad & (\mathbf{S}' \otimes \mathbf{I}_{n_b}) \text{vec}(\mathbf{G}) = \text{vec}(\mathbf{I}_{n_b}) \\ & \hat{\mathbf{y}} - (\hat{\mathbf{y}}' \otimes \mathbf{S}) \text{vec}(\mathbf{G}) = \check{\mathbf{e}} \\ & c_j = \sqrt{\sum_{i=1}^{n_b} g_{i+(j-1)n_b}^2}, \quad j \in [n]. \end{aligned} \quad (10)$$

Equation 10 includes additional auxiliary variables  $c_j \in \mathbb{R}_{\geq 0}$ ,  $j \in [n]$ , and second order cone constraints,

$$c_j = \sqrt{\sum_{i=1}^{n_b} g_{i+(j-1)n_b}^2} \text{ for } j \in [n].$$

Compared to the previous two methods above, the group lasso method is computationally friendlier.

Nonetheless, [Hazimeh et al. \(2023\)](#) demonstrated, both empirically and theoretically, that the group  $L_0$ -regularized method exhibits advantages over its group lasso counterpart across a range of regimes. Group lasso can either be highly dense or possess non-zero coefficients that are overly shrunk. This issue becomes more pronounced when the groups are correlated with each other, as group lasso tends to retain all correlated groups instead of seeking a more concise model.

**Penalty weights and parameter tuning.** In the context of group lasso, the default choice for the penalty weight,  $w_j$ , is  $\sqrt{p_j}$ , where  $p_j$  is the size of each group (in our case,  $p_j = n_b$ ). In our experiments, we allocate different penalty weights to each group using  $w_j = 1/\|\mathbf{G}_{\cdot j}^{\text{bench}}\|_2$ , which allows us to account for variations in scale across different time series in the structure.

We compute the group lasso over  $k+1$  values of the tuning parameter  $\lambda$ , and select the tuning parameter by optimizing the sum of squared reconciled forecast errors on a truncated training set, consisting only of  $\max\{h, s\}$  observations occurred prior to the forecast origin. The collection of candidate values for  $\lambda$  under consideration is  $\{\lambda^1, \dots, \lambda^k, 0\}$ , where  $\lambda^1 = \max_{j=1, \dots, n} \left\| -((\hat{\mathbf{y}}' \otimes \mathbf{S})_{\cdot j^*})' \mathbf{W}^{-1} \hat{\mathbf{y}} \right\|_2 / w_j$ ,  $\lambda^k = 0.0001 \lambda^1$ , and  $\lambda^j = \lambda^1 (\lambda^k / \lambda^1)^{(j-1)/(k-1)}$  for  $j \in [k]$ .

**Proposition 3.** *Ignoring the unbiasedness constraint, we define  $\lambda^1$  as the smallest  $\lambda$  value such that all predictors in the group lasso problem have zero coefficients. Then we have*

$$\lambda^1 = \max_{j=1, \dots, n} \left\| -((\hat{\mathbf{y}}' \otimes \mathbf{S})_{\cdot j^*})' \mathbf{W}^{-1} \hat{\mathbf{y}} \right\|_2 / w_j,$$

where  $j^*$  denotes the column index of  $\hat{\mathbf{y}}' \otimes \mathbf{S}$  that corresponds to the  $j$ th column of  $\mathbf{G}$ .

*Proof.* Denote  $\beta = \text{vec}(\mathbf{G})$ , and the first term in the objective of Equation 9 as  $L(\beta \mid \mathbf{D})$ , where  $\mathbf{D}$  is the working data  $\{\hat{\mathbf{y}}, \hat{\mathbf{y}}' \otimes \mathbf{S}\}$ . Ignoring the unbiasedness constraint, we define  $\lambda^1$  as the smallest  $\lambda$  value such that all predictors in the group lasso problem have zero coefficients, i.e., the solution at  $\lambda^1$  is  $\hat{\beta}^1 = \mathbf{0}$ . (Note that there is no intercept in our problem.) Under the Karush-Kuhn-Tucker conditions, we have

$$\lambda^1 = \max_{j=1, \dots, n} \left\| [\nabla L(\hat{\beta}^1 \mid \mathbf{D})]^{(j)} \right\|_2 / w_j = \max_{j=1, \dots, n} \left\| -((\hat{\mathbf{y}}' \otimes \mathbf{S})_{\cdot j^*})' \mathbf{W}^{-1} \hat{\mathbf{y}} \right\|_2 / w_j.$$

**Computation details.** Due to the incorporation of the unbiasedness constraint, we can not directly use some open-source packages designed for group lasso. Consequently, we employ Gurobi to solve the SOCP problem, configuring it by setting  $\text{OptimalityTol} = 0.0001$ .

### 3.2 Series selection method without unbiasedness constraint

In this section, we relax the unbiasedness constraint,  $\mathbf{G}\mathbf{S} = \mathbf{I}$ , and introduce a reconciliation method with selection that relies on in-sample observations and fitted values. Let  $\mathbf{Y} \in \mathbb{R}^{T \times n}$  denote a matrix comprising observations from all time series on the training set in the structure, and  $\hat{\mathbf{Y}} \in \mathbb{R}^{T \times n}$  denote a matrix of in-sample one-step-ahead forecasts (i.e., fitted values) for all time series. The proposed *empirical group lasso* method considers the optimization problem

$$\min_{\mathbf{G}} \quad \frac{1}{2T} \|\mathbf{Y} - \hat{\mathbf{Y}}\mathbf{G}'\mathbf{S}'\|_F^2 + \lambda \sum_{j=1}^n w_j \|\mathbf{G}_{\cdot j}\|_2,$$

where  $\lambda \geq 0$  is a tuning parameter,  $w_j \neq 0$  is the penalty weight assigned in  $\mathbf{G}_{\cdot j}$  to make a more flexible model. We rewrite the problem as

$$\min_{\text{vec}(\mathbf{G})} \quad \frac{1}{2T} \|\text{vec}(\mathbf{Y}) - (\mathbf{S} \otimes \hat{\mathbf{Y}}) \text{vec}(\mathbf{G}')\|_2^2 + \lambda \sum_{j=1}^n w_j \|\mathbf{G}_{\cdot j}\|_2,$$

which becomes a standard group lasso problem, with  $\text{vec}(\mathbf{Y})$  serving as the dependent variable and  $\mathbf{S} \otimes \hat{\mathbf{Y}}$  as the covariate matrix. We denote this as **Elasso** in the results that follow.

Upon relaxing the unbiasedness constraint, the number of non-zero column entries in the solution for  $\mathbf{G}$  may be less than the number of time series at the bottom level. This differs from the series selection methods with an unbiasedness constraint that we introduced in Section 3.1. In an extreme scenario, it can happen that the solution takes the form of a top-down  $\mathbf{G}_{TD} = [\mathbf{p} \mid \mathbf{O}_{n_b \times (n-1)}]$ , where only the column corresponding to the top level (most aggregated level) retains non-zero values, and  $\mathbf{p} = (p_1, p_2, \dots, p_{n_b})$  is a proportionality vector obtained based on in-sample reconciled forecast errors.



We also explored the empirical version of group best-subset selection with ridge regularization and the intuitive method with  $L_0$  regularization in which we do not impose the unbiasedness constraint. It is worth mentioning that [Hazimeh et al. \(2023\)](#) presented a new algorithmic framework for formulating the group  $L_0$  problem with ridge regularization and provided the **L0Group** Python package for implementation. However, our experiments showed that this algorithm can not terminate within five hours for typical instances with  $p \sim 10^4$ . Therefore, in this paper, we only present the empirical group lasso method for series selection without the unbiasedness constraint.

**Penalty weights and parameter tuning.** Similarly to the setup in the group lasso method, we assign different penalty weights to each group by setting  $w_j = 1/\|\mathbf{G}_{\cdot j}^{\text{OLS}}\|_2$ , where  $\mathbf{G}^{\text{OLS}}$  is the solution obtained by the OLS estimator of  $\mathbf{W}$ . Given a fixed tuning parameter value, we solve the target optimization problem by considering the initial  $T - T_v$  observations, where  $T_v = \max\{h, s\}$  for seasonal time series and  $T_v = \lfloor \frac{1}{10}T \rfloor$  for non-seasonal time series. Then we select the tuning parameter,  $\lambda$ , by minimizing the sum of squared reconciled forecast errors on a truncated training set, comprising only the  $T_v$  observations closest to the forecast origin. Specifically, we form the set of candidate values for  $\lambda$  as  $\{\lambda^1, \dots, \lambda^k, 0\}$ , where  $\lambda^1 = \max_{j=1, \dots, n} \left\| -\frac{1}{N} \left( (\mathbf{S} \otimes \hat{\mathbf{Y}})_{\cdot j*} \right)' \text{vec}(\mathbf{Y}) \right\|_2 / w_j$ ,  $\lambda^k = 0.0001\lambda^1$ , and  $\lambda^j = \lambda^1 (\lambda^k / \lambda^1)^{(j-1)/(k-1)}$  for  $j \in [k]$ . Following the same derivation as in the proof of **Proposition 3**,  $\lambda^1$  is the smallest  $\lambda$  value such that all predictors in the empirical group lasso problem have zero coefficients, i.e.,  $\mathbf{G} = \mathbf{O}$ . Note that we need to resolve the optimization problem based on the whole training set by using the optimal tuning parameter to obtain the final solution.

**Computation details.** While there are open-source packages available for solving group lasso problems, they are still relatively slow when handling large instances. For example, given a specific value for the parameter,  $\lambda$ , our experiments observed that, using the **gglasso** R package, we can not obtain a solution within five hours for typical instances with  $p \sim 10^4$ . Instead, we use Gurobi to solve the problem using the SOCP formulation for the empirical group lasso which aligns with Equation 10 but omits the unbiasedness constraint.

## 4 Monte Carlo simulations

To evaluate the performance of the proposed reconciliation methods with time series selection outlined in Section 3, we carry out two simulations with different designs. In both simulations, we consider a hierarchy comprising two levels of aggregation, as shown in Figure 1. Specifically, the structure has four series at the bottom level, and seven series in total; i.e.,  $n_b = 4$ , and  $n = 7$ . The bottom-level series are first generated and then summed to obtain aggregated series at higher levels.

Section 4.1 considers a setup where the bottom-level series are generated using a structural time series model, but model misspecification exists for some series within the structure. Section 4.2 explores the impact of the correlation between series on the performance of reconciled forecasts.

### 4.1 Setup 1: Exploring the effect of model misspecification

We follow a simulation setup similar to Wickramasuriya et al. (2019), assuming that the bottom-level time series are generated using the basic structural time series model

$$\mathbf{b}_t = \boldsymbol{\mu}_t + \boldsymbol{\gamma}_t + \boldsymbol{\eta}_t,$$

where  $\boldsymbol{\mu}_t$  and  $\boldsymbol{\gamma}_t$  are trend and seasonal components defined by

$$\begin{aligned} \boldsymbol{\mu}_t &= \boldsymbol{\mu}_{t-1} + \mathbf{v}_t + \boldsymbol{\varrho}_t, & \boldsymbol{\varrho}_t &\sim \mathcal{N}(\mathbf{0}, \sigma_{\varrho}^2 \mathbf{I}_4), \\ \mathbf{v}_t &= \mathbf{v}_{t-1} + \boldsymbol{\zeta}_t, & \boldsymbol{\zeta}_t &\sim \mathcal{N}(\mathbf{0}, \sigma_{\zeta}^2 \mathbf{I}_4), \\ \boldsymbol{\gamma}_t &= -\sum_{i=1}^{s-1} \boldsymbol{\gamma}_{t-i} + \boldsymbol{\omega}_t, & \boldsymbol{\omega}_t &\sim \mathcal{N}(\mathbf{0}, \sigma_{\omega}^2 \mathbf{I}_4), \end{aligned}$$

$\boldsymbol{\varrho}_t$ ,  $\boldsymbol{\zeta}_t$ , and  $\boldsymbol{\omega}_t$  are error terms independent of each other and over time, and the error term  $\boldsymbol{\eta}_t$  is generated independently from an  $\text{ARIMA}(p, 0, q)$  process, where  $p$  and  $q$  take values of 0 or 1 with equal probability. The coefficients in the ARIMA process are sampled randomly from a uniform distribution within the range

$[0.5, 0.7]$ , and the contemporaneous error covariance matrix is given by

$$\begin{bmatrix} 5 & 3 & 2 & 1 \\ 3 & 4 & 2 & 1 \\ 2 & 2 & 5 & 3 \\ 1 & 1 & 3 & 4 \end{bmatrix},$$

which enables correlations among time series in a hierarchical structure.

We set  $s = 4$  for quarterly data, with error variances  $\sigma_\rho^2 = 2$ ,  $\sigma_\xi^2 = 0.007$ , and  $\sigma_\omega^2 = 7$ . The initial values for  $\mu_0$ ,  $v_0$ ,  $\gamma_0$ ,  $\gamma_1$ , and  $\gamma_2$  are generated independently from a multivariate normal distribution with zero mean and identity covariance matrix. For each series at the bottom level, we generate a total of  $T + h = 180$  observations, with the last  $h = 16$  observations serving as the test set. Recall that the bottom-level series are aggregated to obtain the data for the aggregated levels. This process is repeated 500 times.

We use ETS models to generate base forecasts for all time series in the hierarchy, using the default settings as implemented in the **forecast** R package (Hyndman et al. 2023). To introduce model misspecification into our experiment, we deliberately undermine the quality of in-sample and out-of-sample forecasts (i.e., fitted values and base forecasts) for some specific time series. Specifically, we investigate three scenarios characterized by artificial model misspecifications, where a 1.5 multiplier is applied to in-sample and out-of-sample forecasts for a single series in each scenario; i.e., series AA at the bottom level, series A at the middle level, and series Total at the top level, resulting in Scenarios A–C, respectively.

The results are summarized in Table 2, Table A1, and Table A2, with each table reporting the average root mean squared error (RMSE) for each level as well as the whole structure (denoted as *Average*). The *Base* row shows the average RMSE of base forecasts, while entries below this row report the percentage decrease (negative) or increase (positive) in the average RMSE of reconciled forecasts compared to base forecasts. The *BU* row uses a “bottom-up” approach, aggregating bottom-level base forecasts to form forecasts for the aggregated series. Notably, in each scenario, the largest improvements occur at the respective hierarchical level where model misspecification is introduced, while slightly deteriorating the performance at other levels.

