

Sparse Bayesian vector autoregressions in huge dimensions

Gregor Kastner¹ | Florian Huber²

¹Institute for Statistics and Mathematics, WU Vienna University of Economics and Business, Vienna, Austria

²Salzburg Centre of European Union Studies (SCEUS), University of Salzburg, Salzburg, Austria

Correspondence

Gregor Kastner, Institute for Statistics and Mathematics, WU Vienna University of Economics and Business, Welthandelsplatz 1, 1020 Vienna, Austria.
Email: gregor.kastner@wu.ac.at

Abstract

We develop a Bayesian vector autoregressive (VAR) model with multivariate stochastic volatility that is capable of handling vast dimensional information sets. Three features are introduced to permit reliable estimation of the model. First, we assume that the reduced-form errors in the VAR feature a factor stochastic volatility structure, allowing for conditional equation-by-equation estimation. Second, we apply recently developed global–local shrinkage priors to the VAR coefficients to cure the curse of dimensionality. Third, we utilize recent innovations to sample efficiently from high-dimensional multivariate Gaussian distributions. This makes simulation-based fully Bayesian inference feasible when the dimensionality is large but the time series length is moderate. We demonstrate the merits of our approach in an extensive simulation study and apply the model to US macroeconomic data to evaluate its forecasting capabilities.

KEYWORDS

Dirichlet-Laplace prior, efficient MCMC, factor stochastic volatility, normal-Gamma prior, shrinkage

1 | INTRODUCTION

Previous research has identified two important features that macroeconomic models should possess: the ability to exploit high-dimensional information sets (Bańbura et al., 2010; Koop et al., 2019; Rockova & McAlinn, 2017; Stock & Watson, 2011) and the possibility to capture nonlinear features of the underlying time series (Bitto & Frühwirth-Schnatter, 2019; Clark, 2011; Clark & Ravazzolo, 2015; Cogley & Sargent, 2001; Huber et al., 2019; Primiceri, 2005). While the literature suggests several paths to estimate large models, the majority of such approaches imply that once nonlinearities are taken into account analytical solutions are no longer available and the computational burden becomes prohibitive. This implies that high-dimensional nonlinear models can practically be estimated only under strong (and often unrealistic) restrictions on the dynamics of the model.

However, especially in forecasting applications or in structural analysis, successful models should generally be able to exploit much information and also control for breaks in the autoregressive parameters or, more importantly, changes in the volatility of economic shocks (Koop et al., 2009; Primiceri, 2005; Sims & Zha, 2006).

Two reasons limit the use of large (or even huge) nonlinear models. The first reason is statistical. Since the number of parameters in a standard vector autoregression rises quadratically with the number of time series included and commonly used macroeconomic time series are rather short, in-sample overfitting turns out to be a serious issue. As a solution, the Bayesian literature on vector autoregressive (VAR) modeling (e.g., Ankargren et al., 2019; Bańbura et al., 2010; Clark, 2011; Clark & Ravazzolo, 2015; Doan et al., 1984; Follett & Yu, 2019; George et al., 2008; Huber & Feldkircher, 2019; Koop, 2013; Korobilis & Pettenuzzo,

2019; Litterman, 1986; Sims & Zha, 1998) suggests shrinkage priors that push the parameter space towards some stylized prior model like a multivariate random walk. On the other hand, Ahelegbey et al. (2016) suggest viewing VARs as graphical models and perform model selection drawing from the literature on sparse directed acyclic graphs. This typically leads to much improved forecasting properties and more meaningful structural inference. Moreover, the majority of the literature on Bayesian VARs imposes conjugate priors on the autoregressive parameters, allowing for analytical posterior solutions and thus avoiding simulation-based techniques such as Markov chain Monte Carlo (MCMC). Frequentist approaches often consider multistep approaches (e.g., Davis et al., 2016).

The second reason is computational. Nonlinear Bayesian models typically have to be estimated by means of MCMC, and computational intensity increases vastly when the number of component series becomes large. This increase stems from the fact that standard algorithms for multivariate regression models call for the inversion of large covariance matrices. Especially for sizable systems, this can quickly turn prohibitive since the inverse of the posterior variance–covariance matrix on the coefficients has to be computed for each sweep of the MCMC algorithm. For natural conjugate models, this step can be vastly simplified because the likelihood possesses a convenient Kronecker structure, implying that all equations in the VAR feature the same set of explanatory variables. This speeds up computation by large margins but restricts the flexibility of the model. Carriero et al. (2016), for instance, exploit this fact and introduce a simplified stochastic volatility specification. Another strand of the literature augments each equation of the VAR by including the residuals of the preceding equations (Carriero et al., 2019), which also provides significant improvements in terms of computational speed. Finally, in a recent contribution, Koop et al. (2019) reduce the dimensionality of the problem at hand by randomly compressing the lagged endogenous variables in the VAR.

All papers mentioned hitherto focus on capturing cross-variable correlation in the conditional mean through the VAR part, and the comovement in volatilities is captured by a rich specification of the error variance (Primiceri, 2005) or by a single factor (Carriero et al., 2016). Another strand of the literature, typically used in financial econometrics, utilizes factor models to provide a parsimonious representation of a covariance matrix, focusing exclusively on the second moment of the predictive density. For instance, Pitt and Shephard (1999) and Aguilar and West (2000) assume that the variance–covariance matrix of a broad panel of time series might be described by a lower dimensional matrix of latent factors featuring

stochastic volatility and a variable-specific idiosyncratic stochastic volatility process.¹

The present paper combines the virtues of exploiting large information sets and allowing for movements in the error variance. The overfitting issue mentioned above is solved as follows. First, we use a Dirichlet–Laplace (DL) prior specification (see Bhattacharya et al., 2015) on the VAR coefficients. This prior is a global–local shrinkage prior in the spirit of Polson and Scott (2011) that enables us to heavily shrink the parameter space but at the same time provides enough flexibility to allow for nonzero regression coefficients if necessary. Second, a factor stochastic volatility model on the VAR errors grants a parsimonious representation of the time-varying error variance–covariance matrix of the VAR. To deal with the computational complexity, we exploit the fact that, conditionally on the latent factors and their loadings, equation-by-equation estimation becomes possible within each MCMC iteration. Moreover, we apply recent advances for fast sampling from high-dimensional multivariate Gaussian distributions (Bhattacharya et al., 2016) that permit estimation of models with hundreds of thousands of autoregressive parameters and an error covariance matrix with tens of thousands of nontrivial time-varying elements on a quarterly US data set in a reasonable amount of time. In a careful analysis, we show to what extent our proposed method improves upon a set of standard algorithms typically used to simulate from the joint posterior distribution of large-dimensional Bayesian VARs.

We first assess the merits of our approach in an extensive simulation study based on a range of different data-generating processes (DGPs). Relative to a set of competing benchmark specifications we show that, in terms of point estimates, the proposed global–local shrinkage prior yields precise parameter estimates and successfully introduces shrinkage in the modeling framework, without overshrinking significant signals.

In an empirical application, we adopt a modified version of the quarterly data set proposed by Stock and Watson (2011) and McCracken and Ng (2016). To illustrate the out-of-sample performance of our model, we forecast important economic indicators such as output, consumer price inflation, and short-term interest rates, amongst others. The proposed model is benchmarked against several alternatives. Our findings suggest that it performs well in terms of one-step-ahead predictive likelihoods. In addition, investigating the time profile of the cumulative log-predictive likelihood reveals that allowing for large information sets in combination with the factor structure especially pays off in times of economic stress.

¹Two recent exceptions are Koop and Korobilis (2013) and Carriero et al. (2016).

The remainder of this paper is structured as follows. Section 2 introduces the econometric framework. Section 3 details the Bayesian estimation approach, including an elaborated account of the (shrinkage) prior setup adopted and the corresponding conditional posterior distributions. Section 4 provides an analysis of the computational gains of our algorithm relative to a set of established algorithms. Section 5 presents the results of an extensive simulation study comparing the performance of carefully selected shrinkage priors for different time series lengths and model dimensions within various (sparse and dense) data-generating scenarios. Section 6, after giving a brief overview of the data set used along with the model specification, illustrates our modeling approach by fitting a single-factor model to 215-dimensional quarterly US data. Moreover, we perform a forecasting exercise to assess the predictive performance of our approach and discuss the choice of the number of latent factors. Finally, Section 7 concludes.

2 | ECONOMETRIC FRAMEWORK

Suppose interest centers on modeling an $m \times 1$ vector of time series denoted by \mathbf{y}_t with $t = 1, \dots, T$. We assume that \mathbf{y}_t follows a heteroskedastic VAR(p) process:²

$$\mathbf{y}_t = \mathbf{A}_1 \mathbf{y}_{t-1} + \dots + \mathbf{A}_p \mathbf{y}_{t-p} + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}_m(\mathbf{0}, \boldsymbol{\Omega}_t). \quad (1)$$

Each \mathbf{A}_j ($j = 1, \dots, p$) is an $m \times m$ matrix of autoregressive coefficients. The error term is assumed to follow a multivariate Gaussian distribution with time-varying variance–covariance matrix $\boldsymbol{\Omega}_t$. To permit reliable and parsimonious estimation when m is large, we decompose the residual covariance matrix into

$$\boldsymbol{\Omega}_t = \boldsymbol{\Lambda} \mathbf{V}_t \boldsymbol{\Lambda} + \boldsymbol{\Sigma}_t, \quad (2)$$

where both $\boldsymbol{\Sigma}_t = \text{diag}(\sigma_{1t}^2, \dots, \sigma_{mt}^2)$ and $\mathbf{V}_t = \text{diag}(e^{h_{1t}}, \dots, e^{h_{qt}})$ are diagonal matrices with dimension m and q , respectively, and $\boldsymbol{\Lambda}$ denotes an $m \times q$ matrix of factor loadings with typical element λ_{ij} ($i = 1, \dots, m$; $j = 1, \dots, q$). The logarithms of the diagonal elements of $\boldsymbol{\Sigma}_t$ and \mathbf{V}_t follow AR(1) processes:

$$h_{jt} = \rho_{hj} h_{j,t-1} + e_{hj,t}, \quad j = 1, \dots, q, \quad (3)$$

$$\log \sigma_{it}^2 = \mu_{\sigma i} + \rho_{\sigma i} (\log \sigma_{i,t-1}^2 - \mu_{\sigma i}) + e_{\sigma i,t}, \quad i = 1, \dots, m. \quad (4)$$

To identify the scaling of the elements of $\boldsymbol{\Lambda}$, the process specified in Equation (3) is assumed to have mean zero, while $\mu_{\sigma j}$ in Equation (4) is the unconditional mean of the log-elements of $\boldsymbol{\Sigma}_t$ to be estimated from the data (cf. Kastner et al., 2017). The parameters ρ_{hj} and $\rho_{\sigma i}$ are a priori

restricted to the interval $(-1, 1)$ and denote the persistence of the latent log variances. The error terms $e_{hj,t}$ and $e_{\sigma i,t}$ constitute independent zero mean innovations with variances ζ_{hj}^2 and $\zeta_{\sigma i}^2$, respectively. This specification implies that the volatilities are mean reverting and thus bounded in the limit.

This error structure is known as the factor stochastic volatility model (see, e.g., Aguilar & West, 2000; Pitt & Shephard, 1999). It can be equivalently written by introducing q conditionally independent latent factors $\mathbf{f}_t \sim \mathcal{N}_q(\mathbf{0}, \mathbf{V}_t)$ and rewriting the error term in Equation (1) as

$$\boldsymbol{\varepsilon}_t = \boldsymbol{\Lambda} \mathbf{f}_t + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \mathcal{N}_m(\mathbf{0}, \boldsymbol{\Sigma}_t). \quad (5)$$

Note that off-diagonal entries of $\boldsymbol{\Omega}_t$ exclusively stem from the volatilities of the q factors, while the diagonal entries of $\boldsymbol{\Omega}_t$ are allowed to feature idiosyncratic deviations driven by the elements of $\boldsymbol{\Sigma}_t$. This specification reduces the number of free elements in $\boldsymbol{\Omega}_t$ from $m(m+1)/2$ to mq , where the latter quantity is typically much smaller than the former. In addition, by conditioning on the latent factors, this representation enables us to derive an efficient Gibbs sampler that allows for conditional equation-by-equation estimation. As will be discussed in more detail in Section 3.2, this constitutes a key feature for computationally feasible Bayesian inference when the dimensionality m becomes large.

The model described by Equations (1) and (2) is related to several alternative specifications commonly used in the literature. For instance, assuming that $\mathbf{V}_t = \mathbf{I}$ and $\boldsymbol{\Sigma}_t \equiv \boldsymbol{\Sigma}$ for all t leads to the specification adopted in Stock and Watson (2005). Setting $q = 1$ and $\boldsymbol{\Sigma}_t \equiv \boldsymbol{\Sigma}$ yields a specification that is similar to the one stipulated in Carriero et al. (2016), with the difference that our model imposes restrictions on the covariances whereas Carriero et al. (2016) estimate a full (but constant) covariance matrix. In addition, our model implies that the stochastic volatility enters $\boldsymbol{\Omega}_t$ in an additive fashion.

Before proceeding to the next subsection it is worth summarizing the key features of the model given by Equations (1)–(5). First, we capture cross-variable movements in the conditional mean through the VAR block of the model and assume that comovement in conditional variances is captured by a factor structure. Second, the model introduces stochastic volatility by assuming that a large panel of volatilities may be efficiently summarized through a set of latent heteroskedastic factors. This choice is more flexible than a single-factor model for the volatility, effectively providing a parsimonious representation of $\boldsymbol{\Omega}_t$ that is flexible enough to replicate the dynamic behavior of the variances of a broad set of macroeconomic quantities.

²For simplicity of exposition we omit the intercept term in the following discussion (which we nonetheless include in the empirical application).

3 | INFERENCE IN LARGE-DIMENSIONAL VAR MODELS

Our approach to estimation and inference is Bayesian. This implies that, after specifying a suitable prior distribution on the model parameters, we can combine this prior with the likelihood implied by the data and the model to obtain the corresponding posterior distribution.

