# Dynamic Shrinkage Processes

Daniel R. Kowal, David S. Matteson, and David Ruppert*

February 27, 2018

## Abstract

We propose a novel class of dynamic shrinkage processes for Bayesian time series and regression analysis. Building upon a global-local framework of prior construction, in which continuous scale mixtures of Gaussian distributions are employed for both desirable shrinkage properties and computational tractability, we model dependence among the local scale parameters. The resulting processes inherit the desirable shrinkage behavior of popular global-local priors, such as the horseshoe prior, but provide additional localized adaptivity, which is important for modeling time series data or regression functions with local features. We construct a computationally efficient Gibbs sampling algorithm based on a Pólya-Gamma scale mixture representation of the proposed process. Using dynamic shrinkage processes, we develop a Bayesian trend filtering model that produces more accurate estimates and tighter posterior credible intervals than competing methods, and apply the model for irregular curve-fitting of minute-by-minute Twitter CPU usage data. In addition, we develop an adaptive time-varying parameter regression model to assess the efficacy of the Fama-French five-factor asset pricing model with momentum added as a sixth factor. Our dynamic analysis of manufacturing and healthcare industry data shows that with the exception of the market risk, no other risk factors are significant except for brief periods.

**KEYWORDS: time series; trend filtering; dynamic linear model; stochastic volatility; asset pricing**

# 1  Introduction

The global-local class of prior distributions is a popular and successful mechanism for providing shrinkage and regularization in a broad variety of models and applications. Global-local priors use continuous scale mixtures of Gaussian distributions to produce desirable shrinkage properties, such as (approximate) sparsity or smoothness, often leading to highly competitive and computationally tractable estimation procedures. For example, in the variable selection context, exact sparsity-inducing priors such as the spike-and-slab prior become intractable for even a moderate number of predictors. By comparison, global-local priors that shrink toward sparsity, such as the horseshoe prior (Carvalho et al., 2010), produce competitive estimators with greater scalability, and are validated by theoretical results, simulation studies, and a variety of applications (Carvalho et al., 2009; Datta and Ghosh, 2013; van der Pas et al., 2014). Unlike non-Bayesian counterparts such as the lasso (Tibshirani, 1996), shrinkage priors also provide adequate uncertainty quantification for parameters of interest (Kyung et al., 2010; van der Pas et al., 2014).

The class of global-local scale mixtures of Gaussian distributions (e.g., Carvalho et al., 2010; Polson and Scott, 2010, 2012a) is defined as follows:

$$[\omega_t|\tau, \lambda_t] \stackrel{indep}{\sim} N(0, \tau^2\lambda_t^2), \quad t = 1, \ldots, T \tag{1}$$

where $\tau > 0$ controls the global shrinkage for all $\{\omega_t\}_{t=1}^T$, while $\lambda_t > 0$ tunes the local shrinkage for a particular $\omega_t$. Such a model is particularly well-suited for sparse data: $\tau$ determines the global level of sparsity for $\{\omega_t\}_{t=1}^T$, while large $\lambda_t$ allows large absolute deviations of $\omega_t$ from its prior mean (zero) and small $\lambda_t$ provides extreme shrinkage to zero. Careful choice of priors for $\lambda_t^2$ and $\tau^2$ provide both the flexibility to accomodate large signals and adequate shrinkage of noise (e.g., Carvalho et al., 2010), so the framework of (1) is widely applicable. The prior in (1) is commonly paired with the likelihood $[y_t|\omega_t, \sigma^2] \stackrel{indep}{\sim} N(\omega_t, \sigma^2)$, but we will consider dynamic generalizations.