**Table 2:** Out-of-sample forecast results for the simulated data in Scenario A, Setup 1.

Method	Top				Middle				Bottom				Average			
	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16
Base	9.6	10.7	12.6	15.6	6.3	7.3	8.6	10.8	6.4	7.5	8.3	9.8	6.8	7.9	9.0	10.9
BU	57.8	68.5	53.7	38.9	58.2	61.8	48.1	34.4	0.0	0.0	0.0	0.0	27.0	29.6	23.8	17.7
OLS	0.6	2.2	1.8	1.4	7.1	6.4	4.6	3.1	-7.6	-8.6	-8.2	-7.3	-2.1	-2.5	-2.7	-2.6
OLS-subset	0.6	<b>1.8</b>	<b>1.5</b>	<b>1.3</b>	7.2	<b>5.2</b>	<b>3.8</b>	<b>2.6</b>	<b>-8.3</b>	<b>-12.9</b>	<b>-11.6</b>	<b>-9.9</b>	<b>-2.4</b>	<b>-5.2</b>	<b>-4.8</b>	<b>-4.1</b>
OLS-intuitive	0.8	2.6	2.1	1.8	7.5	<b>6.1</b>	<b>4.4</b>	<b>3.0</b>	<b>-9.0</b>	<b>-12.8</b>	<b>-11.6</b>	<b>-9.9</b>	<b>-2.7</b>	<b>-4.8</b>	<b>-4.5</b>	<b>-3.8</b>
OLS-lasso	0.6	2.2	1.8	1.6	7.4	6.7	4.8	3.2	-7.6	-8.5	-8.1	-7.2	-2.0	-2.4	-2.6	-2.5
WLSs	7.3	10.6	8.1	5.9	15.6	16.0	11.8	8.0	-6.9	-7.8	-7.4	-6.4	1.9	2.0	1.0	0.2
WLSs-subset	<b>5.0</b>	<b>5.7</b>	<b>4.6</b>	<b>3.6</b>	<b>12.3</b>	<b>10.0</b>	<b>7.5</b>	<b>5.2</b>	<b>-7.6</b>	<b>-10.5</b>	<b>-9.6</b>	<b>-8.2</b>	<b>0.2</b>	<b>-2.0</b>	<b>-2.1</b>	<b>-2.0</b>
WLSs-intuitive	<b>7.1</b>	<b>9.2</b>	<b>7.1</b>	<b>5.2</b>	16.5	<b>15.5</b>	<b>11.5</b>	<b>7.9</b>	-6.8	<b>-9.2</b>	<b>-8.4</b>	<b>-7.3</b>	2.1	<b>0.9</b>	<b>0.1</b>	<b>-0.4</b>
WLSs-lasso	7.3	<b>10.3</b>	<b>8.0</b>	5.9	15.7	16.1	11.8	8.1	<b>-7.0</b>	-7.8	-7.3	-6.4	1.9	2.0	1.0	0.2
WLSv	1.0	2.9	2.3	1.9	4.5	4.3	3.2	2.1	-25.8	-26.4	-22.7	-18.3	-12.4	-12.6	-10.7	-8.4
WLSv-subset	<b>-1.0</b>	<b>0.3</b>	<b>0.4</b>	<b>0.5</b>	<b>0.6</b>	<b>0.6</b>	<b>0.5</b>	<b>0.3</b>	<b>-32.3</b>	<b>-32.2</b>	<b>-27.3</b>	<b>-21.7</b>	<b>-17.3</b>	<b>-17.3</b>	<b>-14.2</b>	<b>-10.9</b>
WLSv-intuitive	<b>-0.5</b>	<b>0.2</b>	<b>0.3</b>	<b>0.5</b>	<b>0.9</b>	<b>0.7</b>	<b>0.5</b>	<b>0.3</b>	<b>-32.3</b>	<b>-32.3</b>	<b>-27.4</b>	<b>-21.7</b>	<b>-17.1</b>	<b>-17.3</b>	<b>-14.2</b>	<b>-10.9</b>
WLSv-lasso	<b>0.4</b>	<b>1.5</b>	<b>1.5</b>	<b>1.4</b>	<b>3.0</b>	<b>2.5</b>	<b>2.0</b>	<b>1.3</b>	<b>-28.5</b>	<b>-29.2</b>	<b>-24.9</b>	<b>-19.9</b>	<b>-14.4</b>	<b>-14.9</b>	<b>-12.3</b>	<b>-9.5</b>
MinT	-0.4	0.7	0.9	0.6	0.7	0.7	0.6	0.3	-32.9	-33.4	-28.3	-22.5	-17.5	-17.8	-14.6	-11.3
MinT-subset	<b>-0.6</b>	0.7	<b>0.8</b>	0.7	<b>0.6</b>	0.8	0.6	0.3	<b>-33.0</b>	-33.1	-28.0	-22.3	<b>-17.6</b>	-17.6	-14.5	-11.2
MinT-intuitive	-0.4	0.7	0.9	0.6	0.7	0.7	0.6	0.3	-32.9	-33.4	-28.3	-22.5	-17.5	-17.8	-14.6	-11.3
MinT-lasso	<b>-0.7</b>	<b>0.3</b>	<b>0.6</b>	<b>0.4</b>	<b>0.3</b>	<b>0.4</b>	<b>0.4</b>	<b>0.1</b>	<b>-33.2</b>	<b>-33.7</b>	<b>-28.5</b>	<b>-22.6</b>	<b>-17.8</b>	<b>-18.1</b>	<b>-14.8</b>	<b>-11.4</b>
MinTs	-0.9	0.6	0.7	0.5	0.6	0.6	0.5	0.2	-32.9	-33.5	-28.3	-22.5	-17.6	-17.9	-14.6	-11.3
MinTs-subset	-0.7	0.9	1.1	1.0	0.7	0.8	0.7	0.4	<b>-33.0</b>	-33.1	-27.9	-22.2	-17.6	-17.5	-14.3	-11.0
MinTs-intuitive	-0.9	0.6	0.7	0.5	0.6	0.6	0.5	0.2	-32.9	-33.5	-28.3	-22.5	-17.6	-17.9	-14.6	-11.3
MinTs-lasso	-0.9	<b>0.4</b>	<b>0.6</b>	0.5	0.6	<b>0.4</b>	<b>0.4</b>	<b>0.1</b>	<b>-33.2</b>	<b>-33.6</b>	<b>-28.4</b>	<b>-22.6</b>	<b>-17.7</b>	<b>-18.0</b>	<b>-14.8</b>	<b>-11.4</b>
EMinT	2.2	2.9	2.5	1.7	2.5	2.9	2.3	1.3	-31.9	-32.3	-27.5	-22.0	-15.9	-16.2	-13.4	-10.5
Elasso	<b>1.5</b>	<b>2.8</b>	<b>2.4</b>	1.7	<b>2.1</b>	<b>2.8</b>	2.3	1.3	<b>-32.1</b>	-32.2	-27.4	-21.9	<b>-16.3</b>	-16.2	-13.3	-10.5

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

Focusing on the results of benchmark reconciliation methods, we find that the BU approach performs the best in Scenarios B and C but ranks as the worst overall in Scenario A. This is not surprising, as bottom-level base forecasts are deteriorated in Scenario A, while higher-level base forecasts are deteriorated in Scenarios B and C. Moreover, WLSv, MinT, and MinTs perform especially well in Setup 1, benefiting from their ability to consider in-sample covariance of base forecast errors, allowing for a larger range of adjustments in reconciliation for base forecasts with higher estimated error variance. EMinT also provides accurate reconciled forecasts in our setup, where the in-sample forecasts for specific series are intentionally undermined, a situation that can be detected by the in-sample information based EMinT method. However, OLS and WLSs significantly underperform other benchmark methods in this simulation design.

In all three scenarios, our proposed methods consistently produce either improved or comparable reconciled forecasts compared to their respective benchmarks. The improvements are particularly pronounced when using OLS and WLSs estimators of  $W$  in the benchmark methods, which do not take into account the

in-sample covariance of base forecast errors. One advantage of using our proposed forecast reconciliation methods with selection is their ability to reduce the difference introduced by using different estimates of  $W$ , thereby mitigating the risk of estimator selection. In some cases, such as Scenarios B and C, we can align the forecast accuracy achieved using different estimators, and make them approach the best results we can obtain. Dropping the unbiasedness assumption, Elastic performs similarly to EMinT overall while achieving improvements at the top level, which is typically the aspect of greatest concern to practitioners.

**Table 3:** *Proportion of time series being selected after using the proposed reconciliation methods with selection in Scenario A, Setup 1.*

	Top	A	B	AA	AB	BA	BB	Summary
OLS-subset	0.52	0.79	0.57	0.79	1	0.91	0.85	
OLS-intuitive	0.80	0.90	0.81	0.80	1	0.85	0.86	
OLS-lasso	0.90	1.00	0.68	1.00	1	1.00	1.00	
WLSs-subset	0.85	0.91	0.86	0.90	1	0.97	0.97	
WLSs-intuitive	0.92	0.95	0.67	0.92	1	0.92	0.95	
WLSs-lasso	0.72	1.00	0.72	1.00	1	1.00	1.00	
WLSv-subset	0.50	0.62	0.42	0.19	1	0.81	0.87	
WLSv-intuitive	0.59	0.55	0.49	0.17	1	0.76	0.86	
WLSv-lasso	0.40	1.00	0.41	0.77	1	1.00	1.00	
MinT-subset	0.66	0.90	0.61	0.72	1	0.91	0.93	
MinT-intuitive	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinT-lasso	0.80	0.96	0.84	0.72	1	0.98	0.97	
MinTs-subset	0.57	0.88	0.52	0.67	1	0.89	0.92	
MinTs-intuitive	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinTs-lasso	0.68	1.00	0.66	0.74	1	1.00	1.00	
Elastic	0.82	0.63	0.69	1.00	1	1.00	1.00	

Note: the last column displays a stacked barplot for each method, based on the total number of selected series data from 500 simulation instances, with a darker sub-bar indicating a larger number.

In addition, we report the proportion of time series being selected from the implementation of our proposed methods in 500 simulation instances, as shown in Table 3, Table A3, and Table A4. Clearly, our proposed methods select fewer time series, and generally improve forecast accuracy, compared to the benchmark methods. Furthermore, we observe that the Subset methods tend to return fewer time series compared to the Intuitive and Lasso methods, which aligns with our expectations that the Intuitive and Lasso methods tend to produce dense estimates. Most importantly, depending on the scenario considered, the time series with model misspecification has been selected less often than others. For example, considering Scenario A, series AA is expected to be removed, while AB is expected to be retained. This allows us to obtain

series AA via operations such as A–AB, Total–B–AB, or Total–AB–BA–BB. The results in Table 3 align with our expectations, and show that series AA is dropped often, whereas AB is selected all the time.

## 4.2 Setup 2: Exploring the effect of correlation

We now simulate a hierarchical structure with correlated series, using a similar simulation to Wickramasuriya (2021), and the same hierarchical structure as shown in Figure 1. We use a stationary VAR(1) data generating process for the time series at the bottom level:

$$\mathbf{b}_t = \mathbf{c} + \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix} \mathbf{b}_{t-1} + \boldsymbol{\varepsilon}_t,$$

where  $\mathbf{c}$  is a constant vector with all entries set to 1,  $\mathbf{A}_1$  and  $\mathbf{A}_2$  are  $2 \times 2$  matrices with eigenvalues  $z_{1,2} = 0.6[\cos(\pi/3) \pm i\sin(\pi/3)]$  and  $z_{3,4} = 0.9[\cos(\pi/6) \pm i\sin(\pi/6)]$ , respectively,  $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ , where

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_2 \end{bmatrix}, \quad \text{and} \quad \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \begin{bmatrix} 2 & \sqrt{6}\rho \\ \sqrt{6}\rho & 3 \end{bmatrix},$$

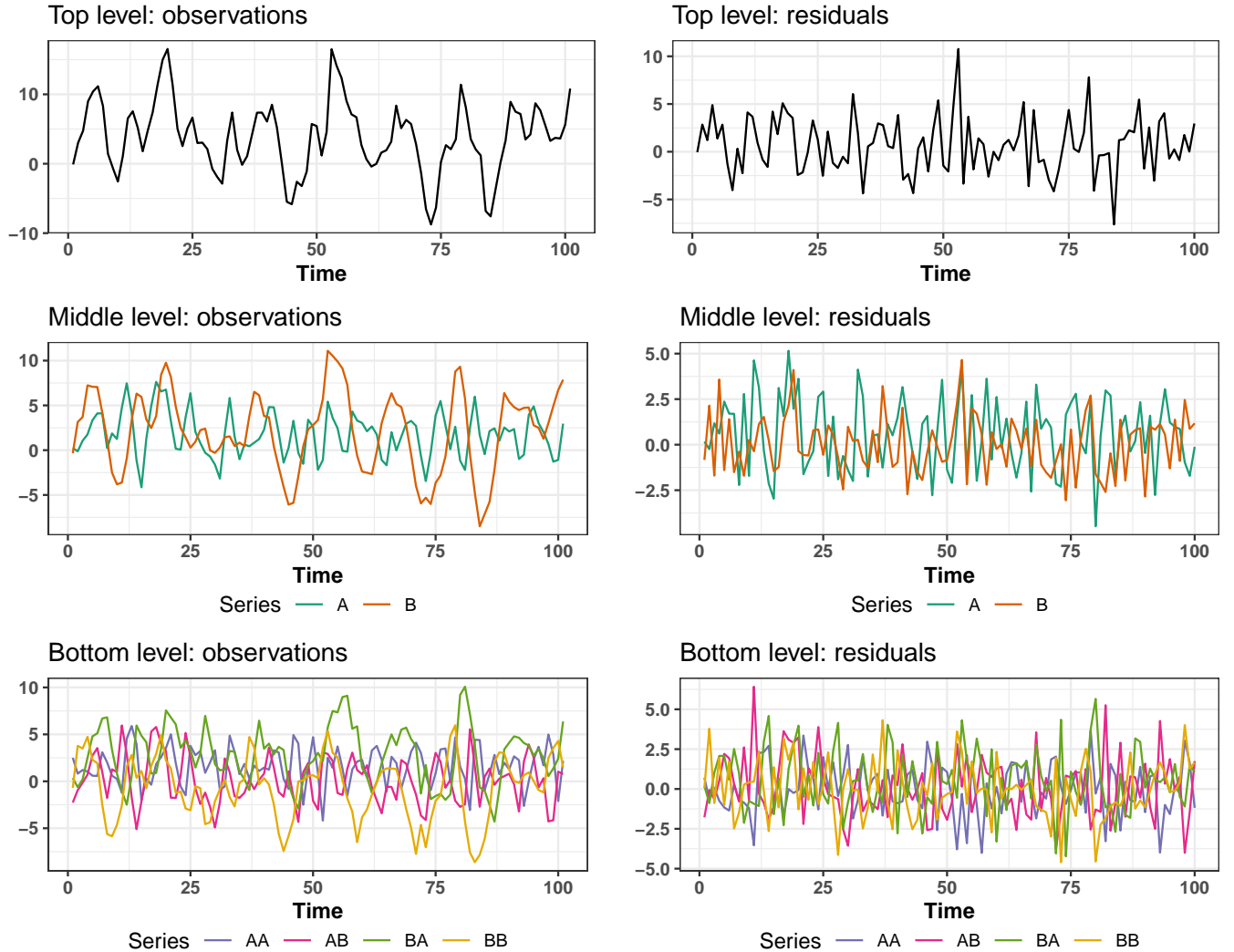
and  $\rho \in \{0, \pm 0.2, \pm 0.4, \pm 0.6, \pm 0.8\}$  controls the error correlation in the simulated hierarchy.