3.1 | A global–local shrinkage prior

For prior implementation, it proves to be convenient to define a $k \times 1$ vector of predictors $\mathbf{x}_t = (\mathbf{y}'_{t-1}, \dots, \mathbf{y}'_{t-p})'$ and an $m \times k$ coefficient matrix $\mathbf{B} = (\mathbf{A}_1, \dots, \mathbf{A}_p)$ with $k = mp$ to rewrite the model in Equation (1) more compactly as $\mathbf{y}_t = \mathbf{B}\mathbf{x}_t + \varepsilon_t$. Stacking the rows of \mathbf{y}_t , \mathbf{x}_t , and ε_t yields

$$\mathbf{Y} = \mathbf{X}\mathbf{B}' + \mathbf{E}, \quad (6)$$

where $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_T)'$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)'$, and $\mathbf{E} = (\varepsilon_1, \dots, \varepsilon_T)'$ denote the corresponding full data matrices.

Typically, the matrix \mathbf{B} is a sparse matrix with nonzero elements mainly located on the main diagonal of \mathbf{A}_1 . In fact, existing priors in the Minnesota tradition tend to strongly push the system towards the prior model in high dimensions. However, especially in large models, an extremely tight prior on \mathbf{B} might lead to severe over-shrinking, effectively zeroing out coefficients that might be important to explain \mathbf{y}_t . If the matrix \mathbf{B} is characterized by a relatively low number of nonzero regression coefficients, a possible solution is a global–local shrinkage prior (Polson & Scott, 2011).

A recent variant that falls within the class of global–local shrinkage priors is the Dirichlet–Laplace (DL) prior put forward in Bhattacharya et al. (2015). This prior possesses convenient shrinkage properties in the presence of a large degree of sparsity of the parameter vector $\mathbf{b} = \text{vec}(\mathbf{B})$. In what follows, we impose the DL prior on each of the $K = mk$ elements of \mathbf{b} , denoted by b_j , for $j = 1, \dots, K$:

$$b_j \sim \mathcal{DE}(\vartheta_j, \zeta) \Leftrightarrow b_j \sim \mathcal{N}(0, \psi_j \vartheta_j^2 \zeta^2), \quad \psi_j \sim \mathcal{E}(1/2), \quad (7)$$

where \mathcal{DE} denotes the double exponential (Laplace) and \mathcal{E} the exponential distribution, ψ_j is an auxiliary scaling parameter to achieve conditional normality, and the elements of $\boldsymbol{\vartheta} = (\vartheta_1, \dots, \vartheta_K)'$ are local auxiliary scaling parameters that are bounded to the $(K-1)$ -dimensional simplex $\mathcal{S}^{K-1} = \{\boldsymbol{\vartheta} : \vartheta_j \geq 0, \sum_{j=1}^n \vartheta_j = 1\}$. A natural prior choice for ϑ_j is the (symmetric) Dirichlet distribution with hyperparameter a : $\vartheta_j \sim \mathcal{D}(a, \dots, a)$. In addition, ζ is a global shrinkage parameter that pushes all elements in \mathbf{B} towards zero and exhibits an important role in determining the tail behavior of the marginal prior distribution on b_j , obtained after integrating out the ϑ_j .

Thus we follow Bhattacharya et al. (2015) and adopt a fully Bayesian approach by specifying a gamma distributed prior on $\zeta \sim \mathcal{G}(Ka, 1/2)$. It is noteworthy that this prior setup has at least two convenient features that appear to be of prime importance for VAR modeling. First, it exerts a strong degree of shrinkage on all elements of \mathbf{B} but still provides additional flexibility such that nonzero regression coefficients are permitted. This critical property is a feature which a large class of global–local shrinkage priors share (Griffin & Brown, 2010; Carvalho et al., 2010; Polson & Scott, 2011) and has been recently adopted in a VAR framework by Huber and Feldkircher (2019) and within the general context of state-space models by Bitto and Frühwirth-Schnatter (2019). Second, implementation is simple and requires relatively little additional input from the researcher. In fact, the prior heavily relies on a single structural hyperparameter that has to be specified with care, namely a .

The hyperparameter a influences the empirical properties of the proposed shrinkage prior along several important dimensions. Smaller values of a lead to heavy shrinkage on all elements of \mathbf{B} . To see this, note that lower values of a imply that more prior mass is placed on small values of ζ a priori. Similarly, when a is small, the Dirichlet prior places more mass on values of ϑ_j close to zero. Since lower values of ζ translate into thicker tails of the marginal prior on b_j , the specific choice of a not only influences the overall degree of shrinkage but also the tail behavior of the prior. Letting \tilde{p} denote the number of predictors, Bhattacharya et al. (2015) show that if a is specified as $\tilde{p}^{-(1+\Delta)}$ for any $\Delta > 0$ to be small, the DL prior displays excellent posterior contraction rates, and Pati et al. (2014), discuss the shrinkage properties of the proposed prior within the context of factor models. In our application, $\tilde{p} = K$ (when considering the total number of predictors) or $\tilde{p} = k$ (when considering the number of predictors per equation).

For the factor loadings we independently use a standard normally distributed prior on each element $\lambda_{ij} \sim \mathcal{N}(0, 1)$ for $i = 1, \dots, m$ and $j = 1, \dots, q$. In the empirical application (Section 6), we consider in addition the row-wise normal-gamma (NG Griffin & Brown, 2010) shrinkage prior discussed in Kastner (2019); that is $\lambda_{ij} | \tau_{ij}^2 \sim \mathcal{N}(0, \tau_{ij}^2)$, $\tau_{ij}^2 | v_i^2 \sim \mathcal{G}(a_\lambda, a_\lambda v_i^2 / 2)$, $\lambda_{ij}^2 \sim \mathcal{G}(c_\lambda, d_\lambda)$. Furthermore, we impose a normally distributed prior on the mean of the log-volatility $\mu_{\sigma_j} \sim \mathcal{N}(0, M_\mu)$ with M_μ denoting the prior variance, and the commonly employed beta distributed prior on the transformed persistence parameter of the log-volatility $\frac{\rho_{sj}+1}{2} \sim \mathcal{B}(a_0, b_0)$ for $s \in \{h, \sigma\}$ and $a_0, b_0 \in \mathbb{R}^+$ to ensure stationarity. Finally, we use a restricted gamma prior on the innovation variances in Equations (3) and (4), $\zeta_{sj}^2 \sim \mathcal{G}(\frac{1}{2}, \frac{1}{2\xi})$. Here, ξ is a hyperparameter used

to control the tightness of the prior. This choice, motivated in Frühwirth-Schnatter and Wagner (2010), implies that if the data are not informative on the degree of time variation of the log-volatilities then we do not bound ζ_{sj}^2 artificially away from zero, effectively applying more shrinkage than the standard inverted gamma prior.

3.2 | Full conditional posterior distributions

Conditional on the latent factors and the corresponding loadings, the model in Equation (1) can be cast as a system of m unrelated regression models for the elements in $\mathbf{z}_t = \mathbf{y}_t - \mathbf{A}\mathbf{f}_t$, labeled z_{it} , with heteroskedastic errors:

$$z_{it} = \mathbf{B}_{i\bullet}\mathbf{x}_t + \eta_{it}, \quad i = 1, \dots, m. \quad (8)$$

Here we let $\mathbf{B}_{i\bullet}$ denote the i th row of \mathbf{B} and η_{it} is the i th element of $\boldsymbol{\eta}_t$. The corresponding posterior distribution of $\mathbf{B}'_{i\bullet}$ is k -variate Gaussian:

$$\mathbf{B}'_{i\bullet} | \bullet \sim \mathcal{N}(\mathbf{b}_i, \mathbf{Q}_i), \quad (9)$$

with \bullet indicating that we condition on the remaining parameters and latent quantities of the model. The posterior variance and mean are given by

$$\mathbf{Q}_i = (\tilde{\mathbf{X}}_i' \tilde{\mathbf{X}}_i + \boldsymbol{\Phi}_i^{-1})^{-1}, \quad (10)$$

$$\mathbf{b}_i = \mathbf{Q}_i (\tilde{\mathbf{X}}_i' \tilde{\mathbf{z}}_i). \quad (11)$$

The diagonal prior covariance matrix of the coefficients related to the i th equation is given by $\boldsymbol{\Phi}_i$, the respective $k \times k$ diagonal submatrix of $\boldsymbol{\Phi} = \zeta \times \text{diag}(\psi_1 \vartheta_1^2, \dots, \psi_K \vartheta_K^2)$. Moreover, $\tilde{\mathbf{X}}_i$ is a $T \times k$ matrix with typical row t given by $\mathbf{X}_t / \sigma_{it}$ and $\tilde{\mathbf{z}}_i$ is a T -dimensional vector with the t th element given by z_{it} / σ_{it} . This normalization renders Equation (8) conditionally homoskedastic with standard normally distributed white noise errors.

The full conditional posterior distribution of ψ_j is inverse Gaussian:

$$\psi_j | \bullet \sim \text{iG}(\vartheta_j \zeta / |b_j|, 1), \quad j = 1, \dots, K. \quad (12)$$

The conditional posterior of the global shrinkage parameter ζ follows a generalized inverse Gaussian (GIG) distribution:

$$\zeta | \bullet \sim \text{GIG} \left(K(a-1), 1, 2 \sum_{j=1}^K |b_j| / \vartheta_j \right). \quad (13)$$

To draw from this distribution, we use the efficient algorithm of Hörmann and Leydold (2013). Moreover, we sample the scaling parameters ϑ_j by first sampling L_j from $L_j | \bullet \sim \text{GIG}(a-1, 1, 2|b_j|)$, and then setting $\vartheta_j = L_j / \sum_{i=1}^K L_i$.

The conditional posterior distributions of the factors are Gaussian and thus straightforward to draw from. The factor loadings are sampled using “deep interweaving” (see

Kastner et al., 2017), and the parameters in Equations (3) and (4) along the full histories of the latent log-volatilities are sampled as in Kastner and Frühwirth-Schnatter (2014) using the R-packages `factorstochvol` (Hosszejni & Kastner, 2019) and `stochvol` (Kastner, 2016).

Our MCMC algorithm iteratively draws from the conditional posterior distributions outlined above and discards the first J draws as burn-in. In terms of computational requirements, the single most intensive step is the simulation from the joint posterior of the autoregressive coefficients in \mathbf{B} . Because this step is implemented on an equation-by-equation basis, speed improvements relative to the standard approach are already quite substantial. However, note that if k is large (i.e., of the order of several thousands), even the commonly employed equation-by-equation sampling fails to deliver a sufficient amount of draws within a reasonable time window. Consequently, we outline an alternative algorithm to draw from a high-dimensional multivariate Gaussian distribution under a Bayesian prior that features a diagonal prior variance–covariance matrix in the upcoming section.

4 | COMPUTATIONAL ASPECTS

The typical approach to sampling from Equation (9) is based on the full system and simultaneously samples from the full conditional posterior of \mathbf{B} , implying that the corresponding posterior distribution is a K -dimensional Gaussian distribution with a $K \times K$ dimensional variance–covariance matrix. Under a nonconjugate prior, the computational difficulties arise from the need to invert the $K \times K$ variance–covariance matrix, which requires operations of order $O(m^6 p^3)$ under Gaussian elimination.

If a conjugate prior in combination with a constant (or vastly simplified heteroskedastic; see Carriero et al., 2016) specification of $\boldsymbol{\Omega}_t$ is used, the corresponding variance–covariance features a Kronecker structure which is computationally cheaper to invert and scales better in large dimensions. Specifically, the manipulations of the corresponding covariance matrix are of order $O(m^3 + k^3)$, a significant gain relative to the standard approach. However, this comes at a cost since all equations have to feature the same set of variables, the prior on the VAR coefficients has to be symmetric, and any stochastic volatility specification that preserves conjugacy is necessarily overly simplistic.

By contrast, recent studies emphasize the computational gains that arise from utilizing a framework that is based on equation-by-equation estimation. Carriero et al. (2019) and Koop et al. (2019) augment each equation of the system by either contemporaneous values of the endogenous variables of the preceding equations or the resid-

uals from the previous equations. Here, our approach renders the equations of the system conditionally independent by conditioning on the factors. From a computational perspective, the differences between using a factor model to disentangle the equations and an approach based on augmenting specific equations by quantities that aim to approximate covariance parameters are negligible. If we sample from Equation (9) directly, the computations involved are of order $O(mk^3) = O(m^4p^3)$. This already poses significant improvements relative to full system estimation.

One contribution of the present paper is the application of the algorithm proposed by Bhattacharya et al. (2016) and developed for univariate regression models under a global–local shrinkage prior. This algorithm is applied to each equation in the system and cycles through the following steps:

1. Sample independently $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}_k, \mathbf{\Phi}_i)$ and $\delta_i \sim \mathcal{N}(\mathbf{0}_T, \mathbf{I}_T)$.
2. Use \mathbf{u}_i and δ_i to construct $\mathbf{v}_i = \tilde{\mathbf{X}}_i \mathbf{u}_i + \delta_i$.
3. Solve $(\tilde{\mathbf{X}}_i \mathbf{\Phi}_i \tilde{\mathbf{X}}_i' + \mathbf{I}_T) \mathbf{w}_i = (\tilde{\mathbf{z}}_i - \mathbf{v}_i)$ for \mathbf{w}_i .
4. Set $\mathbf{B}'_{i\bullet} = \mathbf{u}_i + \mathbf{\Phi}_i \tilde{\mathbf{X}}_i' \mathbf{w}_i$.

This algorithm outperforms all competing variants discussed previously in situations where $k \gg T$, a situation commonly encountered when dealing with large VAR models. In such cases, steps 1–4 can be carried out using $O(pm^2T^2)$ floating point operations. In situations where $k \approx T$, the computational advantages relative to the standard equation-by-equation algorithm mentioned above are modest or even negative. However, note that the cost is quadratic in m and linear in p and thus scales much better when the number of endogenous variables and/or lags thereof is increased. More information on the empirical performance of our algorithm can be found in Section 6.4.