Most commonly, the local scale parameters $\{\lambda_t\}$ are assumed to be *a priori* independent and identically distributed (iid). However, it can be advantageous to forgo the independence assumption. In the dynamic setting, in which the observations $y_t$ are time-ordered and $t$ denotes a time index, it is natural to allow the local scale parameter, $\lambda_t$, to depend on the history of the shrinkage process $\{\lambda_s\}_{s<t}$. As a result, the probability of large (or small) deviations of $\omega_t$ from the prior mean (zero), as determined by $\lambda_t$, is informed by the previous shrinkage behavior $\{\lambda_s\}_{s<t}$. Such model-based dependence may improve the ability of the model to adapt dynamically, which is important for time series estimation, forecasting, and inference.

We propose to model dependence in the process $\{\lambda_t\}$ using a novel log-scale representation of a broad class of global-local shrinkage priors. By considering $\{\lambda_t\}$ on the log-scale, we gain access to a variety of widely successful models for dependent data, such as (vector) autoregressions, linear regressions, Gaussian processes, and factor models, among others. An important contribution of the manuscript is to provide (i) a framework for incorporating dependent data models into popular shrinkage priors and (ii) an accompanying Gibbs sampling algorithm, which relies on new parameter expansion techniques for computational efficiency.

We propose to model dependence of the log-variance process $h_t = \log(\tau^2 \lambda_t^2)$ in (1) using the general dependent data model

$$h_t = \mu + \psi_t + \eta_t, \quad \eta_t \overset{iid}{\sim} Z(\alpha, \beta, 0, 1) \tag{2}$$

where $\mu = \log(\tau^2)$ corresponds to the global scale parameter, $(\psi_t + \eta_t) = \log(\lambda_t^2)$ corresponds to the local scale parameter, and $Z(\alpha, \beta, \mu_z, \sigma_z)$ denotes the *Z-distribution* with density function

$$[z] = \left[\sigma B(\alpha, \beta)\right]^{-1} \left\{ \exp\left[(z - \mu_z)/\sigma_z\right] \right\}^{\alpha} \left\{ 1 + \exp\left[(z - \mu_z)/\sigma_z\right] \right\}^{-(\alpha+\beta)}, z \in \mathbb{R} \tag{3}$$

where $B(\cdot, \cdot)$ is the Beta function. In model (2), the local scale parameter $\lambda_t = \exp[(\psi_t + \eta_t)/2]$ has two components: $\psi_t$, which models dependence (see below), and $\eta_t$, which corresponds to the usual iid (log-) local scale parameter. When $\psi_t = 0$, model (2) reduces to the static setting, and implies an *inverted-Beta* prior for $\lambda_t^2$ (see Section 2.1 for more details). Notably, the class of priors represented in (2) includes the important shrinkage distributions in Table 1, in each case extended to the dependent data setting.

| | | |
|---|---|---|
| $\alpha = \beta = 1/2$ | Horseshoe Prior | Carvalho et al. (2010) |
| $\alpha = 1/2, \beta = 1$ | Strawderman-Berger Prior | Strawderman (1971); Berger (1980) |
| $\alpha = 1, \beta = c - 2, c > 0$ | Normal-Exponential-Gamma Prior | Griffin and Brown (2005) |
| $\alpha = \beta \to 0$ | (Improper) Normal-Jeffreys' Prior | Figueiredo (2003); Bae and Mallick (2004) |

Table 1: Special cases of the inverted-Beta prior.

The role of $\psi_t$ in (2) is to provide locally adaptive shrinkage by modeling dependence. For dynamic dependence, we propose the *dynamic shrinkage process*

$$h_{t+1} = \mu + \phi(h_t - \mu) + \eta_t, \quad \eta_t \overset{iid}{\sim} Z(\alpha, \beta, 0, 1) \tag{4}$$

which is equivalent to (2) with $\psi_t = \phi(h_{t-1} - \mu)$. Relative to static shrinkage priors, model (4) only adds one additional parameter, $\phi$, and reduces to the static setting when $\phi = 0$. Importantly, the proposed Gibbs sampler for the parameters in (4) is linear in the number of time points, $T$, and therefore scalable. Other examples of (2) include linear regression, $\psi_t = \boldsymbol{z}_t' \boldsymbol{\alpha}$ for a vector of predictors $\boldsymbol{z}_t$, Gaussian processes, and various multivariate models (see (11) in Section 4). We focus on dynamic dependence, but our modeling framework and computational techniques may be extended to incorporate more general dependence among shrinkage parameters.