For each time series at the bottom level, we generate a total of 101 observations, with the last observation serving as the test set, i.e.,  $T = 100$  and  $h = 1$ . Once again, the data at the higher levels are obtained by aggregating the bottom-level series. The process is repeated 500 times for each candidate correlation,  $\rho$ .

For each series, base forecasts are generated from ARMA models. We identify the best ARMA model using the automated algorithm implemented in the **forecast** R package (Hyndman & Khandakar 2008). Additionally, when fitting ARMA models for time series Total, A, and BA, we introduce a slight bias by omitting the constant term. Figure 2 presents an illustrative example of a simulated hierarchical time series. The left panels depict time plots for each series at different levels of the structure, while the right panels show the residuals obtained from forecasting each series using the fitted ARMA model. Notably, despite our omission of the constant term when fitting ARMA models to series Total, A, and BA, the

residuals derived from the identified optimal models still exhibit fluctuations around zero and do not display significant deviations in comparison to the residuals from other series. This is because the influence of the constant term is minimal, i.e., it is much smaller compared to the data variability. Thus, it may be challenging to identify the “poor” base forecasts and exclude them from reconciliation in this setup.

Table 4 summarizes the average RMSE of the base forecasts across various error correlations and the percentage relative improvements in RMSE achieved by reconciliation methods relative to the base forecasts. The results show that, for OLS, WLSs, and WLSv estimators, our proposed methods consistently dominate or are equivalent to their respective benchmark methods at all levels. We should highlight the challenge of identifying the “poor” base forecasts in this simulation design, given that the omission of



**Figure 2:** An example hierarchical time series and its in-sample residuals in Setup 2.

**Table 4:** Out-of-sample forecast results across various error correlations for simulation in Setup 2.

Method	Top					Middle					Bottom					Average				
	$\rho=-0.8$	-0.4	0	0.4	0.8	$\rho=-0.8$	-0.4	0	0.4	0.8	$\rho=-0.8$	-0.4	0	0.4	0.8	$\rho=-0.8$	-0.4	0	0.4	0.8
Base	2.4	2.9	3.4	4.1	4.0	1.5	1.8	2.1	2.4	2.5	1.5	1.5	1.5	1.5	1.4	1.6	1.8	2.0	2.1	2.1
BU	-17.0	-9.0	-6.7	-7.0	-7.4	-6.8	0.4	4.8	5.7	2.8	0.0	0.0	0.0	0.0	0.0	-5.3	-1.9	-0.2	-0.1	-1.0
OLS	-11.0	-8.2	-7.7	-8.2	-8.0	-3.5	-0.7	3.1	2.5	0.8	0.7	-0.6	-2.0	-2.3	-2.1	-2.8	-2.4	-1.8	-2.4	-2.7
OLS-subset	<b>-11.4</b>	<b>-8.4</b>	<b>-8.1</b>	<b>-8.4</b>	<b>-8.8</b>	<b>-3.7</b>	-0.7	3.2	2.5	<b>0.4</b>	<b>0.3</b>	<b>-0.8</b>	-2.0	-1.7	<b>-2.6</b>	<b>-3.2</b>	<b>-2.5</b>	<b>-1.9</b>	-2.2	<b>-3.2</b>
OLS-intuitive	<b>-11.6</b>	-8.0	<b>-7.8</b>	-8.0	<b>-8.4</b>	<b>-3.6</b>	-0.4	3.7	2.5	<b>0.3</b>	<b>0.6</b>	-0.2	-1.3	-0.4	-1.5	<b>-3.0</b>	-2.0	-1.3	-1.6	<b>-2.8</b>
OLS-lasso	<b>-19.2</b>	<b>-9.8</b>	-7.2	<b>-8.7</b>	<b>-8.2</b>	<b>-10.5</b>	<b>-1.7</b>	<b>2.9</b>	<b>2.4</b>	0.8	<b>-0.8</b>	<b>-0.8</b>	-1.6	-2.3	-2.1	<b>-7.1</b>	<b>-3.1</b>	-1.6	<b>-2.5</b>	<b>-2.8</b>
WLSs	-16.8	-11.1	-9.6	-10.4	-10.2	-8.1	-2.8	1.5	1.2	-0.4	-0.3	-1.1	-2.4	-2.9	-2.9	-5.7	-3.9	-3.0	-3.6	-4.0
WLSs-subset	<b>-17.3</b>	<b>-11.4</b>	<b>-9.9</b>	<b>-11.1</b>	<b>-10.8</b>	<b>-8.3</b>	-2.8	<b>1.4</b>	<b>0.7</b>	<b>-0.9</b>	<b>-0.7</b>	<b>-1.3</b>	-2.4	<b>-3.2</b>	<b>-3.3</b>	<b>-6.1</b>	<b>-4.0</b>	<b>-3.1</b>	<b>-4.1</b>	<b>-4.5</b>
WLSs-intuitive	<b>-16.9</b>	<b>-11.5</b>	<b>-9.8</b>	-10.0	<b>-10.6</b>	<b>-8.5</b>	-2.8	<b>1.4</b>	1.5	<b>-0.7</b>	<b>-0.7</b>	<b>-1.2</b>	-2.3	-2.7	<b>-3.0</b>	<b>-6.1</b>	<b>-4.0</b>	-3.0	-3.3	<b>-4.3</b>
WLSs-lasso	<b>-18.3</b>	-11.1	-9.2	<b>-10.5</b>	-9.8	<b>-9.3</b>	-2.4	<b>1.4</b>	1.2	-0.1	<b>-0.8</b>	-1.0	-2.4	-2.9	-2.8	<b>-6.6</b>	-3.7	-2.9	<b>-3.7</b>	-3.7
WLSv	-16.5	-11.9	-10.0	-10.6	-10.6	-7.6	-3.4	0.9	1.1	-0.5	-0.5	-1.2	-2.3	-2.9	-3.0	-5.7	-4.3	-3.2	-3.7	-4.2
WLSv-subset	<b>-16.8</b>	<b>-12.1</b>	-9.8	<b>-10.8</b>	<b>-10.7</b>	<b>-7.8</b>	<b>-3.5</b>	1.1	1.2	<b>-1.0</b>	<b>-1.1</b>	<b>-1.3</b>	-2.2	-2.9	<b>-3.2</b>	<b>-6.1</b>	<b>-4.4</b>	-3.0	-3.7	<b>-4.4</b>
WLSv-intuitive	<b>-17.6</b>	<b>-12.6</b>	<b>-10.1</b>	-10.5	-10.6	<b>-8.7</b>	<b>-3.8</b>	<b>0.7</b>	1.1	<b>-0.8</b>	<b>-1.9</b>	<b>-1.5</b>	-2.3	<b>-3.0</b>	-3.0	<b>-7.0</b>	<b>-4.7</b>	<b>-3.3</b>	-3.7	<b>-4.3</b>
WLSv-lasso	<b>-19.8</b>	-11.6	-9.7	-10.5	-10.6	<b>-10.5</b>	-3.0	1.2	1.2	-0.5	<b>-1.2</b>	-1.1	-2.2	-2.9	-3.0	<b>-7.5</b>	-4.1	-3.0	-3.7	-4.2
MinT	-25.4	-18.8	-12.4	<b>-15.3</b>	<b>-12.6</b>	-15.5	-7.0	0.0	-2.0	-2.0	-4.0	-4.6	-4.3	-5.8	-5.1	-11.4	-8.5	-5.0	-7.2	-6.0
MinT-subset	-25.4	-18.8	-12.4	<b>-15.3</b>	<b>-12.6</b>	-15.5	-7.0	0.0	-2.0	-2.0	-4.0	-4.6	-4.3	-5.8	-5.1	-11.4	-8.5	-5.0	-7.2	-6.0
MinT-intuitive	-25.4	-18.8	-12.4	<b>-15.3</b>	<b>-12.6</b>	-15.5	-7.0	0.0	-2.0	-2.0	-4.0	-4.6	-4.3	-5.8	-5.1	-11.4	-8.5	-5.0	-7.2	-6.0
MinT-lasso	-25.4	-18.8	-12.4	<b>-15.3</b>	<b>-12.6</b>	-15.5	-7.0	0.0	-2.0	-2.0	-4.0	-4.6	-4.3	-5.8	-5.1	-11.4	-8.5	-5.0	-7.2	-6.0
MinTs	-25.4	-17.7	-12.1	-14.2	-12.5	-16.1	-6.8	-0.8	-1.6	<b>-2.4</b>	-4.0	-4.6	-4.9	-5.9	<b>-5.2</b>	-11.6	-8.2	-5.4	-6.8	<b>-6.2</b>
MinTs-subset	-25.2	-17.6	-12.1	-14.2	-12.5	-16.1	-6.8	-0.8	-1.6	<b>-2.4</b>	-3.9	-4.6	-4.9	-5.9	<b>-5.2</b>	-11.5	-8.2	-5.4	-6.8	<b>-6.2</b>
MinTs-intuitive	-25.4	-17.7	-12.1	-14.2	-12.5	-16.1	-6.8	-0.8	-1.6	<b>-2.4</b>	-4.0	-4.6	-4.9	-5.9	<b>-5.2</b>	-11.6	-8.2	-5.4	-6.8	<b>-6.2</b>
MinTs-lasso	-25.4	-17.6	-12.1	-14.2	-12.5	-16.1	-6.7	-0.8	-1.6	<b>-2.4</b>	-4.0	-4.6	-4.9	-5.9	<b>-5.2</b>	-11.6	-8.2	-5.4	-6.8	<b>-6.2</b>
EMinT	<b>-31.2</b>	<b>-19.8</b>	<b>-12.5</b>	-14.1	-11.1	<b>-22.9</b>	<b>-10.9</b>	<b>-2.4</b>	<b>-3.2</b>	-1.0	<b>-7.4</b>	<b>-7.3</b>	<b>-6.9</b>	<b>-7.5</b>	-5.1	<b>-16.4</b>	<b>-11.2</b>	<b>-6.9</b>	<b>-7.9</b>	-5.3
Elasso	-31.0	-19.1	-11.1	-13.6	<b>-11.2</b>	-22.7	-9.7	-1.8	-2.4	<b>-1.7</b>	<b>-7.4</b>	-7.2	-6.1	-5.7	-3.5	-16.3	-10.6	-6.0	-6.8	-4.9

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

the constant term has minimal impact relative to the data variability. In addition, we observe that the MinT and MinTs methods perform especially well and our methods provide results similar to these benchmark methods. This is attributed to the use of in-sample covariance by MinT and MinTs, which allows for large adjustments in reconciliation for base forecasts with high estimated error variance. Elasso forecasts are slightly worse than EMinT, possibly due to the difficulty of identifying underperforming base forecasts in this simulation setup. We have also considered alternative error correlation values,  $\rho = -0.6, -0.2, 0.2, 0.4$ , for this simulation setting, but to save space, we do not present all results. The omitted results follow a similar pattern and are available upon request.

In Table 5 and Table A5, we present the proportion of time series being selected by applying our proposed methods. As observed in Table 5, for OLS, WLSs, and WLSv estimators, the Subset and Intuitive methods are still able to exclude the series Total, A, and BA in some instances, in which small biases are introduced in model fitting, while essentially retaining the remaining series in the hierarchy. The Subset methods perform better than the Intuitive method in selection. The Lasso methods typically select all bottom-level series since they tend to yield dense estimates as discussed in Section 3.1. Elasso also selects



all bottom-level series. When dealing with a high positive error correlation, Table A5 shows that our methods still have the potential to do some selection but it becomes somewhat challenging to identify and exclude the series that should be omitted in reconciliation. Hence, our methods are preferred, particularly when the error correlation within the hierarchical structure is negative.

**Table 5:** *Proportion of time series being selected after using the proposed reconciliation methods with selection in Setup 2, with the error correlation being -0.8.*

	Top	A	B	AA	AB	BA	BB	Summary
OLS-subset	0.32	0.34	0.95	0.98	1	0.74	1.00	
OLS-intuitive	0.58	0.52	0.93	0.97	1	0.61	0.97	
OLS-lasso	0.61	0.34	0.38	1.00	1	1.00	1.00	
WLSs-subset	0.27	0.40	0.98	1.00	1	0.73	1.00	
WLSs-intuitive	0.49	0.57	0.96	1.00	1	0.74	0.99	
WLSs-lasso	0.48	0.62	0.72	1.00	1	1.00	1.00	
WLSv-subset	0.30	0.42	1.00	1.00	1	0.68	1.00	
WLSv-intuitive	0.49	0.53	0.99	1.00	1	0.47	1.00	
WLSv-lasso	0.35	0.70	0.85	1.00	1	1.00	1.00	
MinT-subset	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinT-intuitive	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinT-lasso	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinTs-subset	0.87	0.85	1.00	1.00	1	0.85	1.00	
MinTs-intuitive	1.00	1.00	1.00	1.00	1	1.00	1.00	
MinTs-lasso	0.86	0.84	1.00	1.00	1	0.85	1.00	
Elasso	0.94	0.79	0.93	1.00	1	1.00	1.00	

Note: the last column displays a stacked barplot for each method, based on the total number of selected series data from 500 simulation instances, with a darker sub-bar indicating a larger number.

## 5 Applications

In this section we describe two empirical applications: Section 5.1 focuses on a grouped hierarchy built using the Australian labour force survey data released by the Australian Bureau of Statistics, while Section 5.2 considers Australian domestic tourism flows with a natural geographic hierarchy.