5 | SIMULATION STUDY

This section aims at comparing the performance of the DL prior with a range of commonly used alternatives. We investigate sparse, intermediate, and dense DGPs, where $T \in \{50, 100, 150, 200, 250\}$ and $m \in \{10, 20, 50, 100\}$. The probability of an off-diagonal entry to be nonzero is 0.01, 0.1, and 0.8 in each of the respective scenarios. In all scenarios, each intercept entry has a 0.1 probability of being nonzero and all diagonal elements are nonzero with probability 0.8. The nonzero elements are randomly generated from Gaussian distributions roughly tuned to yield stable VARs. More concretely, both the mean μ_I and the standard deviation σ_I of the intercept are set to 0.01, whereas mean and standard deviation of the diago-

nal (D) and the off-diagonal (O) elements are chosen as follows:

- Dense: 80% off-diagonal density level, $\mu_D = \sigma_D = 0.15$ and $\mu_O = \sigma_O = 0.01$.
- Intermediate: 10% off-diagonal density level, $\mu_D = \sigma_D = 0.15$ and $\mu_O = \sigma_O = 0.1$.
- Sparse: 1% off-diagonal density level, $\mu_D = \sigma_D = \mu_O = \sigma_O = 0.3$.

Concerning the errors, we use a single-factor SV specification. The factor loadings are generated from $\mathcal{N}(0.001, 0.001^2)$ to roughly match the above scaling. The AR(1) processes driving the idiosyncratic log-variances are assumed to have mean $\mu_{\sigma_i} = -12$ with persistences ρ_{σ_i} ranging from 0.85 to 0.98 and innovation standard deviations ζ_{σ_i} from 0.3 to 0.1. The process driving the factor log-variance is assumed to be highly persistent, with $\rho_{h1} = 0.99$ and $\zeta_{h1} = 0.1$.

For each of the 60 settings, we simulate 10 data sets. For each of these, we run our MCMC algorithm to obtain 2,000 posterior draws after a burn-in of 1,000. Consequently, the posterior means are compared to the true values and root mean squared errors (RMSEs) are computed. Finally, the median of each of these is reported in Table 1. Alongside the DL prior with weak ($a_{DL} = a = 1/2$) and strong ($a_{DL} = a = 1/k$ and $a_{DL} = a = 1/K$) shrinkage, we also consider the NG prior with a single global shrinkage parameter (see Huber & Feldkircher, 2019, for the exact specification) and a standard conjugate Minnesota prior with a single shrinkage parameter a_M , implemented by using dummy observations. For the NG prior we specify the prior on the global shrinkage parameter to induce heavy shrinkage (by setting both hyperparameters of the gamma prior equal to 0.01) and the prior controlling the excess kurtosis a_{NG} is set equal to 1, corresponding to the Bayesian Lasso (see Park & Casella, 2008), and $a_{NG} = 0.1$. The latter choice places significant prior mass around zero but at the same time leads to a heavy tailed marginal prior. Finally, we report RMSEs of the ordinary least squares (OLS) estimator (if it exists).

As is to be expected, Table 1 reveals strong to severe overfitting of OLS (corresponding to the posterior mode under a flat prior), which can be mitigated to a certain extent when the Minnesota prior with $a_M = 0.001$ is employed instead. Similarly, the DL prior with weak shrinkage ($a_{DL} = 1/2$) displays a tendency to overfit, in particular when T is small. By contrast, the more aggressive DL and NG shrinkage priors show superior performance. Overall, DL($1/k$) and NG(0.1) exhibit lowest RMSEs, where DL($1/k$) performs best in the sparse scenarios, NG(0.1) performs best in the intermediate settings, and no clear winner is to be found in the dense

T	m	Sparse				Intermediate				Dense			
		10	20	50	100	10	20	50	100	10	20	50	100
DL ($a_{DL}=1/2$)													
50		0.079	0.081	0.085	0.088	0.083	0.084	0.086	0.089	0.077	0.082	0.087	0.091
100		0.056	0.056	0.056	0.058	0.060	0.060	0.059	0.061	0.056	0.056	0.058	0.060
150		0.043	0.047	0.045	0.046	0.045	0.050	0.048	0.048	0.044	0.049	0.047	0.049
200		0.040	0.040	0.039	0.038	0.043	0.042	0.042	0.041	0.041	0.040	0.040	0.041
250		0.038	0.034	0.034	0.034	0.039	0.038	0.037	0.036	0.038	0.036	0.036	0.037
DL ($a_{DL}=1/k$)													
50		0.055	0.043	0.037	0.033	0.069	0.053	0.049	0.046	0.057	0.041	0.029	0.026
100		0.042	0.035	0.028	0.025	0.050	0.047	0.044	0.041	0.042	0.033	0.026	0.023
150		0.032	0.032	0.023	0.020	0.042	0.040	0.040	0.038	0.036	0.028	0.023	0.021
200		0.031	0.026	0.020	0.017	0.041	0.038	0.036	0.034	0.032	0.027	0.022	0.020
250		0.030	0.022	0.018	0.015	0.036	0.035	0.033	0.032	0.029	0.025	0.020	0.019
DL ($a_{DL}=1/K$)													
50		0.062	0.043	0.038	0.033	0.072	0.054	0.050	0.047	0.063	0.038	0.028	0.027
100		0.050	0.038	0.028	0.025	0.059	0.049	0.046	0.043	0.047	0.035	0.026	0.024
150		0.044	0.033	0.025	0.021	0.054	0.044	0.041	0.039	0.038	0.033	0.024	0.022
200		0.040	0.031	0.022	0.018	0.051	0.041	0.039	0.036	0.039	0.029	0.022	0.021
250		0.037	0.026	0.019	0.016	0.043	0.039	0.036	0.033	0.039	0.027	0.020	0.020
NG ($a_{NG}=1$)													
50		0.063	0.049	0.044	0.042	0.069	0.053	0.049	0.047	0.060	0.042	0.031	0.027
100		0.051	0.042	0.036	0.033	0.054	0.048	0.042	0.039	0.048	0.037	0.027	0.023
150		0.043	0.037	0.031	0.028	0.046	0.042	0.038	0.036	0.040	0.034	0.026	0.021
200		0.039	0.033	0.028	0.025	0.044	0.039	0.036	0.032	0.038	0.031	0.024	0.020
250		0.037	0.029	0.025	0.023	0.039	0.035	0.032	0.031	0.038	0.029	0.022	0.019
NG ($a_{NG}=0.1$)													
50		0.058	0.043	0.038	0.034	0.066	0.052	0.048	0.045	0.055	0.042	0.029	0.026
100		0.043	0.035	0.028	0.024	0.050	0.045	0.039	0.037	0.044	0.034	0.026	0.022
150		0.035	0.031	0.024	0.020	0.040	0.039	0.034	0.032	0.037	0.031	0.023	0.020
200		0.031	0.026	0.021	0.018	0.038	0.034	0.031	0.028	0.031	0.027	0.022	0.019
250		0.030	0.023	0.018	0.016	0.034	0.030	0.028	0.026	0.031	0.024	0.020	0.018
Minnesota ($a_M=0.001$)													
50		0.135	0.137	0.164	0.105	0.134	0.140	0.153	0.102	0.131	0.143	0.165	0.107
100		0.092	0.105	0.112	0.136	0.094	0.104	0.109	0.119	0.094	0.105	0.116	0.129
150		0.079	0.083	0.088	0.102	0.077	0.082	0.086	0.094	0.080	0.084	0.090	0.102
200		0.070	0.071	0.075	0.082	0.069	0.070	0.073	0.077	0.070	0.072	0.077	0.085
250		0.058	0.063	0.065	0.070	0.059	0.063	0.064	0.067	0.060	0.064	0.068	0.071
Minnesota ($a_M=0.0001$)													
50		0.067	0.052	0.048	0.046	0.073	0.054	0.049	0.045	0.062	0.037	0.028	0.022
100		0.063	0.050	0.047	0.045	0.069	0.052	0.047	0.044	0.059	0.036	0.028	0.022
150		0.061	0.049	0.045	0.044	0.066	0.049	0.045	0.042	0.056	0.035	0.028	0.021
200		0.061	0.047	0.044	0.043	0.065	0.048	0.044	0.041	0.054	0.034	0.027	0.022
250		0.058	0.046	0.042	0.041	0.059	0.046	0.042	0.040	0.053	0.033	0.027	0.021
OLS (if exists)													
50		0.158	0.205	DNE	DNE	0.158	0.211	DNE	DNE	0.160	0.211	DNE	DNE
100		0.106	0.128	0.163	DNE	0.107	0.126	0.157	DNE	0.110	0.128	0.165	DNE
150		0.088	0.099	0.112	0.155	0.087	0.098	0.109	0.155	0.090	0.101	0.115	0.167
200		0.080	0.078	0.087	0.104	0.079	0.078	0.085	0.106	0.080	0.080	0.091	0.114
250		0.065	0.071	0.078	0.092	0.066	0.070	0.077	0.086	0.067	0.071	0.081	0.096

TABLE 1 Median RMSEs stemming from 10 simulations per setting

context. Turning towards NG(1) and DL(1/K) we tend to observe acceptable but slightly inferior overall performance. The Minnesota prior with $a_M = 0.0001$ yields an extreme degree of shrinkage, translating into estimates of autoregressive coefficients that are very close to zero, irrespectively of the contribution from the like-

lihood. In that sense, it overshrinks most of the nonzero coefficients. Nevertheless, in scenarios with extremely low signal-to-noise ratios (such as the dense scenario with $T = 50$ and $m = 100$), this can be beneficial for the overall performance.

For further illustration, we showcase four exemplary scenarios in Figures A1–A4 in the Appendix.

6 | EMPIRICAL FORECASTING APPLICATION

In Section 6.1 we first summarize the data set adopted and present the model specification choices made. Section 6.2 estimates a simple one-factor model to outline the virtues of our proposed framework. Section 6.3 presents the main findings of our forecasting exercise and discusses the choice of the number of factors used for modeling the error covariance structure.

6.1 | Data, model specification and selection issues

The aim of the empirical application is to forecast a set of key US macroeconomic quantities. To this end, we use the quarterly data set provided by McCracken and Ng (2016), a variant of the well-known Stock and Watson (2011) data set for the USA.³ The data span the period ranging from 1959:Q1 to 2015:Q4. We include $m = 215$ quarterly time series, capturing information on 14 important segments of the economy and follow McCracken and Ng in transforming the data to be approximately stationary. Furthermore, we standardize each component series to have zero mean and variance one. In the empirical examples we include $p = 1$ lags of the endogenous variables.⁴ The hyperparameters are chosen as follows: $M_\mu = 10$, $a_0 = 20$, $b_0 = 1.5$, $\xi = 1$, $a_\lambda = 0.1$, $c_\lambda = d_\lambda = 1$.

6.2 | Some empirical key features of the model

To provide some intuition on how our modeling approach works in practice, we first estimate a simple one-factor model (i.e., $q = 1$) and investigate several features of our empirical model. In the next section we will perform an extensive forecasting exercise and discuss the optimal number of factors in terms of forecasting accuracy.

We start by inspecting the posterior distribution of Λ and assess what variables load heavily on the latent factor.

It is worth emphasizing that most quantities⁵ associated with real activity (i.e., industrial production and its components, gross domestic product (GDP) growth, employment measures) load heavily on the factor. Moreover, expectation measures, housing markets, equity prices, and spreads also load heavily on the joint factor.

To assess whether spikes in the volatility associated with the factor coincide with major economic events, the bottom panel of Figure 1 depicts the evolution of the posterior distribution of factor volatility over time. A few findings are worth mentioning. First, volatility spikes sharply during the mid-1970s, a period characterized by the first oil price shock and the bankruptcy of Franklin National Bank in 1974. After declining markedly during the second half of the 1970s, the shift in US monetary policy towards aggressively fighting inflation and the second oil price shock again translate into higher macroeconomic uncertainty. Note that from the mid-1980s onward we observe a general decline in macroeconomic volatility that lasts until the beginning of the 1990s. There we observe a slight increase in volatility possibly caused by the events surrounding the first Gulf War. The remaining years up to the beginning of the 2000s have been relatively unspectacular, with volatility levels being muted most of the time. In 2000/2001, volatility again increases due to the burst of the dot-com bubble and the 9/11 terrorist attacks. Finally, we observe marked spikes in volatility during recessionary episodes like the recent financial crisis in 2008.

Finally, we assess how well the DL prior with $a = 1/k$ performs in shrinking the coefficients in \mathbf{B} to zero. The top panel of Figure 2 depicts a heat map that gives a rough feeling of the size of each regression coefficient based on the posterior median of \mathbf{B} . The bottom panel of Figure 2 depicts the posterior interquartile range, providing some evidence on posterior uncertainty.⁶ The DL prior apparently succeeds in shrinking the vast majority of the approximately 50,000 coefficients towards zero. Even though not discussed in detail to conserve space, we note that at higher lag orders this very strong shrinkage effect is even more pronounced; see also Figures A5–A7 in the Appendix.

The top panel of Figure 3 displays the posterior median estimates when the shrinkage parameter a is chosen to be $1/2$ (cf. Bhattacharya et al., 2015, for a discussion of this choice). While $a = 1/2$ appears to provide a fair amount of shrinkage in other applications, for our huge dimensional example this prior exerts only relatively little shrinkage and tends to lead to overfitting. The diagonal pattern in the first lag appears here as well, but there is a considerable

³In addition to quarterly observations, McCracken and Ng (2016) also provide a subset of the data which is observed monthly. Of course, our method is analogously applicable to higher frequency observations. However, given that the computational cost of the Bhattacharya et al. (2016) approach is quadratic in T , the run-time gains of their approach in comparison to equation-by-equation estimation is then smaller and can, depending on the number of lags, even become negative.

⁴We have also experimented with higher lag orders and also found some evidence of signals at lag two for the data set at hand; see Figures A5–A7 in the Appendix for an illustration. However, out-of-sample predictive studies favored one lag only (cf. Section 6.3).

⁵Hereby we refer to the one-step-ahead forecast error related to a given time series.

⁶Since the corresponding posterior distribution is quite heavy tailed, using posterior standard deviations, while providing a qualitatively similar picture, tends to be slightly exaggerated.