Despite the apparent complexity of the model, we develop a new Gibbs sampling algorithm via a parameter expansion of model (2). The MCMC algorithm combines a log-variance sampler (Kastner and Frühwirth-Schnatter, 2014) and a Pólya-Gamma sampler (Polson et al., 2013) to produce a conditionally Gaussian representation of model (2), which permits flex-

ibility in specification of $\psi_t$. The resulting model is easy to implement, computationally efficient, and widely applicable.

For a motivating example, consider the minute-by-minute Twitter CPU usage data in Figure 1a (James et al., 2016). The data show an overall smooth trend interrupted by irregular jumps throughout the morning and early afternoon, with increased volatility from 16:00-18:00. It is important to identify both abrupt changes as well as slowly-varying intraday trends. To model these features, we combine the likelihood $y_t \overset{indep}{\sim} N(\beta_t, \sigma_t^2)$ with a standard SV model for the observation error variance, $\sigma_t^2$, and a *dynamic horseshoe process* as the prior on the second differences of the conditional mean, $\omega_t = \Delta^2 \beta_t = \Delta\beta_t - \Delta\beta_{t-1}$, given by (4) with $\alpha = \beta = 1/2$ (see Section 3.2 for details). The dynamic horseshoe process either drives $\omega_t$ to zero, in which case $\beta_t$ is locally linear, or leaves $\omega_t$ effectively unpenalized, in which case large changes in slope are permissible (see Figure 1b). The resulting posterior expectation of $\beta_t$ and credible bands for the posterior predictive distribution of $\{y_t\}$ adapt to both irregular jumps and smooth trends (see Figure 1a).

For comparison, Figure 1 provides the posterior expectations of both the observation error standard deviations, $\sigma_t$ (Figure 1c) and the prior standard deviations, $[\tau\lambda_t] = \exp(h_t/2)$ (Figure 1d). The horseshoe-like shrinkage behavior of $\lambda_t$ is evident: values of $\lambda_t$ are either near zero, corresponding to aggressive shrinkage of $\omega_t = \Delta^2 \beta_t$ to zero, or large, corresponding to large absolute changes in the slope of $\beta_t$. Importantly, Figure 1 also provides motivation for a *dynamic* shrinkage process: there is clear *volatility clustering* of $\{\lambda_t\}$, in which the shrinkage induced by $\lambda_t$ persists for consecutive time points. The volatility clustering reflects—and motivates—the temporally adaptive shrinkage behavior of the dynamic shrinkage process.