### 5.1 Forecasting Australian labour force

The dataset from the Labour Force Survey was released by the Australian Bureau of Statistics, and comprises monthly data on the number of unemployed persons in Australia for the period from January

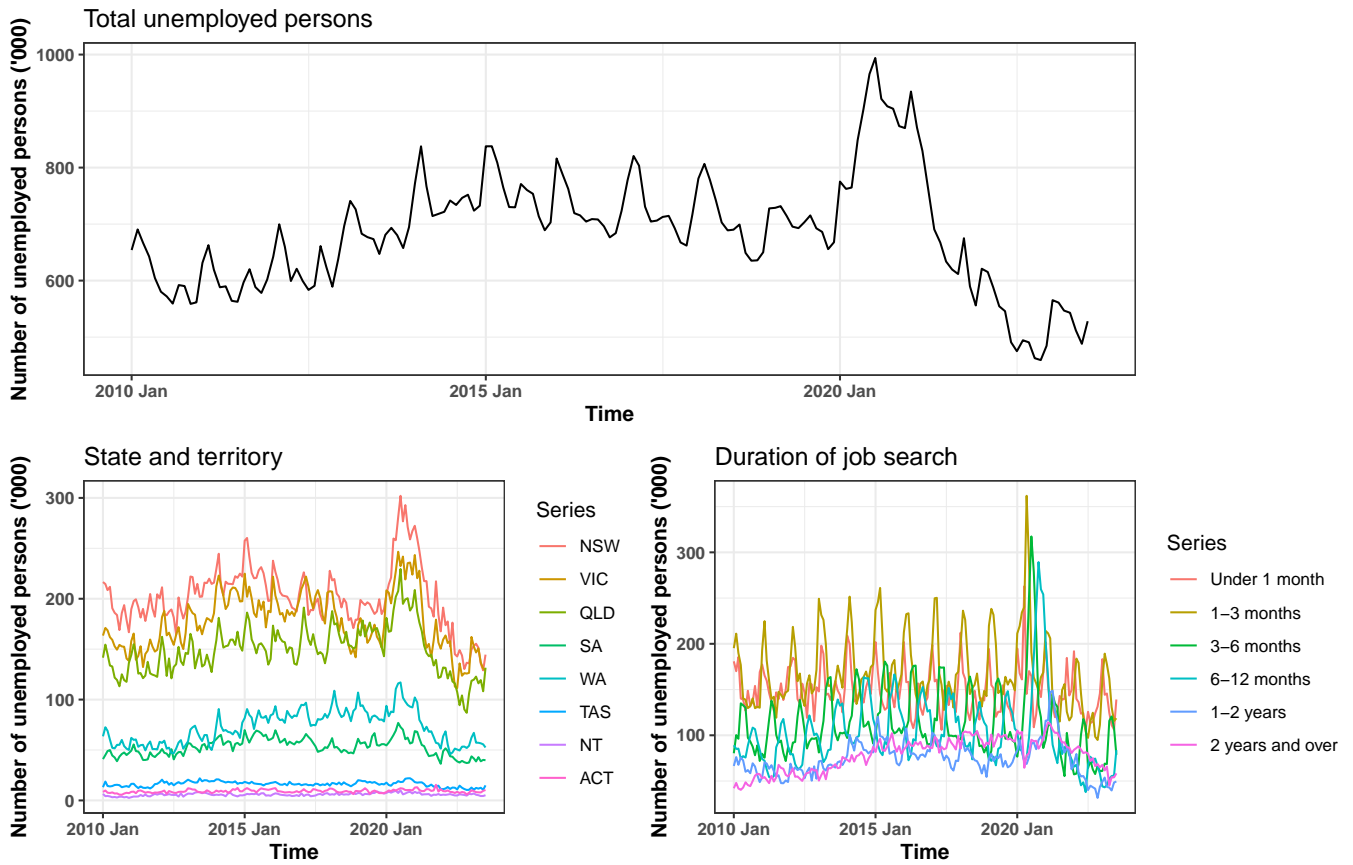
2010 to July 2023<sup>1</sup>. To deal with the few missing observations, we use linear interpolation between observations. Analysis of unemployment data by labour market region and duration of job search can provide valuable insights into regional disparities, and the structural nuances underlying unemployment. Forecast reconciliation is crucial in such a case to ensure aligned decision making.

We construct a grouped hierarchy by disaggregating the number of unemployed persons over two independent attributes, duration of job search (referred to as *Duration*), and State and Territory (referred to as *STT*). At the bottom level, the data are disaggregated by both attributes. We refer to the bottom level as the  $Duration \times STT$  level. Specifically, there are six different groups of job search duration, under 1 month, 1–3 months, 3–6 months, 6–12 months, 1–2 years, and 2 years and over. Additionally, the number of unemployed persons in Australia can be disaggregated by eight states and territories, i.e., NSW (New South Wales), VIC (Victoria), QLD (Queensland), SA (South Australia), WA (Western Australia), TAS (Tasmania), NT (Northern Territory), and ACT (Australian Capital Territory). So the final grouped hierarchy consists of the top series, six series at the Duration level, eight series at the STT level, and 48 series at the  $Duration \times STT$  level, giving 63 time series in total, each of length 163 observations.

The top panel in Figure 3 shows the total number of unemployed persons in Australia from January 2010 to July 2023, representing the top-level series in the hierarchical structure. The monthly series shows strong seasonality within each year, marked by prominent peaks occurring every January, attributable to school-leavers. Lower peaks occur in July, perhaps impacted by the start of the financial year. Amidst the backdrop of COVID-19’s non-essential service shutdowns and trading restrictions, March and April 2020 saw a notable surge in unemployment. However, as coronavirus cases dwindled significantly and restrictions eased in the aftermath, employment made a remarkable recovery, leading to a subsequent decline in unemployment. The bottom-left panel displays the breakdown of unemployed individuals by state and territory, while the bottom-right panel presents the breakdown by the duration of job search. The plots display diverse and rich dynamics both within and between different levels of the hierarchy. For

---

<sup>1</sup>The Labour Force Survey data is publicly available at <https://www.abs.gov.au/statistics/labour/employment-and-unemployment/labour-force-australia-detailed/aug-2023>.



**Figure 3:** *Australia unemployed persons, disaggregated by state and territory, and by duration of job search.*

example, there was noticeable growth observed during 2020 for some states such as NSW, VIC, and QLD, whereas other states did not experience such significant growth. Additionally, there is a resemblance in the seasonal patterns between NSW and QLD, while the seasonal pattern in VIC differs. When comparing the series at the STT level and Duration level, the seasonal patterns in the Duration-level series are more consistent and potentially easier to forecast.

We assess the forecast accuracy of base forecasts and various reconciliation methods through a rolling forecast origin approach. Our aim is to generate 1- to 12-step-ahead forecasts for each of the 63 series while ensuring coherence. Given the limited data compared to the forecast horizon, we initiate the process with a training set of 139 observations for each series. The training set is used to select the optimal ETS model with the automatic algorithm implemented in the **forecast** package for R (Hyndman & Khandakar 2008). Using these fitted ETS models, we generate base forecasts, and perform diverse

forecast reconciliation methods. Then we roll the forecast origin forward by one month and repeat the process, until July 2022. We note that it may be challenging to identify the series with “poor” forecasts due to structural changes in the data caused by the COVID-19 pandemic, which affect the accuracy of forecasts across all time series.

**Table 6:** Average out-of-sample forecast results for Australian labour force data.

Method	Top				Duration				STT				Duration x STT				Average			
	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12
Base	29.4	44.9	58.6	67.6	10.1	14.2	16.3	18.1	6.6	8.4	9.9	10.7	2.3	2.9	3.1	3.3	4.0	5.3	6.1	6.6
BU	46.7	34.1	29.3	24.2	7.4	2.3	0.8	0.8	5.1	9.8	10.5	10.4	0.0	0.0	0.0	0.0	8.4	7.1	6.8	6.3
OLS	2.0	1.7	1.5	1.0	0.6	-4.2	-4.3	-3.3	-0.7	0.4	0.0	-0.1	1.9	0.7	0.8	0.7	1.0	-0.5	-0.6	-0.5
OLS-subset	2.1	<b>1.0</b>	<b>-1.2</b>	<b>-2.0</b>	0.6	-4.1	<b>-5.2</b>	<b>-4.4</b>	<b>-1.0</b>	0.6	<b>-0.6</b>	<b>-1.2</b>	1.9	0.8	<b>0.3</b>	<b>0.2</b>	1.0	-0.5	<b>-1.5</b>	<b>-1.6</b>
OLS-intuitive	<b>-1.3</b>	<b>1.5</b>	<b>1.0</b>	<b>0.2</b>	<b>-0.9</b>	-3.9	-4.2	-3.2	<b>-1.3</b>	0.8	0.5	0.6	<b>1.8</b>	1.2	1.3	1.2	<b>0.1</b>	-0.1	-0.3	-0.2
OLS-lasso	2.0	1.7	1.5	1.0	0.6	-4.2	-4.3	-3.3	-0.7	0.4	0.0	-0.1	1.9	0.7	0.8	0.7	1.0	-0.5	-0.6	-0.5
WLSs	16.9	10.8	9.0	7.0	-0.7	-4.8	-5.0	-4.3	-3.0	0.6	1.4	1.6	-1.3	<b>-1.7</b>	<b>-1.5</b>	<b>-1.6</b>	0.6	-0.3	-0.2	-0.3
WLSs-subset	<b>14.9</b>	<b>9.7</b>	<b>6.4</b>	<b>4.2</b>	<b>-1.6</b>	-3.7	-4.1	-3.4	-2.3	1.5	1.6	<b>1.0</b>	-0.6	-0.2	0.5	0.3	0.6	0.6	0.4	0.1
WLSs-intuitive	16.9	10.8	9.0	7.0	-0.7	-4.8	-5.0	-4.3	-3.0	0.6	1.4	1.6	-1.3	<b>-1.7</b>	<b>-1.5</b>	<b>-1.6</b>	0.6	-0.3	-0.2	-0.3
WLSs-lasso	16.9	10.8	9.0	7.0	-0.7	-4.8	-5.0	-4.3	-3.0	0.6	1.4	1.6	-1.3	<b>-1.7</b>	<b>-1.5</b>	<b>-1.6</b>	0.6	-0.3	-0.2	-0.3
WLSv	15.6	9.8	8.4	6.6	1.1	-5.0	<b>-5.6</b>	-4.5	-2.6	-0.5	-0.1	0.2	<b>-1.5</b>	-1.5	<b>-1.5</b>	-1.3	0.9	-0.7	-0.8	-0.6
WLSv-subset	<b>10.2</b>	<b>5.1</b>	<b>2.9</b>	<b>1.8</b>	<b>-0.6</b>	<b>-5.4</b>	-5.5	<b>-4.6</b>	-1.8	<b>-0.8</b>	<b>-1.1</b>	<b>-1.0</b>	-1.1	-1.1	-0.7	-0.5	<b>0.2</b>	<b>-1.3</b>	<b>-1.5</b>	<b>-1.3</b>
WLSv-intuitive	15.6	9.8	8.4	6.6	1.1	-5.0	<b>-5.6</b>	-4.5	-2.6	-0.5	-0.1	0.2	<b>-1.5</b>	-1.5	<b>-1.5</b>	-1.3	0.9	-0.7	-0.8	-0.6
WLSv-lasso	15.6	9.8	8.4	6.6	1.1	-5.0	<b>-5.6</b>	-4.5	-2.6	-0.5	-0.1	0.2	<b>-1.5</b>	-1.5	<b>-1.5</b>	-1.3	0.9	-0.7	-0.8	-0.6
MinTs	9.9	5.8	6.4	5.3	0.6	-5.0	-5.0	-3.7	-3.4	-2.4	-1.4	-0.8	-0.4	-0.9	-0.8	-0.7	0.4	<b>-1.3</b>	-0.9	-0.5
MinTs-subset	<b>9.1</b>	6.1	6.5	<b>4.5</b>	<b>0.3</b>	-4.8	<b>-5.1</b>	<b>-3.9</b>	<b>-3.6</b>	-2.4	-1.3	<b>-0.9</b>	<b>-0.6</b>	-0.8	-0.7	-0.7	<b>0.1</b>	-1.2	-0.9	<b>-0.7</b>
MinTs-intuitive	9.9	5.8	6.4	5.3	0.6	-5.0	-5.0	-3.7	-3.4	-2.4	-1.4	-0.8	-0.4	-0.9	-0.8	-0.7	0.4	<b>-1.3</b>	-0.9	-0.5
MinTs-lasso	9.9	5.8	6.4	5.3	0.6	-5.0	-5.0	-3.7	-3.4	-2.4	-1.4	-0.8	-0.4	-0.9	-0.8	-0.7	0.4	<b>-1.3</b>	-0.9	-0.5
EMinT	43.1	17.6	9.9	10.2	36.8	25.2	27.9	24.1	16.8	15.1	6.2	6.2	32.3	27.9	29.8	27.9	31.4	23.3	21.4	19.6
Elasso	<b>-5.8</b>	<b>-2.0</b>	<b>-2.5</b>	<b>-2.3</b>	<b>33.5</b>	<b>11.1</b>	<b>1.1</b>	<b>-3.6</b>	<b>-17.4</b>	<b>-8.9</b>	<b>-10.0</b>	<b>-8.8</b>	<b>20.6</b>	<b>6.4</b>	<b>2.6</b>	<b>0.5</b>	<b>12.7</b>	<b>3.4</b>	<b>-1.1</b>	<b>-2.9</b>

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

The average results are presented in Table 6. The poor performance of the MinT method and associated methods can be attributed to the poor sample covariance estimator when the sample size is only slightly larger than the number of series in the structure. The Subset methods using different estimators of  $G$  generally improve forecast accuracy over their benchmark methods, particularly when focusing on aggregation levels, which are typically of paramount concern to practitioners. The only exception is the WLSs-subset method, which returns reduced accuracy for longer horizons. However, it still demonstrates improvements in top-level forecasts. Moreover, the Intuitive and Lasso methods almost always yield results identical to the corresponding benchmark methods, because they tend to provide dense estimates, and ETS models typically do not result in extremely poor forecasts. The only exception is OLS-intuitive, which shows improved forecast accuracy at the top level but deterioration at other levels. When we drop the unbiasedness assumption, EMinT is the worst-performing method across all levels because it relies on

the assumption that the series in the hierarchy are jointly weakly stationary, which is evidently not the case in the application. Elasso significantly improves the quality of forecasts over EMinT, with the most accurate coherent forecasts observed at the top level and STT level. Overall, Elasso performs well for longer forecast horizons, but it is less effective for one-step-ahead forecasts.