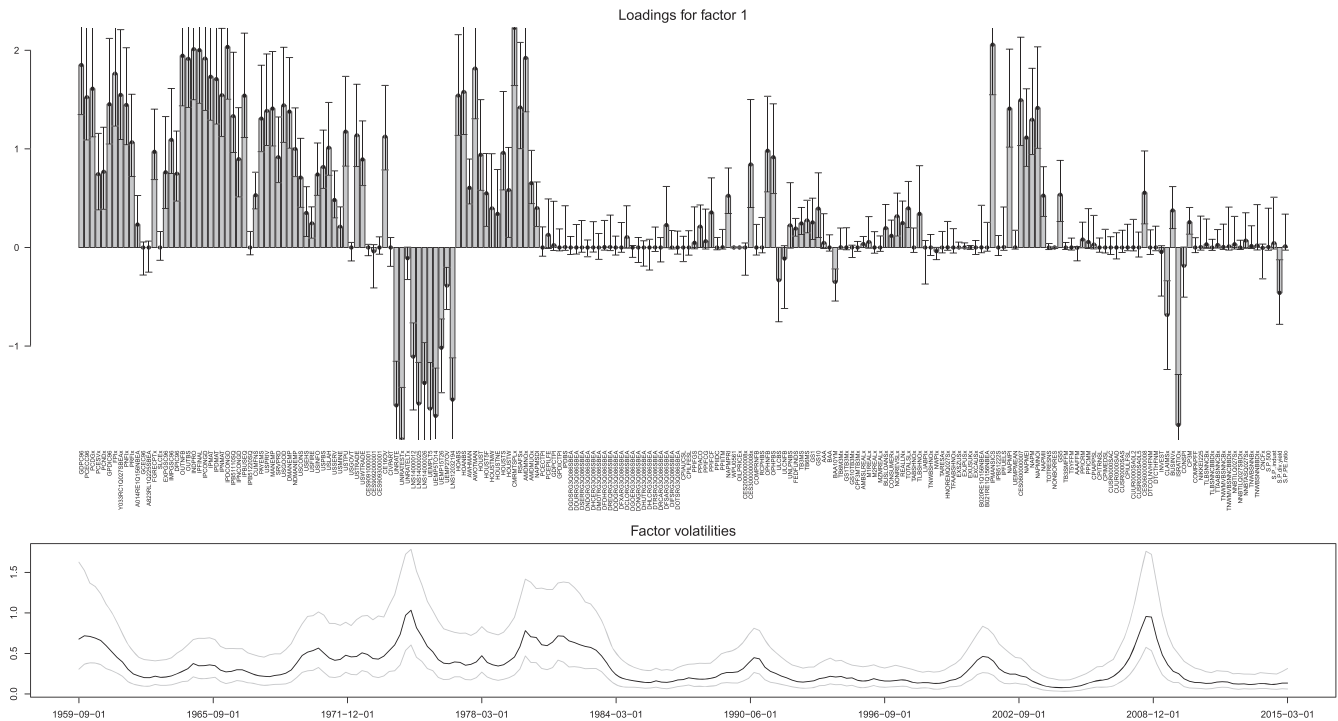


FIGURE 1 5th, 50th, and 95th posterior percentiles of factor loadings (upper panel) and factor volatility (lower panel)

amount of nonzero medians elsewhere. Correspondingly, the interquartile ranges visualized in the bottom panel of Figure 3 are also very large compared to those obtained with $a = 1/k$.

Interestingly, for selected time series measuring inflation (both consumer and producer price inflation) we find that lags of monetary aggregates are allowed to load on the respective inflation series. This result points towards a big advantage of our proposed prior relative to standard VAR priors in the Minnesota tradition: While these priors have been shown to work relatively well in huge dimensions (see Bańbura et al., 2010), they also display a tendency to overshrink when the overall tightness of the prior is integrated out in a Bayesian framework, effectively pushing the posterior distribution of \mathbf{B} towards the prior mean and thus ruling out patterns observed under the DL prior.

Inspection of the interquartile range also indicates that the proposed shrinkage prior succeeds in reducing posterior uncertainty markedly. Note that the pattern found for the posterior median of \mathbf{B} can also be found in terms of the posterior dispersion. We again observe that the coefficients associated with the first, own lag of a given variable are allowed to be nonzero whereas in most other cases the associated posterior is strongly concentrated around zero.

6.3 | Predictive evidence

We focus on forecasting gross domestic product (GDPC96), industrial production (INPRO), total nonfarm payroll (PAYEMS), civilian unemployment rate (UNRATE), new

privately owned housing units started (HOUST), consumer price index inflation (CPIAUCSL), producer price index for finished goods inflation (PPIFGS), effective federal funds rate (FEDFUNDS), 10-year Treasury constant maturity rate (GS10), US/UK exchange rate (EXUSUKx), and the S&P 500 (SP500). This choice includes the variables investigated by Koop et al. (2019) and some additional important macroeconomic indicators that are commonly monitored by practitioners, resulting in a total of 11 series.

To assess the forecasting performance of our model, we conduct a pseudo out-of-sample forecasting exercise with initial estimation sample ranging from 1959:Q3 to 1990:Q2. Based on this estimation period, we compute one-quarter-ahead predictive densities for the first period in the hold-out (i.e., 1990:Q3). After obtaining the corresponding predictive densities and evaluating the corresponding log-predictive likelihoods, we expand the estimation period and reestimate the model. This procedure is repeated 100 times until the final point of the full sample is reached. The quarterly scores obtained this way are then accumulated.

Our model with $q \in \{0, 1, \dots, 4\}$ factors is benchmarked against the prior model, a pure factor stochastic volatility (FSV) model with conditional mean equal to zero (i.e., $\mathbf{B} = \mathbf{0}_{m \times k}$). In what follows we label this specification FSV

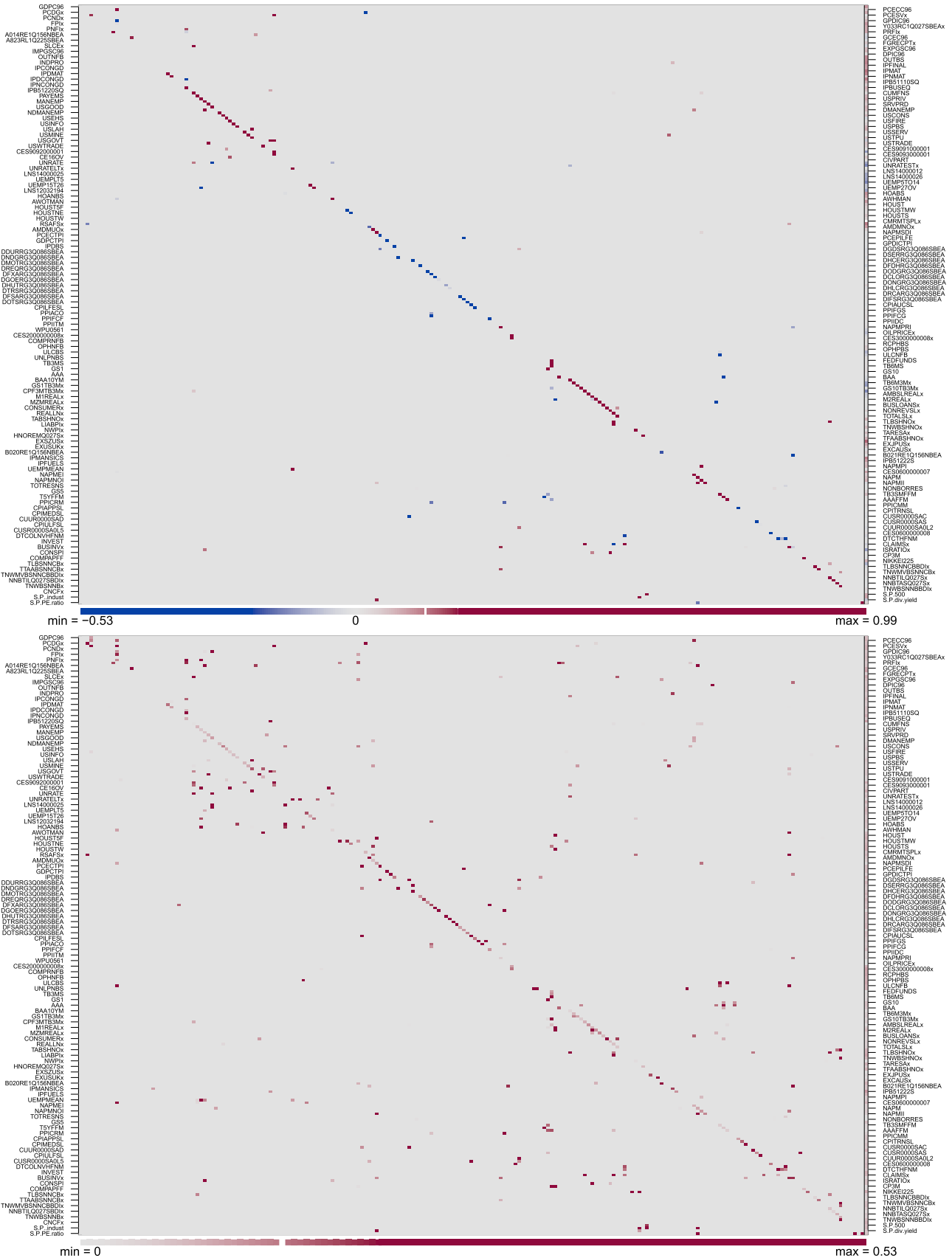


FIGURE 2 Posterior medians (top) and posterior interquartile ranges (bottom) of VAR coefficients, $a = 1/k = 1/216$ [Colour figure can be viewed at wileyonlinelibrary.com]

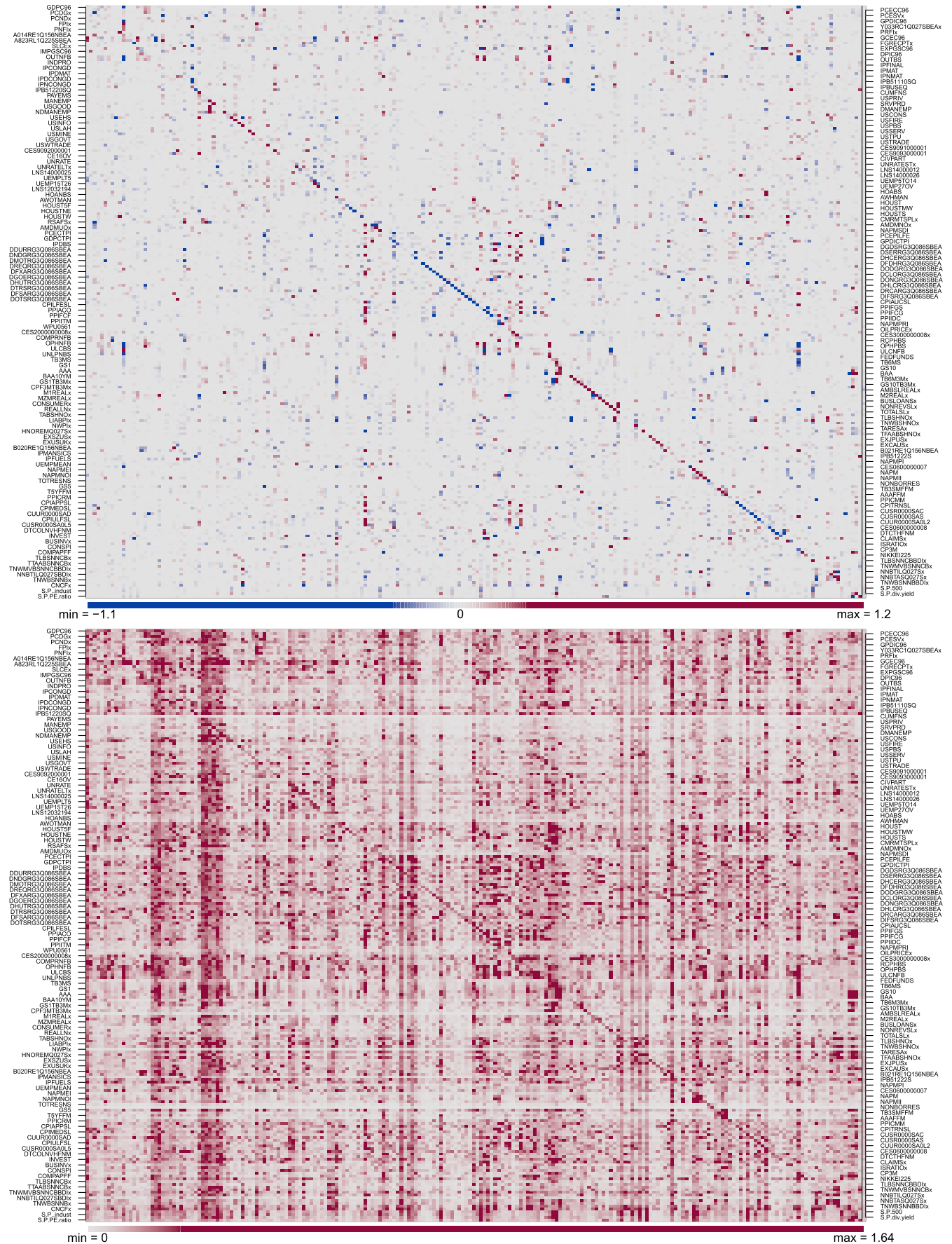


FIGURE 3 Posterior medians (top) and posterior interquartile ranges (bottom) of VAR coefficients, $a = 1/2$ [Colour figure can be viewed at wileyonlinelibrary.com]

0. To assess the merits of the proposed shrinkage prior vis-à-vis a Minnesota prior and an NG shrinkage prior we also include the models described in Section 5. Moreover, we include two models that impose the restriction that $\mathbf{A}_1 = \mathbf{I}_m$ and $\mathbf{A}_1 = 0.8 \times \mathbf{I}_m$, while \mathbf{A}_j for $j > 1$ are set equal to zero matrices in both cases. The first model, labeled FSV 1, assumes that the conditional mean of \mathbf{y}_t follows a random walk process, and the second specification, denoted by FSV 0.8, imposes the restriction that the variables in \mathbf{y}_t feature a rather strong degree of persistence but are stationary. The exercise serves to evaluate whether it pays to impose a VAR structure on the first moment of the joint density of our data and to assess how many factors are needed to obtain precise multivariate density predictions for our 11 variables of interest.

Overall log-predictive scores (LPSs) are summarized in Table 2. An immediate finding is that ignoring the error covariance structure (using zero factors) produces rather inaccurate forecasts for all models considered. While a single factor model improves predictive accuracy by a large margin, allowing for more factors (i.e., even more flexible modeling of the covariance structure) further increases the forecasting performance. For this specific exercise, we identify two or three factors to be a reasonable choice for most models when the joint log-predictive scores of the aforementioned variables are considered. We would like

to stress that this choice critically depends on the number of variables we include in our prediction set. If we focus attention on the marginal predictive densities (i.e., the univariate predictive densities obtained after integrating out the remaining elements in \mathbf{y}_t), we find that fewer or even no factors receive more support (see Table 3), whereas in the case of higher dimensional prediction sets more than two factors lead to more accurate density predictions (cf. Kastner et al., 2017, for an investigation of this issue in the context of a standard FSV model). As a general remark, we note that identifying the optimal number of factors in high-dimensional FSV models is a challenging problem in practice. Using the deviance information criterion (DIC; cf. Chan & Grant, 2016) may be an option but is likely to be unstable in very high dimensions. The approach adopted in the paper at hand, namely the decomposition of the marginal likelihood into predictive likelihoods (cf. Geweke & Amisano, 2010) tends to be more stable, in particular when interest is placed on predicting subsets only. Moreover, it can be trivially parallelized, thus becoming computationally feasible on high-performance computing infrastructures.

Considering forecasting accuracy across models reveals that our proposed VAR(1)-FSV with a DL(1/k) prior displays excellent forecasting capabilities, outperforming all competitors. Among the VAR(1) models, DL(1/K) and

TABLE 2 Average log-predictive scores for the number of factors $q \in \{0, 1, \dots, 5\}$ in various VAR-FSV specifications as well as pure FSV models

	0	1	2	3	4	5
VAR(1)-FSV DL(1/2)	-14.79	-13.45	-12.74	-12.53	-11.85	-11.44
VAR(1)-FSV DL(1/k)	-10.51	-9.75	-9.15	-9.15	-9.17	-9.43
VAR(1)-FSV DL(1/K)	-10.54	-9.63	-9.22	-9.26	-9.22	-9.38
VAR(1)-FSV NG(1)	-10.76	-10.26	-9.72	-9.60	-9.68	-9.86
VAR(1)-FSV NG(0.1)	-10.55	-9.95	-9.31	-9.26	-9.47	-9.46
VAR(1)-FSV Min(0.01)	-10.99	-10.37	-9.89	-9.88	-9.98	-10.24
VAR(1)-FSV Min(0.001)	-12.15	-11.21	-11.08	-10.62	-10.56	-10.83
VAR(2)-FSV DL(1/k)	-10.50	-9.88	-9.22	-9.38	-9.28	-9.43
VAR(2)-FSV NG(0.1)	-10.47	-10.04	-9.44	-9.29	-9.40	-9.49
VAR(5)-FSV DL(1/k)	-10.70	-9.97	-9.78	-9.21	-9.44	-9.53
VAR(5)-FSV NG(0.1)	-10.61	-10.01	-9.64	-9.51	-9.50	-9.52
FSV 0	-12.08	-11.08	-11.02	-10.66	-10.61	-10.71
FSV 0.8	-11.78	-11.32	-11.06	-10.87	-10.96	-11.25
FSV 1	-11.71	-11.28	-10.95	-10.95	-11.07	-11.18

Note. Estimation and prediction are conducted on all $m = 215$ component series; the predictive density is then evaluated on the set of 11 variables of interest. Larger numbers indicate better joint predictive density performance.

TABLE 3 Average univariate log-predictive scores for inflation (CPIAUCSL), short-term interest rates (FEDFUNDS), and output growth (GDPC96) with $q \in \{0, 1, 2\}$ factors

	CPIAUCSL			FEDFUNDS			GDPC96		
	0	1	2	0	1	2	0	1	2
VAR(1)-FSV DL(1/k)	-1.03	-1.11	-1.13	-1.26	-1.26	-1.24	0.08	0.05	0.03
VAR(1)-FSV NG(0.1)	-1.00	-1.07	-1.10	-1.28	-1.26	-1.22	-0.10	-0.12	-0.15
VAR(2)-FSV DL(1/k)	-1.02	-1.12	-1.14	-1.25	-1.27	-1.22	0.04	-0.01	-0.01
VAR(2)-FSV NG(0.1)	-1.00	-1.11	-1.12	-1.26	-1.23	-1.25	-0.14	-0.16	-0.14
VAR(5)-FSV DL(1/k)	-1.05	-1.16	-1.16	-1.29	-1.26	-1.26	-0.02	-0.10	-0.13
VAR(5)-FSV NG(0.1)	-1.00	-1.09	-1.13	-1.29	-1.27	-1.25	-0.19	-0.18	-0.21

NG(0.1) also do well, and the Bayesian Lasso (NG(1)) as well as the Minnesota prior with medium shrinkage (Min(0.01)) show decent performance. Clearly, DL(1/2) overfits and Min(0.001) overshrinks. Note that higher lag orders seem rarely to increase predictive accuracy. However, comparing the differences between the benchmark pure FSV models and the VAR-FSV models considered, we find that explicitly modeling the conditional mean improves the forecasting accuracy in practically all cases.

To investigate whether forecasting performance is homogeneous over time, Figure 4 visualizes the cumulative LPSs relative to the zero-factor FSV model over time. The benefit of the flexible SV structure in the VAR residuals is particularly pronounced during the 2008 financial crisis which can be seen by comparing the solid lines to the broken lines. During this period, time-varying covariance modeling appears to be of great importance and the performance of models that ignore contemporaneous dependence deteriorates. This finding is in line with Kastner (2019), who reports analogous results for US asset returns. The increase in predictive accuracy can be traced back to the fact that within an economic downturn the correlation structure of our data set changes markedly, with most indicators that measure real activity sharply declining in lockstep. A model that takes contemporaneous cross-variable linkages seriously is thus able to fully exploit such behavior, which in turn improves predictions.

Up to this point, we have focused exclusively on the joint performance of our model for the specific set of variables considered. To gain a deeper understanding on how our model performs for relevant selected quantities, Table 3 displays marginal LPSs for the two most promising prior specifications with one, two, and five lags. The variables we consider are inflation (CPIAUCSL), short-term interest rates (FEDFUNDS), and output growth (GDPC96).

In contrast to the findings based on joint LPSs, we observe that models without a factor structure tend to perform better than models that set $q > 0$, with the exception of interest rates where all models predict more or less equally badly. This finding corroborates our conjecture stated above, implying that if the set of focus variables is subsequently enlarged, more factors are necessary in order to obtain precise density predictions. Here, we only focus on marginal model performance, implying that for each variable, contemporaneous relations between the elements in \mathbf{y}_t are integrated out. This, in turn, implies that the additional gain in model flexibility is offset by the comparatively larger number of parameters. Concerning the difference between VAR priors, it appears that NG slightly outperforms DL for inflation, whereas DL is superior when it comes to predicting output growth.

6.4 | A note on the computational burden

Even though the efficient sampling schemes outlined in this paper help to overcome absolutely prohibitive com-

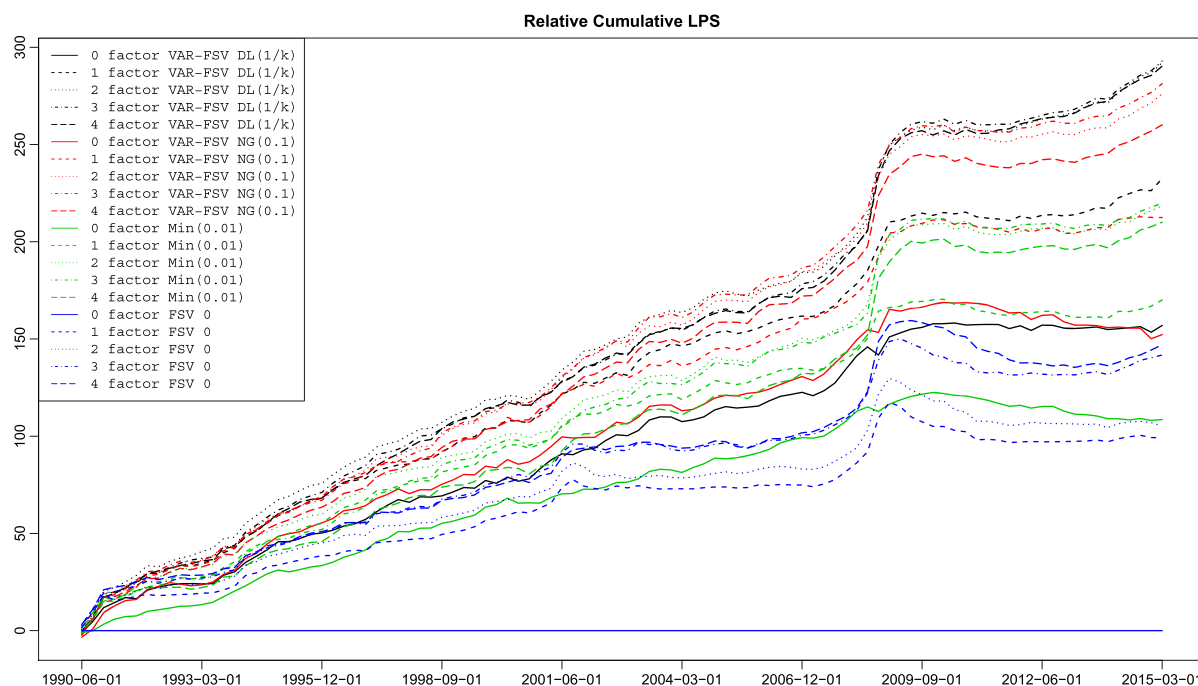
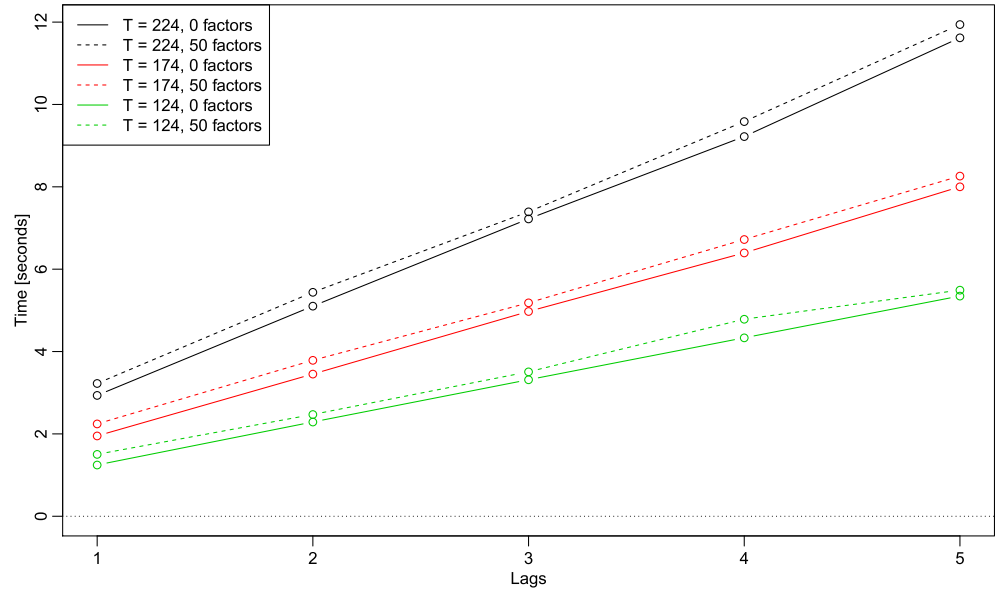


FIGURE 4 Cumulative log predictive scores, relative to a zero-mean model with independent stochastic volatility components for all component series. Higher values correspond to better one-quarter-ahead density predictions up to the corresponding point in time [Colour figure can be viewed at wileyonlinelibrary.com]

FIGURE 5 Empirical CPU times for each MCMC iteration on a standard laptop computer using one core. Time series lengths: $T \in \{124, 174, 224\}$; numbers of latent factors: $q \in \{0, 50\}$ [Colour figure can be viewed at wileyonlinelibrary.com]



putational burdens, the CPU time needed to perform fully Bayesian inference in a model of this size can still be considered substantial. In what follows we shed light on the estimation time required and how it is related to the length of the time series T , the lag length p , and the number of latent factors $q \in \{0, 50\}$. Figure 5 shows the time needed to perform a single draw from the joint posterior distribution of the $215 + 215^2 p$ coefficients and their corresponding $2(215 + 215^2 p) + 1$ auxiliary shrinkage quantities, the qT factor realizations and the associated $215q$ loadings, alongside $(T + 1)(215 + q)$ latent volatilities with their corresponding $645 + 2q$ parameters. This amounts to 166,841 random draws for the smallest model considered (one lag, no factors, $T = 124$) and 776,341 random draws for the largest model (5 lags, 50 factors, $T = 224$) at each MCMC iteration.

As mentioned above, the computation time rises approximately linearly with the number of lags included. Dotted lines indicate the time in seconds needed to perform a single draw from a model with 50 factors included, while solid lines refer to the time needed to estimate a model without factors and a diagonal time-varying variance–covariance matrix Ω_t . Interestingly, the additional complexity when moving from a model without factors to a highly parametrized model with 50 factors appears to be negligible, increasing the time needed by a fraction of a second on average. The important role of the length of the sample can be seen by comparing the green, red, and black lines. The time necessary to perform a simple MCMC draw quickly rises with the length of our sample, consistent with the statements made in Section 4. This feature of our algorithm, however, is convenient especially when researchers are interested in combining many short time series or performing recursive forecasting based on a tiny initial estimation sample.

7 | CLOSING REMARKS

In this paper we propose an alternative route to estimate huge-dimensional VAR models that allow for time variation in the error variances. The Dirichlet–Laplace prior, a recent variant of a global–local shrinkage prior, enables us to heavily shrink the parameter space towards the prior model while providing enough flexibility that individual regression coefficients are allowed to be unrestricted. This prior setup alleviates overfitting issues generally associated with large VAR models. To cope with computational issues we assume that the one-step-ahead forecast errors of the VAR feature a factor stochastic volatility structure that enables us to perform equation-by-equation estimation, conditional on the loadings and the factors. Since posterior simulation of each equation's autoregressive parameters involves manipulating large matrices, we implement an alternative recent algorithm that improves upon existing methods by large margins, rendering a fully fledged Bayesian estimation of truly huge systems possible.

In an empirical application we first present various key features of our approach based on a single-factor model. This single factor, which summarizes the joint dynamics of the VAR errors, can be interpreted as an uncertainty measure that closely tracks observed factors such as the volatility index. The question whether such a simplistic structure proves to be an adequate representation of the time-varying covariance matrix naturally arises, and we therefore provide a detailed forecasting exercise to evaluate the merits of our approach relative to the prior model and a set of competing models with a different number of latent factors in the errors.