**Table 7:** Out-of-sample forecast results on a single test set (from August 2022 to July 2023) for Australian labour force data.

Method	Top				Duration				STT				Duration x STT				Average			
	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12
Base	18.5	13.6	18.3	28.3	11.8	12.7	13.9	16.9	6.7	6.0	6.0	6.3	2.3	2.6	2.7	2.9	4.1	4.1	4.4	5.1
BU	<b>-81.5</b>	33.4	-19.9	-45.0	<b>-30.7</b>	-9.2	-7.9	-10.1	-12.9	-10.4	-13.4	-13.5	0.0	0.0	0.0	0.0	<b>-17.1</b>	-2.8	-5.9	-9.3
OLS	-16.2	-14.2	-13.4	-10.4	2.5	-2.6	-2.7	-0.6	-1.8	-0.9	-1.9	0.3	6.7	5.1	5.1	4.9	2.1	0.7	0.4	1.1
OLS-subset	<b>-17.0</b>	-2.1	<b>-31.2</b>	<b>-38.4</b>	<b>2.0</b>	-1.7	<b>-5.0</b>	<b>-2.7</b>	<b>-2.7</b>	<b>-4.2</b>	<b>-8.6</b>	<b>-7.3</b>	<b>6.7</b>	5.2	<b>3.3</b>	<b>3.7</b>	<b>1.7</b>	1.1	<b>-3.5</b>	<b>-3.8</b>
OLS-intuitive	<b>-79.6</b>	<b>-23.9</b>	<b>-31.9</b>	<b>-32.1</b>	<b>-13.0</b>	0.4	-0.8	0.3	<b>-8.9</b>	4.9	7.0	13.2	<b>6.3</b>	12.3	12.4	11.6	<b>-8.5</b>	5.6	4.7	4.4
OLS-lasso	-16.2	-14.2	-13.4	-10.4	2.5	-2.6	-2.7	-0.6	-1.8	-0.9	-1.9	0.3	6.7	5.1	5.1	4.9	2.1	0.7	0.4	1.1
WLSs	-60.6	-29.4	-44.0	-38.6	-12.0	-7.6	-6.9	-5.9	-6.5	-8.1	-9.4	-8.0	3.2	1.7	1.7	1.6	-7.7	-4.4	-5.8	-5.9
WLSs-subset	<b>-61.6</b>	-22.4	<b>-47.3</b>	<b>-50.4</b>	-12.0	<b>-8.0</b>	<b>-10.5</b>	<b>-7.8</b>	<b>-6.6</b>	<b>-10.7</b>	<b>-14.3</b>	<b>-12.8</b>	3.2	5.6	4.1	5.9	<b>-7.8</b>	-2.8	<b>-6.7</b>	<b>-6.4</b>
WLSs-intuitive	-60.6	-29.4	-44.0	-38.6	-12.0	-7.6	-6.9	-5.9	-6.5	-8.1	-9.4	-8.0	3.2	1.7	1.7	1.6	-7.7	-4.4	-5.8	-5.9
WLSs-lasso	-60.6	-29.4	-44.0	-38.6	-12.0	-7.6	-6.9	-5.9	-6.5	-8.1	-9.4	-8.0	3.2	1.7	1.7	1.6	-7.7	-4.4	-5.8	-5.9
WLSv	-60.6	-29.1	-41.4	-36.6	-14.5	-8.7	-5.6	-4.8	-3.3	-6.8	-8.0	-7.0	5.5	2.6	2.6	3.1	-6.7	-4.0	-4.5	-4.6
WLSv-subset	-51.6	<b>-32.7</b>	-36.6	-29.6	<b>-18.3</b>	<b>-9.8</b>	<b>-10.5</b>	<b>-10.9</b>	-1.1	-4.3	<b>-8.1</b>	<b>-7.3</b>	<b>2.5</b>	<b>2.1</b>	<b>2.3</b>	<b>1.8</b>	<b>-7.9</b>	<b>-4.3</b>	<b>-5.8</b>	<b>-6.5</b>
WLSv-intuitive	-60.6	-29.1	-41.4	-36.6	-14.5	-8.7	-5.6	-4.8	-3.3	-6.8	-8.0	-7.0	5.5	2.6	2.6	3.1	-6.7	-4.0	-4.5	-4.6
WLSv-lasso	-60.6	-29.1	-41.4	-36.6	-14.5	-8.7	-5.6	-4.8	-3.3	-6.8	-8.0	-7.0	5.5	2.6	2.6	3.1	-6.7	-4.0	-4.5	-4.6
MinTs	-27.0	-21.1	-22.9	-21.8	-9.1	-7.9	-6.7	-4.6	-3.6	-9.0	-10.5	-7.6	7.7	3.7	3.0	3.4	-1.9	-3.3	-3.9	-3.1
MinTs-subset	<b>-41.4</b>	-9.3	-17.8	<b>-45.1</b>	<b>-12.2</b>	-5.0	<b>-8.2</b>	<b>-6.8</b>	<b>-6.1</b>	<b>-9.8</b>	-7.1	<b>-9.3</b>	<b>5.3</b>	4.7	3.8	3.6	<b>-5.3</b>	-1.5	-3.0	<b>-6.1</b>
MinTs-intuitive	-27.0	-21.1	-22.9	-21.8	-9.1	-7.9	-6.7	-4.6	-3.6	-9.0	-10.5	-7.6	7.7	3.7	3.0	3.4	-1.9	-3.3	-3.9	-3.1
MinTs-lasso	-27.0	-21.1	-22.9	-21.8	-9.1	-7.9	-6.7	-4.6	-3.6	-9.0	-10.5	-7.6	7.7	3.7	3.0	3.4	-1.9	-3.3	-3.9	-3.1
EMinT	-60.4	-14.0	1.4	-29.9	-6.0	12.0	10.7	-6.7	16.7	-0.9	-12.4	<b>-21.0</b>	23.3	17.2	16.7	10.1	7.7	10.8	9.0	-3.7
Elasso	-4.2	-3.3	<b>-22.3</b>	-8.0	<b>-19.7</b>	<b>-9.9</b>	<b>-19.9</b>	<b>-25.3</b>	<b>-24.6</b>	<b>-24.3</b>	<b>-22.6</b>	-14.6	<b>-10.8</b>	<b>-3.8</b>	<b>-0.2</b>	<b>-4.9</b>	<b>-15.7</b>	<b>-9.3</b>	<b>-11.4</b>	<b>-13.2</b>

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

We provide results based on the final test set spanning from August 2022 to July 2023 in Table 7. This is the latest available data, enabling us to use more data for model training and explore the post-COVID pattern. All Subset methods (using various estimators of  $W$ ) produce improved or comparable reconciled forecasts compared to their benchmarks. The accuracy improvements become more noticeable for longer forecast horizons. Similar to the average results in Table 6, the Intuitive and Lasso methods yield results identical to the benchmark methods due to their tendency to offer dense estimates. Surprisingly, when relaxing the unbiasedness constraint, Elasso ranks the best and demonstrates significant improvement over EMinT, and outperforms other methods across almost all levels except for the top level.

Table 8 presents the number of series selected at each level and the optimal tuning parameter values obtained using different proposed methods. We only show results from the Subset and Elasso methods, as they had the best RMSE results. The variation in the scale of the optimal parameters for different methods

**Table 8:** Number of time series selected using different proposed methods and the optimal parameter values identified in the labour application, considering a single test set (from August 2022 to July 2023). The None row shows the original number of series in the structure.

	Number of time series retained					Optimal parameters		
	Top	Duration	STT	Duration x STT	Total	$\lambda$	$\lambda_0$	$\lambda_2$
None	1	6	8		48	63	-	-
OLS-subset	0	5	1		48	54	-	4.16
WLSs-subset	0	5	1		46	52	-	0.38
WLSv-subset	1	5	7		48	61	-	0.51
MinTs-subset	0	1	1		47	49	-	0.03
Elasso	1	5	2		3	11	213.59	-

comes from the difference in the scales of the objective function. Table 8 shows that all Subset methods exclude some series. Remarkably, the Elasso method consistently outperforms the others overall, even though it uses only 11 series for forecast reconciliation. Most of the series at the STT level are removed, while the majority of series at the Duration level are retained. This aligns with our data description, highlighting that the seasonal patterns in the Duration level series are more consistent and potentially easier to forecast compared to those at the STT level.

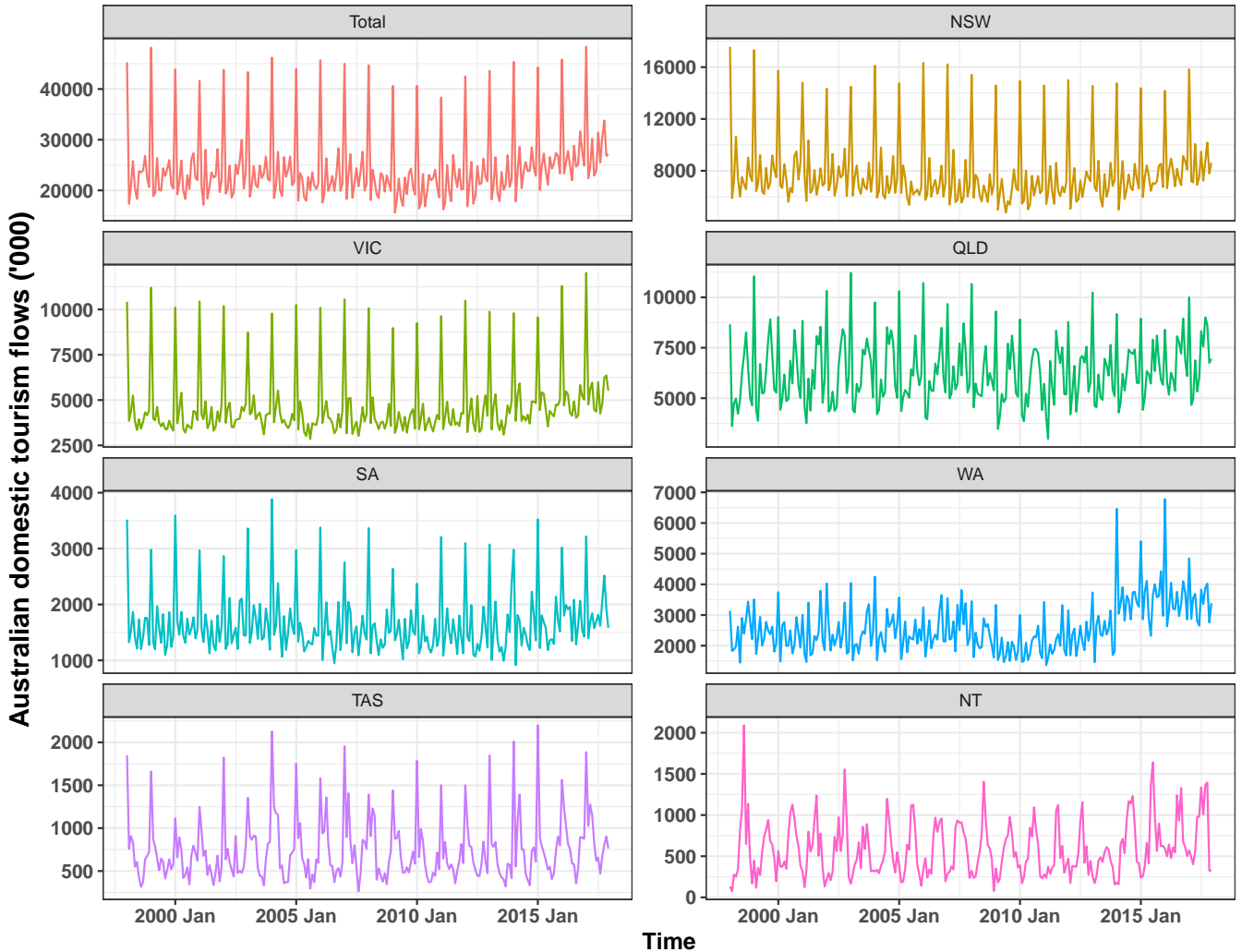
## 5.2 Forecasting Australian domestic tourism

Australian domestic tourism flows are measured as the number of overnight trips Australians spend away from home. The data are sourced from the National Visitor Survey and collected through computer-assisted telephone interviews involving approximately 120,000 Australian residents aged 15 years and older. The data follow a geographic structure, with national total tourism flows in the top-level aggregation, then disaggregated into seven states and territories (referred to as *State* level hereafter), further dividing into 27 zones, and finally, into 76 regions. Thus,  $n_b = 76$  and  $n = 111$ . Each series in the hierarchy spans the period from January 1998 to December 2017, with a total of 240 observations.

Figure 4 shows the aggregate tourism flows for Australia as well as individual states, revealing pronounced seasonal patterns across the national total and states, albeit with varying seasonal patterns among the series. Notably, there was significant growth starting from around 2010 for the national total flow and some states such as NSW, VIC, QLD, and WA. While flows are relatively flat for SA, TAS, and NT. Moreover, the

time plot displays that a large decrease in tourism flows for WA occurred in 2016.

Our objective is to forecast tourism flows for each series in the geographic hierarchy while ensuring coherence across all levels. We use a rolling forecast origin to evaluate the forecast accuracy of different methods. We start with a training set of 216 months for each series, and compute base forecasts from optimal ETS models. We then roll the forecast origin forward, month by month, until December 2016. The base forecasts are reconciled using our proposed methods and some existing reconciliation methods. Table 9 reports the average RMSE values for base forecasts generated by ETS models, along with the percentage relative improvements obtained by each reconciliation method. Similar to Section 5.1, the MinT method and the respective proposed methods are not considered due to their poor performance. The



**Figure 4:** Domestic tourism flows from January 1998 to December 2017 for the whole of Australia as well as the states.



results show that the OLS method outperforms other benchmark methods like WLSs, WLSv and MinTs, despite the fact that WLSv and MinTs account for the in-sample covariance of base forecast errors. This highlights the effectiveness of the OLS method despite its simplicity.