Finally, three potential extensions are worth mentioning. First, given the fact that systematic and in-depth

empirical comparisons of the various recently developed roads towards handling high-dimensional VARs with time-varying contemporaneous covariance in a Bayesian framework (VAR-FSV, VAR-Cholesky-SV, compressed VAR-SV, etc.) are still missing and it is not clear whether one of these models turns out to dominate the others for all points in time, one could consider averaging/selecting dynamically. Second, note that it is trivial to relax the assumption of symmetry for the DL components. In the context of VARs, this might be of particular interest for distinguishing diagonal (a_D large) from off-diagonal (a_O small) elements in the spirit of the Minnesota prior or increasing the amount of shrinkage with increasing lag order (cf. Huber & Feldkircher, 2019, for a similar setup in the context of the normal-gamma shrinkage prior). Third, we would like to stress that our approach could also be used to estimate huge-dimensional time-varying parameter VAR models with stochastic volatility. To cope with the computational difficulties associated with the vast state space, a possible approach could be to rely on an additional layer of hierarchy that imposes a (dynamic) factor structure on the time-varying autoregressive coefficients in the spirit of Eisenstat et al. (2018) and thus reduce the computational burden considerably.

ACKNOWLEDGMENTS

The authors acknowledge funding from the Austrian Science Fund (FWF) for the project “High-dimensional statistical learning: New methods to advance economic and sustainability policies” (ZK 35), jointly carried out by WU Vienna University of Economics and Business, Paris Lodron University Salzburg, TU Wien, and the Austrian Institute of Economic Research (WIFO).

DATA AVAILABILITY STATEMENT

The data analyzed in this manuscript are available from the corresponding author on request.

REFERENCES

- Aguilar, O., & West, M. (2000). Bayesian dynamic factor models and portfolio allocation. *Journal of Business and Economic Statistics*, 18(3), 338–357. <https://doi.org/10.2307/1392266>
- Ahelegbey, D. F., Billio, M., & Casarin, R. (2016). Sparse graphical vector autoregression: A Bayesian approach. *Annals of Economics and Statistics*, 123–124, 333–361. <https://doi.org/10.15609/annaeconstat2009.123-124.0333>
- Ankargren, S., Unosson, M., & Yang, Y. (2019). A flexible mixed-frequency vector autoregression with a steady-state prior. arXiv: 1911.09151 [econ.EM].
- Bañbura, M., Giannone, D., & Reichlin, L. (2010). Large Bayesian vector autoregressions. *Journal of Applied Econometrics*, 25(1), 71–92. <https://doi.org/10.1002/jae.1137>
- Bhattacharya, A., Chakraborty, A., & Mallick, B. K. (2016). Fast sampling with Gaussian scale mixture priors in high-dimensional regression. *Biometrika*, 103(4), 985–991. <https://doi.org/10.1093/biomet/asw042>
- Bhattacharya, A., Pati, D., Pillai, N. S., & Dunson, D. B. (2015). Dirichlet–Laplace priors for optimal shrinkage. *Journal of the American Statistical Association*, 110(512), 1479–1490. <https://doi.org/10.1080/01621459.2014.960967>
- Bitto, A., & Frühwirth-Schnatter, S. (2019). Achieving shrinkage in a time-varying parameter model framework. *Journal of Econometrics*, 210(1), 75–97. <https://doi.org/10.1016/j.jeconom.2018.11.006>
- Carriero, A., Clark, T. E., & Marcellino, M. (2016). Common drifting volatility in large Bayesian VARs. *Journal of Business and Economic Statistics*, 34(3), 375–390. <https://doi.org/10.1080/07350015.2015.1040116>
- Carriero, A., Clark, T. E., & Marcellino, M. (2019). Large Bayesian vector autoregressions with stochastic volatility and non-conjugate priors. *Journal of Econometrics*, 212(1), 137–154. <https://doi.org/10.1016/j.jeconom.2019.04.024>
- Carvalho, C. M., Polson, N. G., & Scott, J. G. (2010). The horseshoe estimator for sparse signals. *Biometrika*, 97(2), 465–480. <https://doi.org/10.1093/biomet/asq017>
- Chan, J. C. C., & Grant, A. L. (2016). Fast computation of the deviance information criterion for latent variable models. *Computational Statistics and Data Analysis*, 100, 847–859. <https://doi.org/10.1016/j.csda.2014.07.018>
- Clark, T. E. (2011). Real-time density forecasts from Bayesian vector autoregressions with stochastic volatility. *Journal of Business and Economic Statistics*, 29(3), 327–341. <https://doi.org/10.1198/jbes.2010.09248>
- Clark, T. E., & Ravazzolo, F. (2015). Macroeconomic forecasting performance under alternative specifications of time-varying volatility. *Journal of Applied Econometrics*, 30(4), 551–575. <https://doi.org/10.1002/jae.2379>
- Cogley, T., & Sargent, T. J. (2001). Evolving post-World War II US inflation dynamics. *NBER macroeconomics annual* (pp. 331–373), Vol. 16. Cambridge, MA: National Bureau of Economic Research.
- Davis, R. A., Zang, P., & Zheng, T. (2016). Sparse vector autoregressive modeling. *Journal of Computational and Graphical Statistics*, 25(4), 1077–1096. <https://doi.org/10.1080/10618600.2015.1092978>
- Doan, T., Litterman, R. B., & Sims, C. A. (1984). Forecasting and conditional projection using realistic prior distributions. *Econometric Reviews*, 3(1), 1–100. <https://doi.org/10.1080/07474938408800053>
- Eisenstat, E., Chan, J. C. C., & Strachan, R. W. (2018). Reducing dimensions in a large TVP-VAR. (*Working Paper 18-37*). Waterloo ON, Canada: Rimini Centre for Economic Analysis.
- Follett, L., & Yu, C. (2019). Achieving parsimony in Bayesian vector autoregressions with the horseshoe prior. *Econometrics and Statistics*, 11, 130–144. <https://doi.org/10.1016/j.ecosta.2018.12.004>
- Frühwirth-Schnatter, S., & Wagner, H. (2010). Stochastic model specification search for Gaussian and partial non-Gaussian state space models. *Journal of Econometrics*, 154(1), 85–100. <https://doi.org/10.1016/j.jeconom.2009.07.003>
- George, E. I., Sun, D., & Ni, S. (2008). Bayesian stochastic search for VAR model restrictions. *Journal of Econometrics*, 142(1), 553–580. <https://doi.org/10.1016/j.jeconom.2007.08.017>

- Geweke, J., & Amisano, G. (2010). Comparing and evaluating Bayesian predictive distributions of asset returns. *International Journal of Forecasting*, 26(2), 216–230. <https://doi.org/10.1016/j.ijforecast.2009.10.007>
- Griffin, J. E., & Brown, P. J. (2010). Inference with normal-gamma prior distributions in regression problems. *Bayesian Analysis*, 5(1), 171–188. <https://doi.org/10.1214/10-BA507>
- Hörmann, W., & Leydold, J. (2013). Generating generalized inverse Gaussian random variates. *Statistics and Computing*, 24(4), 1–11. <https://doi.org/10.1007/s11222-013-9387-3>
- Hosszejni, D., & Kastner, G. (2019). Modeling univariate and multivariate stochastic volatility in R with stochvol and factorstochvol. R package vignette. Retrieved from <https://CRAN.R-project.org/package=factorstochvol/vignettes/paper.pdf>
- Huber, F., & Feldkircher, M. (2019). Adaptive shrinkage in Bayesian vector autoregressive models. *Journal of Business and Economic Statistics*, 37(1). <https://doi.org/10.1080/07350015.2016.1256217>
- Huber, F., Kastner, G., & Feldkircher, M. (2019). Should I stay or should I go? A latent threshold approach to large-scale mixture innovation models. *Journal of Applied Econometrics*, 34(5), 621–640. <https://doi.org/10.1002/jae.2680>
- Kastner, G. (2016). Dealing with stochastic volatility in time series using the R package stochvol. *Journal of Statistical Software*, 69(5), 1–30. <https://doi.org/10.18637/jss.v069.i05>
- Kastner, G. (2019). Sparse Bayesian time-varying covariance estimation in many dimensions. *Journal of Econometrics*, 210(1), 98–115. <https://doi.org/10.1016/j.jeconom.2018.11.007>
- Kastner, G., & Frühwirth-Schnatter, S. (2014). Ancillarity-sufficiency interweaving strategy (ASIS) for boosting MCMC estimation of stochastic volatility models. *Computational Statistics and Data Analysis*, 76, 408–423. <https://doi.org/10.1016/j.csda.2013.01.002>
- Kastner, G., Frühwirth-Schnatter, S., & Lopes, H. F. (2017). Efficient Bayesian inference for multivariate factor stochastic volatility models. *Journal of Computational and Graphical Statistics*, 26(4), 905–917. <https://doi.org/10.1080/10618600.2017.1322091>
- Koop, G. M. (2013). Forecasting with medium and large Bayesian VARs. *Journal of Applied Econometrics*, 28(2), 177–203. <https://doi.org/10.1002/jae.1270>
- Koop, G., & Korobilis, D. (2013). Large time-varying parameter VARs. *Journal of Econometrics*, 177(2), 185–198. <https://doi.org/10.1016/j.jeconom.2013.04.007>
- Koop, G., Korobilis, D., & Pettenuzzo, D. (2019). Bayesian compressed vector autoregressions. *Journal of Econometrics*, 210(1), 135–154. <https://doi.org/10.1016/j.jeconom.2018.11.009>
- Koop, G., Leon-Gonzalez, R., & Strachan, R. W. (2009). On the evolution of the monetary policy transmission mechanism. *Journal of Economic Dynamics and Control*, 33(4), 997–1017. <https://doi.org/10.1016/j.jedc.2008.11.003>
- Korobilis, D., & Pettenuzzo, D. (2019). Adaptive hierarchical priors for high-dimensional vector autoregressions. *Journal of Econometrics*, 212(1). <https://doi.org/10.1016/j.jeconom.2019.04.029>
- Litterman, R. B. (1986). Forecasting with Bayesian vector autoregressions: Five years of experience. *Journal of Business and Economic Statistics*, 4(1), 25–38. <https://doi.org/10.2307/1391384>
- McCracken, M. W., & Ng, S. (2016). FRED-MD: A monthly database for macroeconomic research. *Journal of Business and Economic Statistics*, 34(4), 574–589. <https://doi.org/10.1080/07350015.2015.1086655>
- Park, T., & Casella, G. (2008). The Bayesian Lasso. *Journal of the American Statistical Association*, 103(452), 681–686. <https://doi.org/10.1198/016214508000000337>
- Pati, D., Bhattacharya, A., Pillai, N. S., & Dunson, D. B. (2014). Posterior contraction in sparse Bayesian factor models for massive covariance matrices. *Annals of Statistics*, 42(3), 1102–1130. <https://doi.org/10.1214/14-AOS1215>
- Pitt, M. K., & Shephard, N. (1999). Time-varying covariances: A factor stochastic volatility approach, Bayesian Statistics 6: Proceedings of the Sixth Valencia International Meeting, pp. 547–570.
- Polson, N. G., & Scott, J. G. (2011). Shrink globally, act locally: Sparse Bayesian regularization and prediction. In *Bayesian Statistics 9: Proceedings of the Ninth Valencia International Meeting* (Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P., Heckerman, D., Smith, A. F. M., & West, M., Eds.), Clarendon Press, Oxford, UK, pp. 501–538.
- Primiceri, G. E. (2005). Time varying structural vector autoregressions and monetary policy. *Review of Economic Studies*, 72(3), 821–852. <https://doi.org/10.1111/j.1467-937X.2005.00353.x>
- Rockova, V., & McAlinn, K. (2017). Dynamic variable r spike-and-slab process priors. arXiv: 1708.00085 [stat.ME].
- Sims, C. A., & Zha, T. (1998). Bayesian methods for dynamic multivariate models. *International Economic Review*, 39(4), 949–968. <https://doi.org/10.2307/2527347>
- Sims, C. A., & Zha, T. (2006). Were there regime switches in US monetary policy? *American Economic Review*, 96(1), 54–81. <https://doi.org/10.1257/000282806776157678>
- Stock, J. H., & Watson, M. W. (2005). Understanding changes in international business cycle dynamics. *Journal of the European Economic Association*, 3(5), 968–1006. <https://doi.org/10.1162/1542476054729446>
- Stock, J. H., & Watson, M. W. (2011). Dynamic factor models, *Oxford handbook of economic forecasting* (pp. 35–59). Oxford, UK: Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780195398649.013.0003>

AUTHOR BIOGRAPHIES

Gregor Kastner researches the Bayesian modeling of economic time series, in particular the efficient estimation of (factor) stochastic volatility and vector autoregressive models. Currently, he serves as coordinator and principal investigator on the project “High-dimensional statistical learning: New methods to advance economic and sustainability policies” funded by the Austrian Science Fund (FWF). He is on the board of the European Seminar on Bayesian Econometrics (ESOB) and works as reproduction and copy editor for the Journal of Statistical Software. His research appeared in outlets such as the Journal of Econometrics, the Journal of Applied Econometrics, the Journal of Computational and Graphical Statistics, Computational

Statistics & Data Analysis, and the Journal of Statistical Software.

Florian Huber is Professor of Empirical Macroeconomics at the University of Salzburg. His research interests focus primarily on macroeconomic modeling and forecasting, time series analysis, and Bayesian econometrics. He has published in journals such as the Journal of Business and Economic Statistics, the Journal of Applied Econometrics, the European Economic Review, the Journal of Economic Dynamics and Control, the Journal of the Royal Statistical Society: Series A, the Journal of Banking and Finance, the Journal of Economic Behavior and Organization, the Oxford Bulletin of Economics and Statistics and the International Journal of Forecasting, inter alia.

How to cite this article: Kastner G, Huber F. Sparse Bayesian vector autoregressions in huge dimensions. *Journal of Forecasting*. 2020;1–24. <https://doi.org/10.1002/for.2680>

APPENDIX A: FURTHER ILLUSTRATIONS

First, we showcase four selected data-generating scenarios (small + sparse, small + dense, large + sparse, large + dense) and visualize the posterior distribution of the VAR coefficients under seven different prior choices in Figures A1–A4. For a comprehensive overview, see Table 1 in the main part of the paper.

Second, we illustrate results from a VAR(2)-FSV DL(1/ k) and a VAR(5)-FSV DL(1/ k) model for the US data in Figures A5–A7.

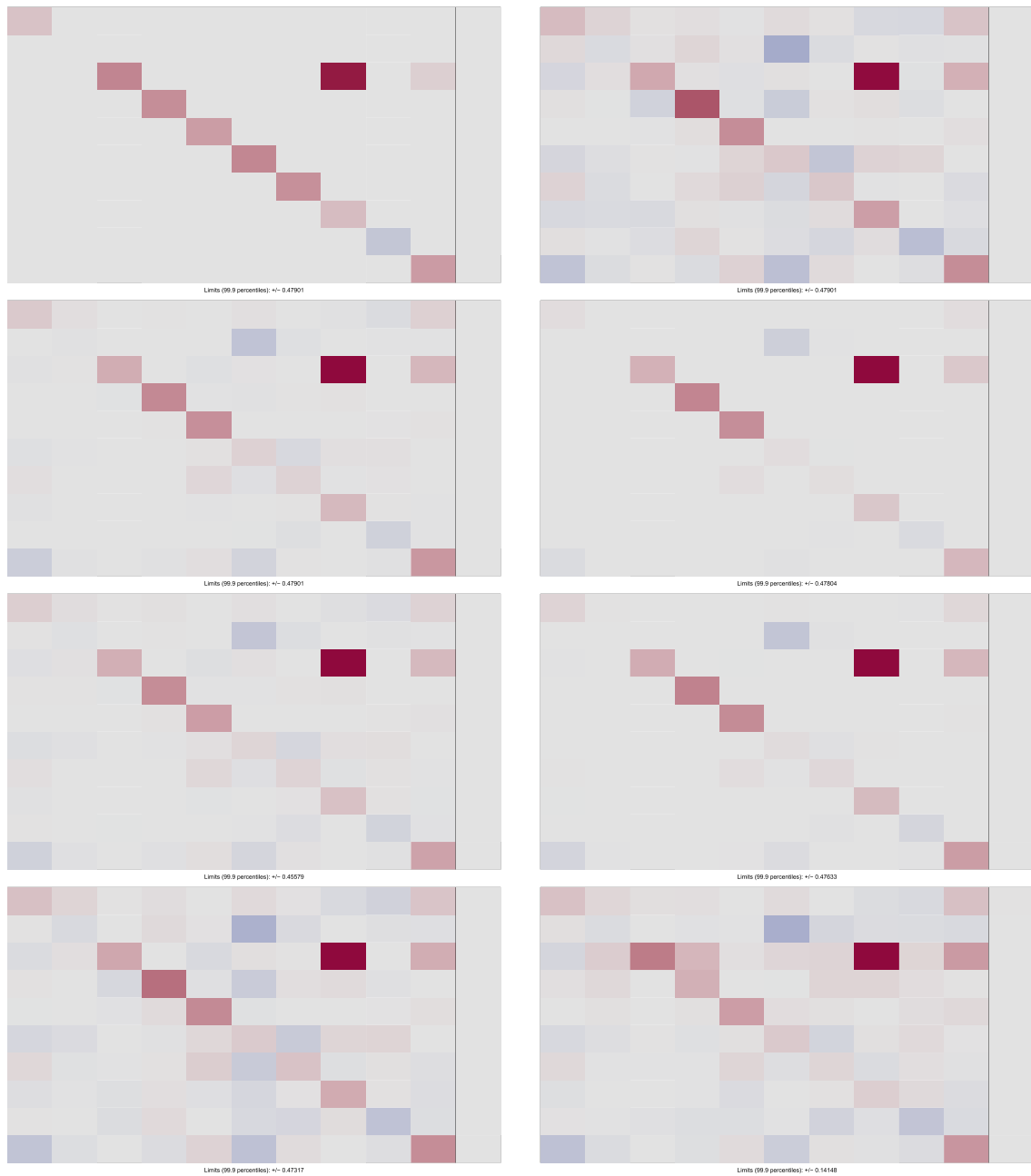


FIGURE A1 Exemplary visualization of the true and estimated VAR coefficients in the *sparse* scenario where $T = 250$ and $m = 10$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/k = 1/11$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/1,000$ (left) and $a_M = 1/10,000$ (right)

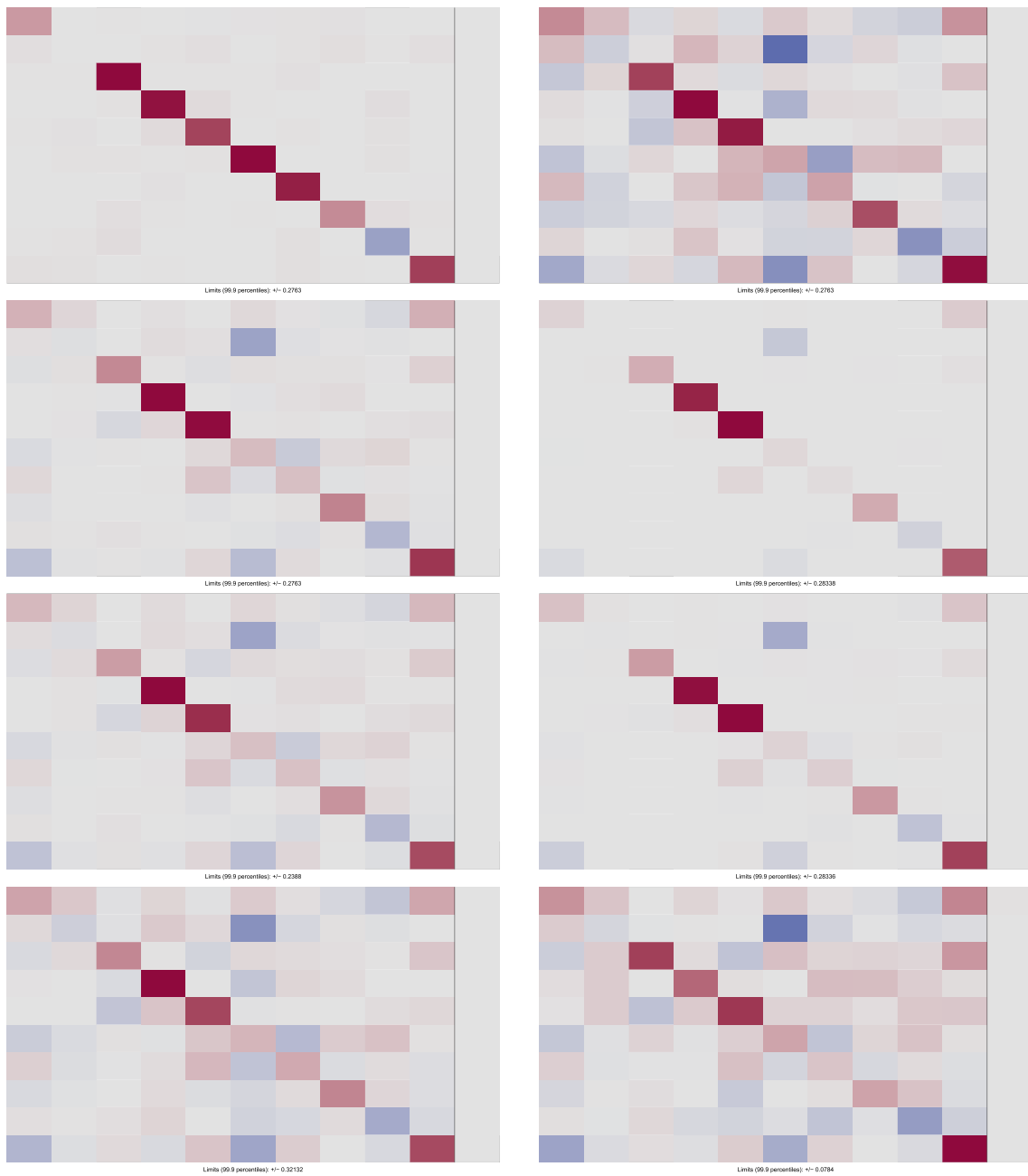


FIGURE A2 Exemplary visualization of the true and estimated VAR coefficients in the *dense* scenario where $T = 250$ and $m = 10$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/k = 1/11$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/1,000$ (left) and $a_M = 1/10,000$ (right) [Colour figure can be viewed at wileyonlinelibrary.com]

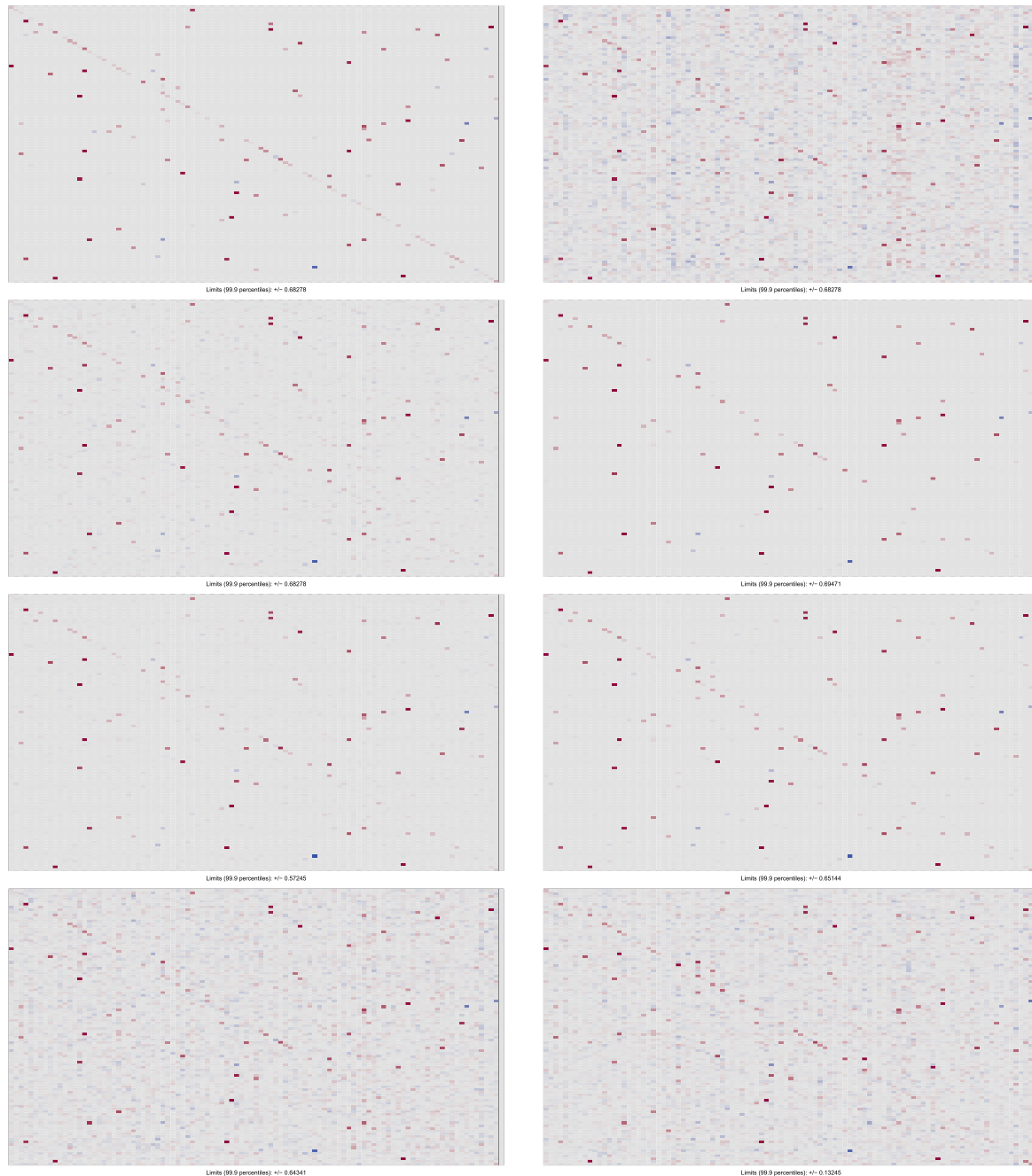


FIGURE A3 Exemplary visualization of the true and estimated VAR coefficients in the *sparse* scenario where $T = 250$ and $m = 100$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/k = 1/101$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/1,000$ (left) and $a_M = 1/10,000$ (right) [Colour figure can be viewed at wileyonlinelibrary.com]