Overall, the Subset methods outperform their respective benchmark methods, especially for aggregation levels and longer forecast horizons. The only exception is the OLS-subset method, which slightly reduces overall accuracy while still improving top-level forecasts. Moreover, the Intuitive and Lasso methods produce results almost identical to the corresponding benchmark methods, which is not surprising as ETS models typically do not yield extremely poor forecasts, making them challenging to be selected out using methods that tend to return dense estimates. When we relax the unbiasedness constraint, EMinT consistently performs the worst across all levels due to the evident lack of joint weak stationarity among the series in the hierarchy. The Elasso method presents significant improvement compared to the EMinT method, and it also outperforms other methods across almost all levels except for the bottom level.

**Table 9:** Average out-of-sample forecast results for Australian domestic tourism data.

Method	Top				State				Zone				Region				Average			
	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12
Base	1565.8	1520.2	1548.3	1773.1	366.0	406.1	421.4	442.0	142.5	170.8	178.5	185.3	72.1	86.4	90.9	94.4	121.2	140.0	146.2	153.6
BU	14.3	38.8	42.0	38.8	4.6	10.3	13.8	15.7	-0.1	0.9	1.3	1.6	0.0	0.0	0.0	0.0	2.5	6.0	6.9	7.4
OLS	-0.6	1.1	1.8	1.9	-1.2	-1.0	-1.0	-1.3	-2.8	-4.0	-4.8	-5.6	-0.1	-0.8	-1.6	-2.4	-1.1	-1.6	-2.1	-2.7
OLS-subset	-0.6	<b>-1.9</b>	<b>-4.9</b>	<b>-3.0</b>	-1.2	<b>-1.2</b>	1.5	0.9	-2.6	-0.5	-0.2	-1.1	0.2	2.1	1.7	0.6	-1.0	0.3	0.5	-0.2
OLS-intuitive	-0.6	1.1	1.8	1.9	-1.2	-1.0	-1.0	-1.3	-2.8	-4.0	-4.8	-5.6	-0.1	-0.8	-1.6	-2.4	-1.1	-1.6	-2.1	-2.7
OLS-lasso	<b>-0.8</b>	2.1	2.8	2.9	<b>-1.3</b>	-0.4	0.5	0.3	-2.3	-3.5	-4.2	-4.9	0.0	<b>-0.9</b>	-1.5	-2.2	-1.0	-1.3	-1.5	-2.0
WLSs	4.1	16.5	19.0	18.1	0.6	2.0	4.1	5.2	-2.7	-3.1	-3.3	-3.4	-0.5	-1.0	-1.4	-1.8	-0.4	0.7	1.0	1.1
WLSs-subset	4.1	<b>6.9</b>	<b>8.9</b>	<b>10.2</b>	0.6	<b>1.7</b>	<b>4.0</b>	<b>4.3</b>	-2.7	-3.0	-2.1	-2.1	-0.5	-0.1	-0.1	-0.5	-0.4	<b>0.0</b>	<b>0.9</b>	<b>1.0</b>
WLSs-intuitive	4.1	16.5	19.0	18.1	0.6	2.0	4.1	5.2	-2.7	-3.1	-3.3	-3.4	-0.5	-1.0	-1.4	-1.8	-0.4	0.7	1.0	1.1
WLSs-lasso	<b>3.6</b>	17.1	19.5	18.5	<b>0.3</b>	2.1	4.3	5.5	-2.7	-2.9	-3.2	-3.3	<b>-0.6</b>	-1.0	-1.4	-1.7	<b>-0.5</b>	0.8	1.1	1.2
WLSv	6.2	22.3	25.0	23.6	1.1	3.7	6.3	7.9	-2.6	-2.4	-2.4	-2.2	-1.6	-1.6	-1.8	-1.9	-0.5	1.5	2.1	2.4
WLSv-subset	<b>0.9</b>	<b>4.5</b>	<b>4.9</b>	<b>5.7</b>	<b>-1.2</b>	<b>-0.3</b>	<b>-0.1</b>	<b>0.7</b>	<b>-3.2</b>	<b>-3.8</b>	<b>-4.5</b>	<b>-4.9</b>	-1.3	-1.2	-1.7	<b>-2.3</b>	<b>-1.6</b>	<b>-1.2</b>	<b>-1.6</b>	<b>-1.7</b>
WLSv-intuitive	6.2	22.3	25.0	23.6	1.1	3.7	6.3	7.9	-2.6	-2.4	-2.4	-2.2	-1.6	-1.6	-1.8	-1.9	-0.5	1.5	2.1	2.4
WLSv-lasso	6.2	22.3	25.0	23.6	1.1	3.7	6.3	7.9	-2.6	-2.4	-2.4	-2.2	-1.6	-1.6	-1.8	-1.9	-0.5	1.5	2.1	2.4
MinTs	5.1	17.2	19.7	18.6	0.1	1.9	4.4	5.7	-3.5	-3.3	-3.4	-3.4	<b>-1.9</b>	<b>-2.0</b>	<b>-2.4</b>	<b>-2.7</b>	-1.2	0.2	0.7	0.9
MinTs-subset	<b>1.8</b>	<b>1.3</b>	<b>2.0</b>	<b>3.2</b>	<b>-2.2</b>	<b>-2.1</b>	<b>-1.3</b>	<b>-0.7</b>	<b>-4.2</b>	<b>-4.5</b>	<b>-4.9</b>	<b>-5.4</b>	-1.5	-1.3	-1.9	-2.5	<b>-2.0</b>	<b>-2.2</b>	<b>-2.3</b>	<b>-2.5</b>
MinTs-intuitive	5.1	17.2	19.7	18.6	0.1	1.9	4.4	5.7	-3.5	-3.3	-3.4	-3.4	<b>-1.9</b>	<b>-2.0</b>	<b>-2.4</b>	<b>-2.7</b>	-1.2	0.2	0.7	0.9
MinTs-lasso	5.1	17.2	19.7	18.6	0.1	1.9	4.4	5.7	-3.5	-3.3	-3.4	-3.4	<b>-1.9</b>	<b>-2.0</b>	<b>-2.4</b>	<b>-2.7</b>	-1.2	0.2	0.7	0.9
EMinT	-2.3	24.3	58.8	59.7	36.9	56.0	68.4	70.4	51.4	64.6	75.8	81.4	65.9	72.3	81.9	85.9	48.3	62.3	75.4	79.0
Elasso	<b>-17.0</b>	<b>-19.4</b>	<b>-19.8</b>	<b>-18.7</b>	<b>-21.6</b>	<b>-17.3</b>	<b>-19.3</b>	<b>-19.6</b>	<b>-6.5</b>	<b>-9.4</b>	<b>-11.5</b>	<b>-12.6</b>	<b>2.2</b>	<b>0.4</b>	<b>-1.0</b>	<b>-1.8</b>	<b>-7.0</b>	<b>-7.7</b>	<b>-9.2</b>	<b>-9.9</b>

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

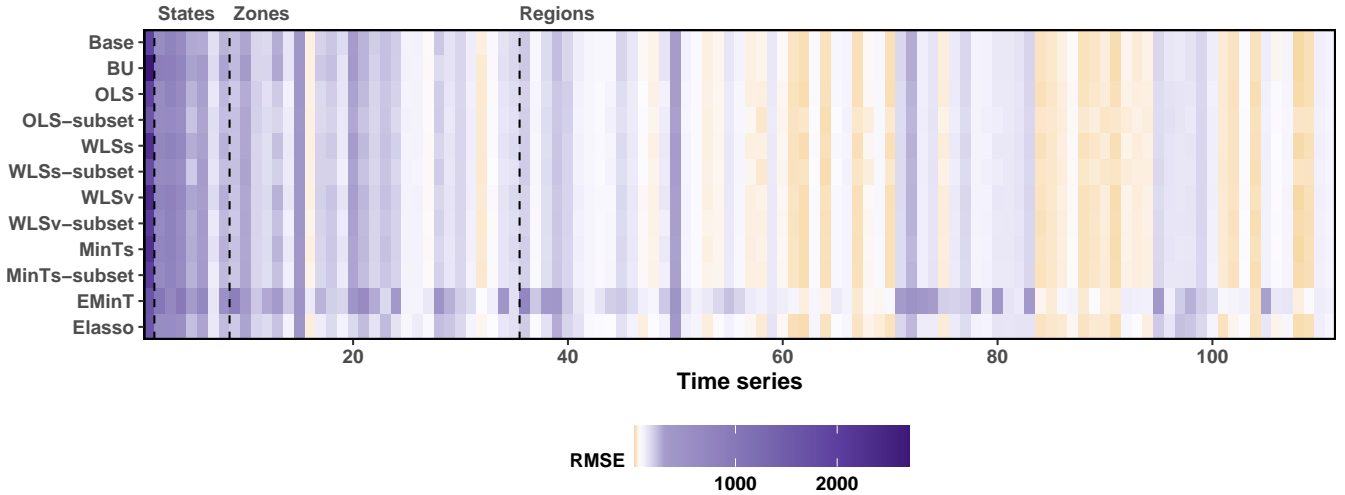
We also present the results based on the last one training set spanning from January 2017 to December 2017 in Table 10. The reconciliation errors across each of the 111 series and across the four levels in the hierarchy are displayed in Figure 5. The results show a similar performance to the average results described above, indicating relatively high-quality forecasts from the Subset and Elasso methods.



**Table 10:** Out-of-sample forecast results on a single test set (from January 2017 to December 2017) for Australian domestic tourism data.

Method	Top				State				Zone				Region				Average			
	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12	h=1	1-4	1-8	1-12
Base	1158.2	716.6	1279.5	1907.6	452.7	323.3	349.9	424.8	165.5	163.6	160.7	179.7	100.8	89.4	88.2	94.1	148.3	127.9	133.1	152.1
BU	89.1	132.8	53.4	42.0	-4.6	10.3	17.0	19.7	1.1	-2.4	0.4	1.0	0.0	0.0	0.0	0.0	5.7	7.6	7.6	8.5
OLS	-4.7	-0.4	0.5	1.4	-3.0	-3.9	-1.6	-1.5	-2.1	-4.2	-5.6	-7.5	1.0	-0.4	-1.9	-3.2	-1.0	-2.1	-2.7	-3.6
OLS-subset	-4.7	8.0	<b>-1.4</b>	<b>-14.1</b>	-3.0	5.5	0.3	<b>-7.9</b>	-2.1	-1.5	-3.7	<b>-8.7</b>	1.0	1.7	-0.1	-2.3	-1.0	1.7	-1.2	<b>-6.5</b>
OLS-intuitive	-4.7	-0.4	0.5	1.4	-3.0	-3.9	-1.6	-1.5	-2.1	-4.2	-5.6	-7.5	1.0	-0.4	-1.9	-3.2	-1.0	-2.1	-2.7	-3.6
OLS-lasso	-4.7	-0.4	0.5	1.4	-3.0	-3.9	-1.6	-1.5	-2.1	-4.2	-5.6	-7.5	1.0	-0.4	-1.9	-3.2	-1.0	-2.1	-2.7	-3.6
WLSs	25.1	55.2	20.8	19.1	-15.8	-5.0	3.5	6.2	-5.9	-5.4	-4.7	-5.0	-0.2	-0.8	-1.6	-2.2	-3.0	-0.1	0.3	0.9
WLSs-subset	25.1	<b>18.7</b>	<b>0.8</b>	<b>-7.8</b>	-15.8	-2.7	<b>-2.1</b>	<b>-6.2</b>	-5.9	-4.1	<b>-4.8</b>	<b>-8.5</b>	-0.2	0.3	-1.0	<b>-2.5</b>	-3.0	<b>-0.6</b>	<b>-2.1</b>	<b>-5.5</b>
WLSs-intuitive	25.1	55.2	20.8	19.1	-15.8	-5.0	3.5	6.2	-5.9	-5.4	-4.7	-5.0	-0.2	-0.8	-1.6	-2.2	-3.0	-0.1	0.3	0.9
WLSs-lasso	25.1	55.2	20.8	19.1	-15.8	-5.0	3.5	6.2	-5.9	-5.4	-4.7	-5.0	-0.2	-0.8	-1.6	-2.2	-3.0	-0.1	0.3	0.9
WLSv	38.2	76.2	29.6	25.6	-17.4	-3.1	7.0	9.9	-5.0	-4.3	-3.1	-3.2	-4.2	-1.6	-1.8	-2.1	-3.9	1.3	2.0	2.8
WLSv-subset	38.2	<b>34.5</b>	<b>10.7</b>	<b>8.5</b>	-17.4	<b>-8.8</b>	<b>-0.8</b>	<b>1.4</b>	-5.0	<b>-5.5</b>	<b>-5.3</b>	<b>-6.7</b>	-4.1	<b>-2.0</b>	<b>-2.6</b>	<b>-3.4</b>	-3.9	<b>-2.3</b>	<b>-2.0</b>	<b>-2.2</b>
WLSv-intuitive	38.2	76.2	29.6	25.6	-17.4	-3.1	7.0	9.9	-5.0	-4.3	-3.1	-3.2	-4.2	-1.6	-1.8	-2.1	-3.9	1.3	2.0	2.8
WLSv-lasso	38.2	76.2	29.6	25.6	-17.4	-3.1	7.0	9.9	-5.0	-4.3	-3.1	-3.2	-4.2	-1.6	-1.8	-2.1	-3.9	1.3	2.0	2.8
MinTs	20.6	53.6	21.6	19.0	<b>-22.2</b>	-7.2	3.5	6.3	<b>-12.1</b>	-6.6	-5.1	-5.3	<b>-5.3</b>	-2.6	-2.8	-3.1	-8.6	-1.8	-0.3	0.4
MinTs-subset	20.6	<b>20.0</b>	<b>6.4</b>	<b>5.6</b>	<b>-22.2</b>	<b>-11.3</b>	<b>-2.5</b>	<b>-0.1</b>	<b>-12.1</b>	-7.5	<b>-6.4</b>	<b>-7.8</b>	<b>-5.3</b>	<b>-2.9</b>	<b>-3.2</b>	<b>-3.9</b>	-8.6	<b>-4.5</b>	<b>-3.2</b>	<b>-3.3</b>
MinTs-intuitive	20.6	53.6	21.6	19.0	<b>-22.2</b>	-7.2	3.5	6.3	<b>-12.1</b>	-6.6	-5.1	-5.3	<b>-5.3</b>	-2.6	-2.8	-3.1	-8.6	-1.8	-0.3	0.4
MinTs-lasso	20.6	53.6	21.6	19.0	<b>-22.2</b>	-7.2	3.5	6.3	<b>-12.1</b>	-6.6	-5.1	-5.3	<b>-5.3</b>	-2.6	-2.8	-3.1	-8.6	-1.8	-0.3	0.4
EMinT	116.5	97.8	-15.8	-13.7	149.4	114.5	63.5	47.5	108.4	68.4	60.6	54.2	122.1	103.1	90.2	78.2	123.2	93.9	67.9	55.5
Elasso	<b>-84.5</b>	<b>-50.4</b>	<b>-16.3</b>	<b>-16.4</b>	<b>-18.3</b>	<b>0.6</b>	<b>-9.0</b>	<b>-11.4</b>	<b>-7.8</b>	<b>-8.8</b>	<b>-7.5</b>	<b>-10.4</b>	<b>2.9</b>	<b>1.6</b>	<b>4.1</b>	<b>0.3</b>	<b>-10.2</b>	<b>-4.4</b>	<b>-3.2</b>	<b>-6.7</b>

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.