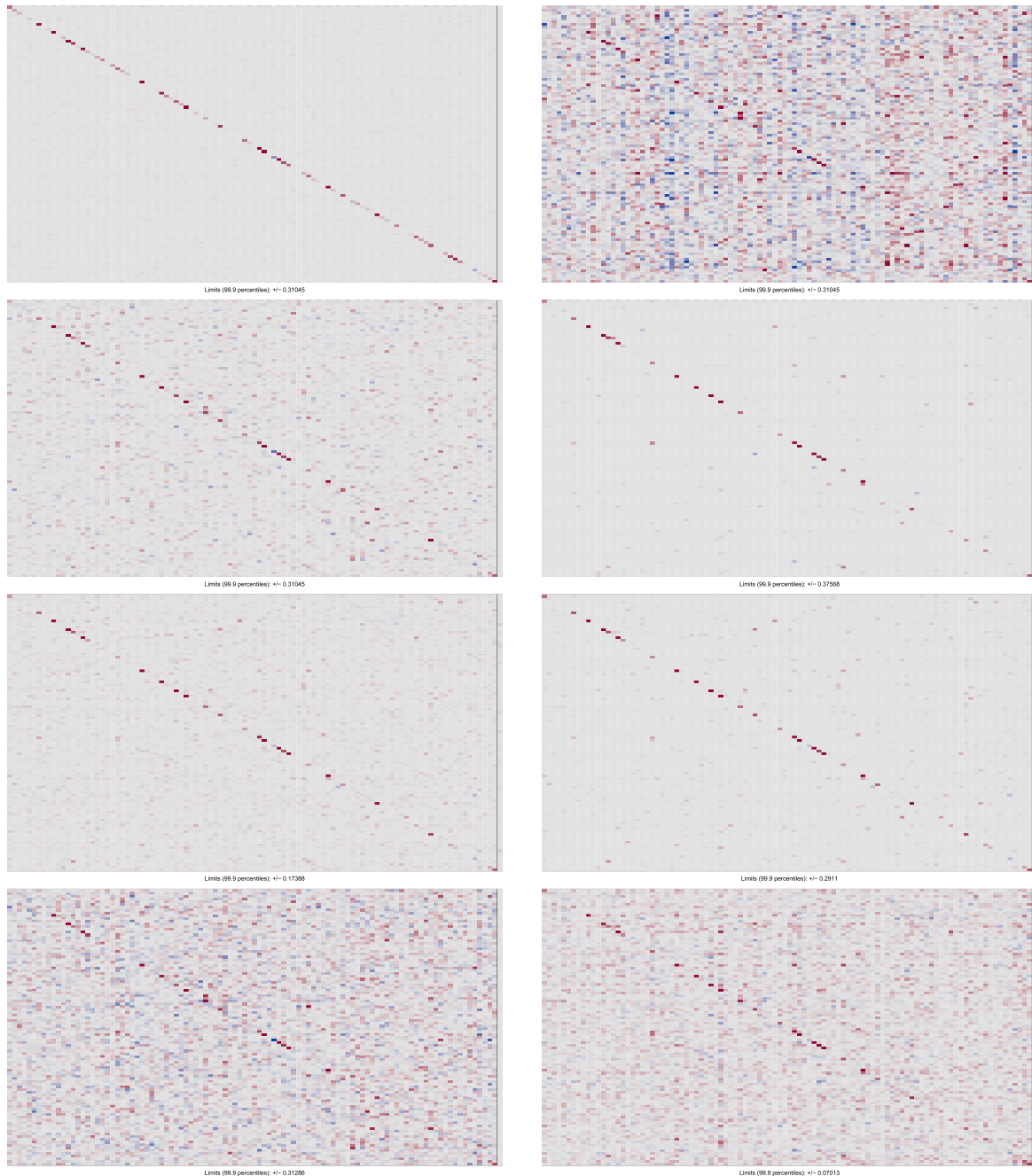


FIGURE A4 Exemplary visualization of the true and estimated VAR coefficients in the *dense* scenario where $T = 250$ and $m = 100$. Top left: DGP. Top right: OLS estimates. Second row: DL prior with $a_{DL} = 1/2$ (left) and $a_{DL} = 1/k = 1/101$ (right). Third row: NG prior with $a_{NG} = 1$ (left) and $a_{NG} = 1/10$ (right). Fourth row: Minnesota prior with $a_M = 1/1,000$ (left) and $a_M = 1/10,000$ (right) [Colour figure can be viewed at wileyonlinelibrary.com]

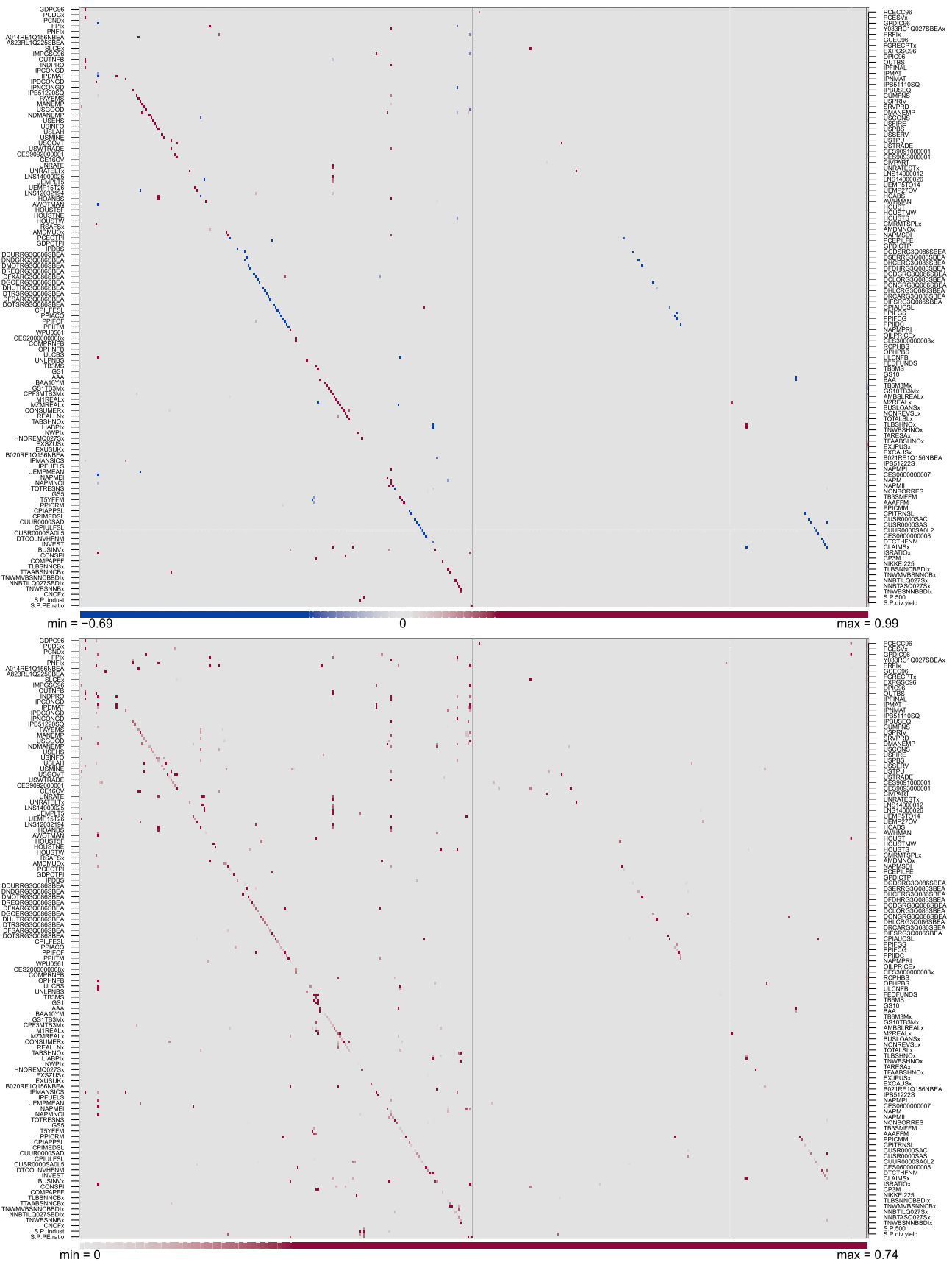


FIGURE A5 Posterior medians (top) and posterior interquartile ranges (bottom) of VAR(2) coefficients, $a = 1/k = 1/431$ [Colour figure can be viewed at wileyonlinelibrary.com]

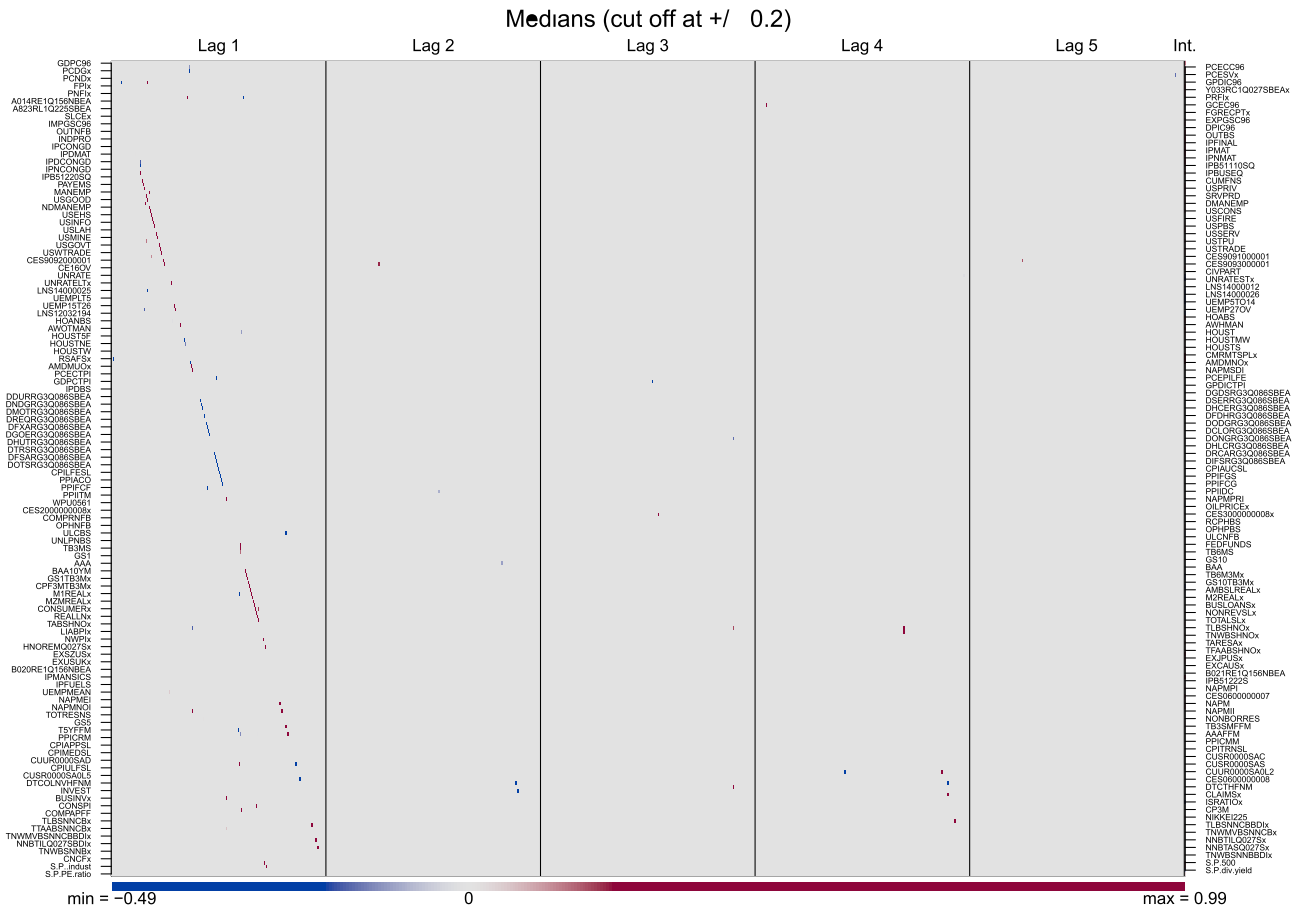


FIGURE A6 Posterior medians of VAR(5) coefficients, $\alpha = 1/k = 1/1076$ [Colour figure can be viewed at wileyonlinelibrary.com]

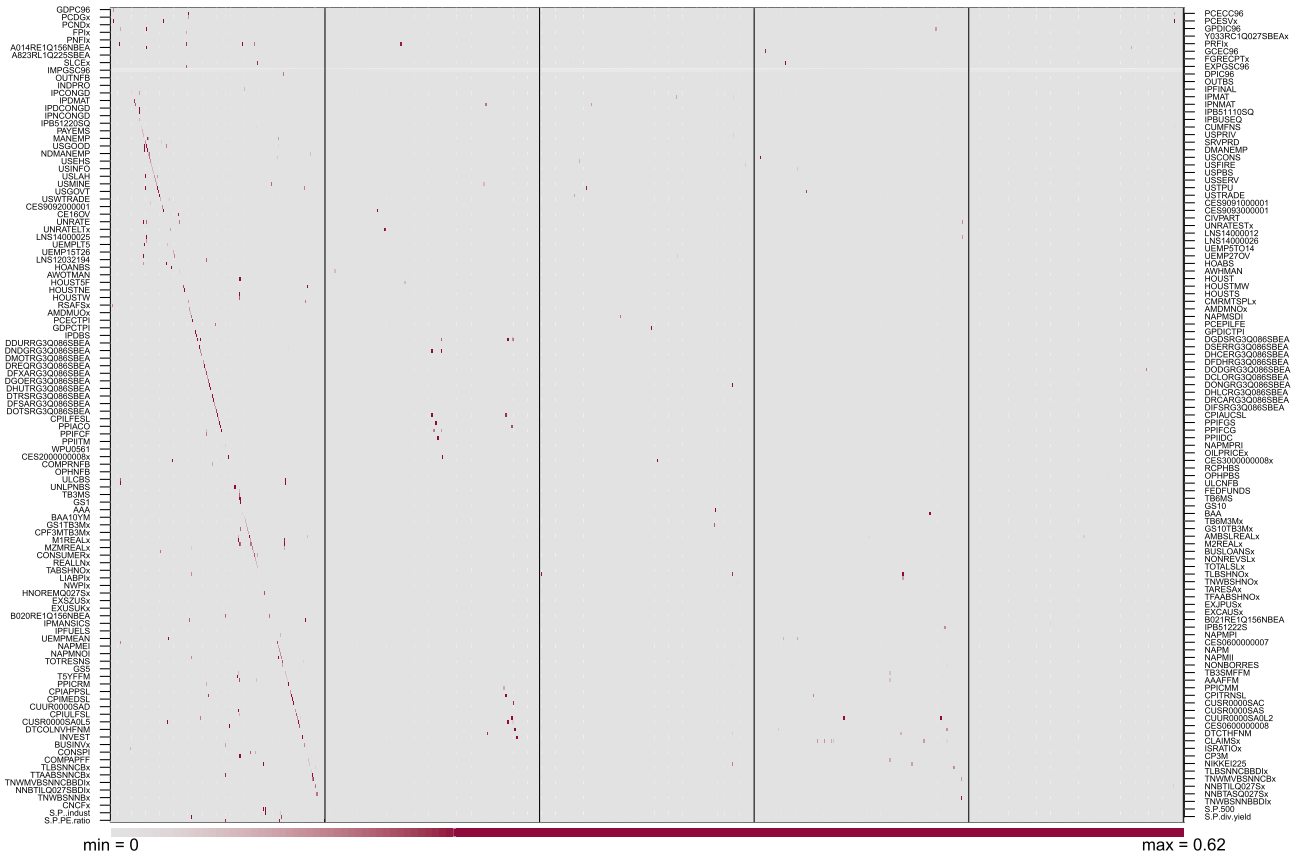


FIGURE A7 Posterior interquartile ranges of VAR(5) coefficients, $a = 1/k = 1/1,076$ [Colour figure can be viewed at wileyonlinelibrary.com]