**Figure 5:** Average out-of-sample forecasting performance, measured in terms of RMSE (from 1- to 12-step-ahead), for each series across different reconciliation methods. Time series are arranged along the horizontal axis.

Additionally, Table 11 presents a summary of the number of series selected using different proposed methods for each level as well as the optimal tuning parameter values identified. Here we only give the results of the Subset and Elasso methods since they are useful in the tourism application. Note that the variation in the scale of the optimal parameters for different methods comes from the difference in the scales of the objective. We observe that the OLS-subset and WLSs-subset methods exclude some series at the State and Zone levels for forecast reconciliation. In contrast, the WLSv and MinTs methods retain all

series, which is reasonable because they take into account the in-sample covariance, making themselves allow for larger adjustments made to series with large in-sample forecast error variances in forecast reconciliation. Nonetheless, the WLSv and MinTs methods can still enhance the quality of reconciled forecasts due to the inclusion of shrinkage through additional ridge regularization. It is surprising that Elasso performs exceptionally well despite using only 13 series for reconciliation.

**Table 11:** *Number of time series selected using different proposed methods and the optimal parameter values identified in the tourism application, considering a single test set (from January 2017 to December 2017). The None row shows the original number of series in the structure.*

	Number of time series retained					Optimal parameters		
	Top	State	Zone	Region	Total	$\lambda$	$\lambda_0$	$\lambda_2$
None	1	7	27	76	111	-	-	-
OLS-subset	1	2	13	76	92	-	27.98	10.00
WLSs-subset	1	1	15	76	93	-	18.73	10.00
WLSv-subset	1	7	27	76	111	-	0.03	0.01
MinTs-subset	1	7	27	76	111	-	0.05	0.01
Elasso	1	4	0	8	13	71759.21	-	-

## 6 Conclusion

In the existing literature on forecast reconciliation, we map all base forecasts into bottom-level disaggregated forecasts, which are then summed to yield coherent forecasts for the entire structure. The mapping step can be conceptually regarded as a forecast combination. It is common that the base forecasts for some time series perform poorly. This may reduce the overall effectiveness of forecast reconciliation methods. In this paper, we have addressed this issue by introducing a selection mechanism to forecast reconciliation; i.e., incorporating time series selection when reconciling forecasts, while ensuring the generation of coherent forecasts for all series.

Under the unbiasedness constraint, we developed three reconciliation methods with selection mechanisms to automatically remove some base forecasts when forming reconciled forecasts. These methods include group best-subset selection with ridge regularization (Subset), an intuitive method with  $L_0$  regularization (Intuitive), and a group lasso method (Lasso). These methods use different penalty functions designed to penalize the columns of the weighting matrix,  $\mathbf{G}$ , towards zero. Additionally, we relaxed the unbiasedness

constraint and proposed the empirical group lasso method (Elasso) which selects series based on in-sample observations and fitted values.

Simulation experiments and two empirical applications have demonstrated the superiority of the proposed methods over existing reconciliation methods that do not involve series selection. When model misspecification was introduced for some series in the hierarchy, our proposed methods guaranteed coherent forecasts that outperformed or, at least, matched their respective benchmark methods in the minimum trace reconciliation framework. In both empirical applications, where no apparent model misspecification was present, the Subset and Elasso methods were always preferred, particularly for aggregation levels and longer forecast horizons, while the Intuitive and Lasso methods yield results identical to the corresponding benchmark methods, as they tend to provide dense estimates.

A feature of the proposed methods is their ability to reduce the disparities arising from using different estimates of the base forecast error covariance matrix, thereby mitigating the challenges associated with estimator selection, which is a prominent issue within the field of forecast reconciliation research.

As the number of series grows, solving these problems efficiently becomes challenging, and the exact computation of these estimators remains a hurdle. In our study, we have used Gurobi, one of the most widely used commercial solvers, to address NP-hard MIP problems. Despite various efforts to develop MIP-based approaches for solving  $L_0$ -regularized regression problems, extending these methods to incorporate additional constraints remains a challenge. We leave this aspect to be addressed in future research.

## References

- Athanasopoulos, G., Hyndman, R. J., Kourentzes, N. & Panagiotelis, A. (2024), ‘Forecast reconciliation: A review’, *International J Forecasting* **forthcoming**.
- Athanasopoulos, G., Hyndman, R. J., Kourentzes, N. & Petropoulos, F. (2017), ‘Forecasting with temporal hierarchies’, *European J Operational Research* **262**(1), 60–74.

- Ben Taieb, S. & Koo, B. (2019), Regularized regression for hierarchical forecasting without unbiasedness conditions, in ‘Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining’, KDD ’19, Association for Computing Machinery, New York, NY, USA, pp. 1337–1347.
- Bertsimas, D., King, A. & Mazumder, R. (2016), ‘Best subset selection via a modern optimization lens’, *The Annals of Statistics* **44**(2), 813–852.
- Di Fonzo, T. & Girolimetto, D. (2023), ‘Cross-temporal forecast reconciliation: Optimal combination method and heuristic alternatives’, *International J Forecasting* **39**(1), 39–57.
- Hazimeh, H. & Mazumder, R. (2020), ‘Fast best subset selection: Coordinate descent and local combinatorial optimization algorithms’, *Operations Research* **68**(5), 1517–1537.
- Hazimeh, H., Mazumder, R. & Radchenko, P. (2023), ‘Grouped variable selection with discrete optimization: Computational and statistical perspectives’, *The Annals of Statistics* **51**(1), 1–32.
- Hazimeh, H., Mazumder, R. & Saab, A. (2022), ‘Sparse regression at scale: branch-and-bound rooted in first-order optimization’, *Mathematical Programming* **196**(1), 347–388.
- Hyndman, R. J., Ahmed, R. A., Athanasopoulos, G. & Shang, H. L. (2011), ‘Optimal combination forecasts for hierarchical time series’, *Computational Statistics & Data Analysis* **55**(9), 2579–2589.
- Hyndman, R. J. & Athanasopoulos, G. (2021), *Forecasting: principles and practice*, 3rd edn, OTexts, Melbourne, Australia. <https://OTexts.com/fpp3>.
- Hyndman, R. J., Athanasopoulos, G., Bergmeir, C., Caceres, G., Chhay, L., O’Hara-Wild, M., Petropoulos, F., Razbash, S., Wang, E. & Yasmeen, F. (2023), *forecast: Forecasting functions for time series and linear models*. R package version 8.21.1. <https://pkg.robjhyndman.com/forecast/>.
- Hyndman, R. J. & Khandakar, Y. (2008), ‘Automatic time series forecasting: the forecast package for R’, *J Statistical Software* **26**(3), 1–22.

- Hyndman, R. J., Lee, A. J. & Wang, E. (2016), ‘Fast computation of reconciled forecasts for hierarchical and grouped time series’, *Computational Statistics & Data Analysis* **97**, 16–32.
- Mazumder, R., Radchenko, P. & Dedieu, A. (2022), ‘Subset selection with shrinkage: Sparse linear modeling when the SNR is low’, *Operations Research* .
- Panagiotelis, A., Athanasopoulos, G., Gamakumara, P. & Hyndman, R. J. (2021), ‘Forecast reconciliation: A geometric view with new insights on bias correction’, *International J Forecasting* **37**(1), 343–359.
- Wickramasuriya, S. L. (2021), ‘Properties of point forecast reconciliation approaches’. arXiv:2103.11129.
- Wickramasuriya, S. L., Athanasopoulos, G. & Hyndman, R. J. (2019), ‘Optimal forecast reconciliation for hierarchical and grouped time series through trace minimization’, *J American Statistical Association* **114**(526), 804–819.
- Yuan, M. & Lin, Y. (2006), ‘Model selection and estimation in regression with grouped variables’, *J Royal Statistical Society. Series B, Statistical Methodology* **68**(1), 49–67.

# Appendix

The section provides additional results for the simulation data in Section 4.

**Table A1:** Out-of-sample forecast results for the simulated data in Scenario B, Setup 1.

Method	Top				Middle				Bottom				Average			
	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16
Base	9.6	10.7	12.6	15.6	12.1	14.4	15.3	17.0	4.2	4.9	5.9	7.5	7.2	8.5	9.6	11.4
BU	<b>-1.0</b>	0.4	0.6	0.7	<b>-47.7</b>	<b>-49.6</b>	<b>-43.6</b>	<b>-36.2</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>-23.0</b>	<b>-24.0</b>	<b>-19.8</b>	<b>-15.3</b>
OLS	8.5	13.9	10.4	7.6	-28.2	-29.4	-26.7	-23.1	22.9	23.9	17.0	11.3	-4.2	-3.8	-4.2	-4.1
OLS-subset	<b>-0.5</b>	<b>0.5</b>	<b>0.6</b>	<b>0.7</b>	<b>-46.3</b>	<b>-49.0</b>	<b>-43.2</b>	<b>-35.9</b>	<b>2.2</b>	<b>1.0</b>	<b>0.7</b>	<b>0.5</b>	<b>-21.5</b>	<b>-23.4</b>	<b>-19.4</b>	<b>-15.0</b>
OLS-intuitive	<b>-0.5</b>	<b>0.5</b>	<b>0.6</b>	<b>0.6</b>	<b>-46.5</b>	<b>-49.0</b>	<b>-43.2</b>	<b>-36.0</b>	<b>2.2</b>	<b>1.2</b>	<b>0.7</b>	<b>0.5</b>	<b>-21.6</b>	<b>-23.4</b>	<b>-19.4</b>	<b>-15.0</b>
OLS-lasso	<b>-0.2</b>	<b>1.5</b>	<b>1.4</b>	<b>1.3</b>	<b>-46.9</b>	<b>-48.9</b>	<b>-43.1</b>	<b>-35.8</b>	<b>0.9</b>	<b>0.8</b>	<b>0.5</b>	<b>0.3</b>	<b>-22.1</b>	<b>-23.3</b>	<b>-19.3</b>	<b>-14.9</b>
WLSs	12.1	18.6	14.0	10.2	-34.4	-35.1	-31.7	-26.9	15.6	17.0	12.0	8.0	-9.0	-8.0	-7.6	-6.5
WLSs-subset	<b>-0.1</b>	<b>1.2</b>	<b>1.1</b>	<b>1.1</b>	<b>-46.7</b>	<b>-48.8</b>	<b>-43.1</b>	<b>-35.8</b>	<b>1.5</b>	<b>1.1</b>	<b>0.8</b>	<b>0.6</b>	<b>-21.8</b>	<b>-23.2</b>	<b>-19.2</b>	<b>-14.8</b>
WLSs-intuitive	<b>0.0</b>	<b>1.2</b>	<b>1.0</b>	<b>0.9</b>	<b>-46.5</b>	<b>-48.8</b>	<b>-43.1</b>	<b>-35.9</b>	<b>1.7</b>	<b>1.3</b>	<b>0.9</b>	<b>0.6</b>	<b>-21.6</b>	<b>-23.1</b>	<b>-19.2</b>	<b>-14.9</b>
WLSs-lasso	<b>-0.1</b>	<b>1.5</b>	<b>1.5</b>	<b>1.3</b>	<b>-46.7</b>	<b>-48.9</b>	<b>-43.1</b>	<b>-35.8</b>	<b>0.9</b>	<b>0.8</b>	<b>0.5</b>	<b>0.3</b>	<b>-22.0</b>	<b>-23.2</b>	<b>-19.3</b>	<b>-14.9</b>
WLSv	-0.8	2.3	1.8	1.6	-46.3	-47.9	-42.3	-35.2	1.6	1.9	1.2	0.8	-21.7	-22.2	-18.6	-14.4
WLSv-subset	-0.7	<b>1.3</b>	<b>1.4</b>	<b>1.4</b>	<b>-46.9</b>	<b>-48.7</b>	<b>-42.9</b>	<b>-35.6</b>	<b>1.0</b>	<b>1.0</b>	<b>0.8</b>	<b>0.6</b>	<b>-22.2</b>	<b>-23.1</b>	<b>-19.1</b>	<b>-14.7</b>
WLSv-intuitive	-0.4	<b>1.5</b>	<b>1.4</b>	<b>1.2</b>	<b>-46.9</b>	<b>-48.6</b>	<b>-42.8</b>	<b>-35.6</b>	<b>0.9</b>	<b>1.2</b>	<b>0.9</b>	<b>0.7</b>	<b>-22.2</b>	<b>-23.0</b>	<b>-19.0</b>	<b>-14.7</b>
WLSv-lasso	-0.6	<b>1.3</b>	<b>1.3</b>	<b>1.3</b>	<b>-47.2</b>	<b>-48.9</b>	<b>-43.0</b>	<b>-35.7</b>	<b>0.6</b>	<b>0.8</b>	<b>0.5</b>	<b>0.4</b>	<b>-22.4</b>	<b>-23.3</b>	<b>-19.2</b>	<b>-14.8</b>
MinT	0.2	0.5	0.6	0.5	-47.5	-49.4	-43.5	-36.1	1.1	0.5	0.3	0.1	-22.3	-23.7	-19.6	<b>-15.3</b>
MinT-subset	<b>-0.1</b>	0.8	0.9	0.9	-46.9	-49.1	-43.3	-36.0	1.7	0.9	0.5	0.3	-21.9	-23.4	-19.4	-15.1
MinT-intuitive	0.2	0.5	0.6	0.5	-47.5	-49.4	-43.5	-36.1	1.1	0.5	0.3	0.1	-22.3	-23.7	-19.6	<b>-15.3</b>
MinT-lasso	<b>-0.3</b>	<b>0.3</b>	0.6	0.5	<b>-47.6</b>	-49.4	-43.5	-36.1	<b>0.8</b>	<b>0.3</b>	<b>0.2</b>	0.1	<b>-22.5</b>	<b>-23.9</b>	<b>-19.7</b>	<b>-15.3</b>
MinTs	-0.3	0.3	<b>0.4</b>	<b>0.4</b>	-47.6	-49.5	<b>-43.6</b>	<b>-36.2</b>	0.7	0.2	0.1	<b>0.0</b>	-22.6	-23.9	<b>-19.8</b>	<b>-15.3</b>
MinTs-subset	<b>-0.8</b>	0.5	0.8	0.8	-47.2	-49.2	-43.4	-36.0	1.0	0.7	0.4	0.3	-22.3	-23.6	-19.5	-15.1
MinTs-intuitive	-0.3	0.3	<b>0.4</b>	<b>0.4</b>	-47.6	-49.5	<b>-43.6</b>	<b>-36.2</b>	0.7	0.2	0.1	<b>0.0</b>	-22.6	-23.9	<b>-19.8</b>	<b>-15.3</b>
MinTs-lasso	<b>-0.9</b>	<b>0.2</b>	0.5	0.5	<b>-47.7</b>	-49.5	<b>-43.6</b>	<b>-36.2</b>	<b>0.5</b>	0.2	0.1	0.1	<b>-22.8</b>	<b>-24.0</b>	<b>-19.8</b>	<b>-15.3</b>
EMinT	2.2	2.9	2.5	1.7	-46.2	-48.1	-42.4	-35.3	3.6	2.9	2.0	1.1	-20.5	-21.9	-18.2	-14.3
Elasso	<b>1.4</b>	<b>2.7</b>	<b>2.4</b>	<b>1.6</b>	<b>-46.4</b>	<b>-48.2</b>	-42.4	<b>-35.4</b>	<b>3.1</b>	3.2	2.1	1.2	<b>-20.9</b>	-21.9	-18.2	-14.3

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

**Table A2:** Out-of-sample forecast results for the simulated data in Scenario C, Setup 1.

Method	Top				Middle				Bottom				Average			
	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16	h=1	1-4	1-8	1-16
Base	25.0	30.3	30.9	32.3	6.3	7.3	8.6	10.8	4.2	4.9	5.9	7.5	7.8	9.2	10.3	12.0
BU	-62.0	<b>-64.4</b>	<b>-59.0</b>	-51.5	<b>-0.3</b>	<b>0.0</b>	<b>0.1</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>-28.5</b>	<b>-30.2</b>	<b>-25.3</b>	<b>-19.8</b>
OLS	-34.8	-35.5	-33.5	-30.1	45.3	50.6	37.7	25.1	27.7	29.9	21.2	13.7	3.1	3.8	1.6	-0.2
OLS-subset	<b>-35.3</b>	<b>-41.9</b>	<b>-39.2</b>	<b>-35.0</b>	<b>43.9</b>	<b>39.5</b>	<b>29.5</b>	<b>19.6</b>	<b>27.1</b>	<b>23.6</b>	<b>16.8</b>	<b>10.9</b>	<b>2.4</b>	<b>-3.5</b>	<b>-4.2</b>	<b>-4.5</b>
OLS-intuitive	<b>-41.2</b>	<b>-49.2</b>	<b>-45.5</b>	<b>-40.0</b>	<b>35.1</b>	<b>26.8</b>	<b>20.3</b>	<b>13.7</b>	<b>21.9</b>	<b>15.9</b>	<b>11.5</b>	<b>7.6</b>	<b>-4.0</b>	<b>-12.2</b>	<b>-10.9</b>	<b>-9.1</b>
OLS-lasso	<b>-61.8</b>	<b>-63.6</b>	<b>-58.1</b>	<b>-50.9</b>	<b>0.4</b>	<b>1.3</b>	<b>1.3</b>	<b>0.7</b>	<b>0.3</b>	<b>0.8</b>	<b>0.6</b>	<b>0.4</b>	<b>-28.2</b>	<b>-29.3</b>	<b>-24.5</b>	<b>-19.2</b>
WLSs	-50.9	-52.4	-48.7	-43.3	17.6	20.0	14.5	9.3	9.6	11.3	7.7	4.9	-16.3	-16.7	-14.9	-12.5
WLSs-subset	<b>-61.8</b>	<b>-63.6</b>	<b>-58.1</b>	<b>-50.7</b>	<b>0.3</b>	<b>1.4</b>	<b>1.4</b>	<b>0.9</b>	<b>0.3</b>	<b>0.9</b>	<b>0.7</b>	<b>0.6</b>	<b>-28.2</b>	<b>-29.3</b>	<b>-24.4</b>	<b>-19.0</b>
WLSs-intuitive	<b>-61.8</b>	<b>-63.8</b>	<b>-58.3</b>	<b>-50.9</b>	<b>0.0</b>	<b>1.0</b>	<b>1.0</b>	<b>0.7</b>	<b>0.3</b>	<b>0.7</b>	<b>0.6</b>	<b>0.5</b>	<b>-28.3</b>	<b>-29.5</b>	<b>-24.6</b>	<b>-19.2</b>
WLSs-lasso	<b>-61.7</b>	<b>-63.5</b>	<b>-58.0</b>	<b>-50.7</b>	<b>0.5</b>	<b>1.5</b>	<b>1.4</b>	<b>0.9</b>	<b>0.3</b>	<b>0.9</b>	<b>0.7</b>	<b>0.5</b>	<b>-28.1</b>	<b>-29.2</b>	<b>-24.4</b>	<b>-19.1</b>
WLSv	-61.1	-63.4	-58.1	-50.8	1.0	1.7	1.3	0.8	0.7	1.0	0.6	0.4	-27.6	-29.1	-24.5	-19.2
WLSv-subset	<b>-61.9</b>	<b>-63.6</b>	<b>-58.2</b>	<b>-50.9</b>	<b>0.2</b>	<b>1.3</b>	<b>1.2</b>	<b>0.8</b>	<b>0.1</b>	<b>0.8</b>	<b>0.6</b>	<b>0.5</b>	<b>-28.3</b>	<b>-29.3</b>	<b>-24.5</b>	<b>-19.2</b>
WLSv-intuitive	<b>-61.8</b>	<b>-63.8</b>	<b>-58.3</b>	<b>-51.0</b>	<b>0.0</b>	<b>1.1</b>	<b>1.1</b>	<b>0.6</b>	<b>0.1</b>	<b>0.6</b>	<b>0.5</b>	<b>0.4</b>	<b>-28.4</b>	<b>-29.5</b>	<b>-24.7</b>	<b>-19.3</b>
WLSv-lasso	<b>-61.8</b>	<b>-63.9</b>	<b>-58.4</b>	<b>-51.1</b>	<b>0.2</b>	<b>0.9</b>	<b>0.9</b>	<b>0.5</b>	<b>0.1</b>	<b>0.5</b>	<b>0.4</b>	<b>0.3</b>	<b>-28.3</b>	<b>-29.6</b>	<b>-24.8</b>	<b>-19.4</b>
MinT	-62.1	-64.3	-58.9	<b>-51.6</b>	-0.2	0.6	0.5	0.2	0.8	0.5	0.3	0.1	-28.3	-29.9	-25.1	<b>-19.8</b>
MinT-subset	-61.8	-63.7	-58.2	-50.9	0.4	1.2	1.3	0.8	0.8	1.0	0.7	0.5	-28.0	-29.3	-24.5	-19.2
MinT-intuitive	-62.1	-64.3	-58.9	<b>-51.6</b>	-0.2	0.6	0.5	0.2	0.8	0.5	0.3	0.1	-28.3	-29.9	-25.1	<b>-19.8</b>
MinT-lasso	-62.1	<b>-64.4</b>	-58.9	-51.5	<b>-0.3</b>	<b>0.3</b>	<b>0.4</b>	<b>0.1</b>	<b>0.6</b>	<b>0.3</b>	<b>0.1</b>	<b>0.1</b>	<b>-28.4</b>	<b>-30.1</b>	<b>-25.2</b>	<b>-19.8</b>
MinTs	<b>-62.2</b>	<b>-64.4</b>	<b>-59.0</b>	<b>-51.6</b>	<b>-0.3</b>	0.3	0.4	0.1	0.4	0.3	0.1	<b>0.0</b>	<b>-28.5</b>	-30.1	-25.2	<b>-19.8</b>
MinTs-subset	-62.0	-63.8	-58.4	-51.1	0.4	1.1	1.2	0.7	0.5	0.9	0.7	0.5	-28.2	-29.5	-24.6	-19.3
MinTs-intuitive	<b>-62.2</b>	<b>-64.4</b>	<b>-59.0</b>	<b>-51.6</b>	<b>-0.3</b>	0.3	0.4	0.1	0.4	0.3	0.1	<b>0.0</b>	<b>-28.5</b>	-30.1	-25.2	<b>-19.8</b>
MinTs-lasso	<b>-62.2</b>	<b>-64.4</b>	-58.9	-51.5	-0.2	0.3	0.4	0.1	<b>0.2</b>	<b>0.2</b>	0.1	<b>0.0</b>	<b>-28.5</b>	-30.1	-25.2	<b>-19.8</b>
EMinT	-60.7	-63.5	-58.2	-51.0	2.5	2.9	2.3	1.3	3.6	2.9	2.0	1.1	-26.2	-28.3	-23.8	-18.9
Elasso	<b>-60.9</b>	<b>-63.6</b>	-58.2	<b>-51.1</b>	<b>2.3</b>	<b>2.8</b>	2.3	1.3	<b>3.1</b>	3.1	2.1	1.2	<b>-26.5</b>	-28.3	-23.8	-18.9

Note: The Base row shows the average RMSE of the base forecasts. Entries below this row indicate the percentage decrease (negative) or increase (positive) in the average RMSE of the reconciled forecasts compared to the base forecasts. The entries with the lowest values in each column are highlighted in blue. In each panel, the proposed methods are indicated with a gray background, and methods that outperform the benchmark method are marked in bold.

**Table A3:** Proportion of time series being selected after using the proposed reconciliation methods with selection in Scenario B, Setup 1.

	Top	A	B	AA	AB	BA	BB	Summary
OLS-subset	0.55	0.04	0.41	0.74	0.78	0.79	0.83	
OLS-intuitive	0.61	0.04	0.52	0.75	0.69	0.69	0.83	
OLS-lasso	0.04	0.35	0.02	1.00	1.00	1.00	1.00	
WLSs-subset	0.45	0.06	0.36	0.81	0.84	0.81	0.87	
WLSs-intuitive	0.61	0.06	0.48	0.75	0.71	0.73	0.84	
WLSs-lasso	0.02	0.33	0.02	1.00	1.00	1.00	1.00	
WLSv-subset	0.54	0.29	0.46	0.91	0.94	0.86	0.89	
WLSv-intuitive	0.59	0.32	0.53	0.82	0.86	0.77	0.86	
WLSv-lasso	0.27	0.42	0.26	1.00	1.00	1.00	1.00	
MinT-subset	0.69	0.64	0.66	0.95	0.96	0.90	0.90	
MinT-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinT-lasso	0.82	0.74	0.83	1.00	0.99	0.97	0.97	
MinTs-subset	0.62	0.63	0.58	0.95	0.96	0.90	0.86	
MinTs-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinTs-lasso	0.68	0.75	0.68	1.00	1.00	1.00	1.00	
Elasso	0.78	0.95	0.68	1.00	1.00	1.00	1.00	

Note: the last column displays a stacked barplot for each method, based on the total number of selected series data from 500 simulation instances, with a darker sub-bar indicating a larger number.

**Table A4:** Proportion of time series being selected after using the proposed reconciliation methods with selection in Scenario C, Setup 1.

	Top	A	B	AA	AB	BA	BB	Summary
OLS-subset	0.75	0.45	0.44	0.82	0.79	0.83	0.80	
OLS-intuitive	0.47	0.70	0.69	0.86	0.92	0.90	0.89	
OLS-lasso	0.38	0.01	0.01	1.00	1.00	1.00	1.00	
WLSs-subset	0.08	0.42	0.41	0.87	0.85	0.84	0.89	
WLSs-intuitive	0.06	0.55	0.50	0.66	0.87	0.69	0.88	
WLSs-lasso	0.35	0.03	0.03	1.00	1.00	1.00	1.00	
WLSv-subset	0.31	0.67	0.65	0.88	0.90	0.91	0.90	
WLSv-intuitive	0.34	0.63	0.60	0.80	0.89	0.84	0.87	
WLSv-lasso	0.45	0.35	0.36	1.00	1.00	1.00	1.00	
MinT-subset	0.69	0.78	0.80	0.91	0.91	0.91	0.91	
MinT-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinT-lasso	0.75	0.89	0.86	0.97	0.97	0.97	0.97	
MinTs-subset	0.67	0.74	0.76	0.90	0.89	0.88	0.91	
MinTs-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinTs-lasso	0.77	0.72	0.73	1.00	1.00	1.00	1.00	
Elasso	0.95	0.64	0.64	1.00	1.00	1.00	1.00	

Note: the last column displays a stacked barplot for each method, based on the total number of selected series data from 500 simulation instances, with a darker sub-bar indicating a larger number.

**Table A5:** Proportion of time series being selected after using the proposed reconciliation methods with selection in Setup 2, with the error correlation being 0.8.

	Top	A	B	AA	AB	BA	BB	Summary
OLS-subset	0.33	0.52	0.96	0.95	0.98	0.96	0.78	
OLS-intuitive	0.54	0.77	0.93	0.89	0.97	0.83	0.85	
OLS-lasso	0.69	0.53	0.60	1.00	1.00	1.00	1.00	
WLSs-subset	0.29	0.60	1.00	1.00	1.00	0.98	0.86	
WLSs-intuitive	0.63	0.67	0.99	0.98	1.00	0.93	0.86	
WLSs-lasso	0.69	0.76	0.91	1.00	1.00	1.00	1.00	
WLSv-subset	0.32	0.55	1.00	1.00	1.00	0.99	0.76	
WLSv-intuitive	0.58	0.56	1.00	1.00	0.98	1.00	0.75	
WLSv-lasso	0.77	0.84	0.99	1.00	1.00	1.00	1.00	
MinT-subset	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinT-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinT-lasso	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinTs-subset	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinTs-intuitive	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
MinTs-lasso	1.00	1.00	1.00	1.00	1.00	1.00	1.00	
Elasso	0.73	0.65	0.98	0.98	0.86	1.00	0.99	

Note: the last column displays a stacked barplot for each method, based on the total number of selected series data from 500 simulation instances, with a darker sub-bar indicating a larger number.