Shrinkage priors and variable selection have been used successfully for time series modeling in a broad variety of settings. Belmonte et al. (2014) propose a Bayesian Lasso prior and Bitto and Frühwirth-Schnatter (2016) use a normal-gamma prior for shrinkage in dynamic linear models, while Korobilis (2013) consider several (non-dynamic) scale mixture priors for time series regression. In each case, the lack of a local (dynamic) scale parameter implies a
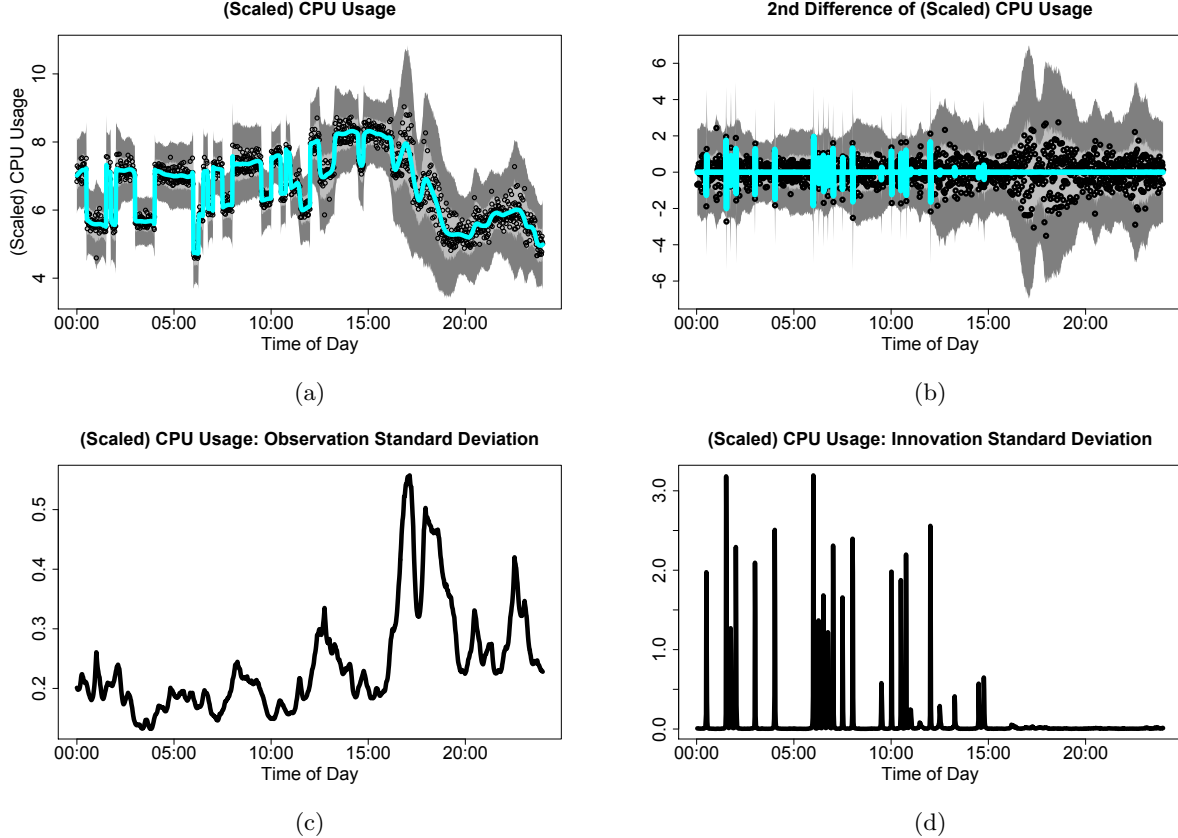
Figure 1: Bayesian trend filtering ($D = 2$) with dynamic horseshoe process innovations of minute-by-minute Twitter CPU usage data. (a) Observed data $y_t$ (points), posterior expectation (cyan) of $\beta_t$, and 95% pointwise highest posterior density (HPD) credible intervals (light gray) and 95% simultaneous credible bands (dark gray) for the posterior predictive distribution of $y_t$. (b) Second difference of observed data $\Delta^2 y_t$ (points), posterior expectation of $\omega_t = \Delta^2 \beta_t$ (cyan), and 95% pointwise HPD intervals (light gray) and simultaneous credible bands (dark gray) for the posterior predictive distribution of $\Delta^2 y_t$. (c) Posterior expectation of time-dependent observation standard deviations, $\sigma_t$. (d) Posterior expectation of time-dependent innovation (prior) standard deviations, $\tau \lambda_t$.

time-invariant rate of shrinkage for each variable. Frühwirth-Schnatter and Wagner (2010) introduce indicator variables to discern between static and dynamic parameters, but the model cannot shrink adaptively for local time periods. Nakajima and West (2013) provide a procedure for local thresholding of dynamic coefficients, but the computational challenges of model implementation are significant. Chan et al. (2012) propose a class of time-varying dimension models, but due to the computational complexity of the model, only consider inclusion or exclusion of a variable for all times, which produces non-dynamic variable selection and a limited set of models. Rockova and McAlinn (2017) develop an optimization approach

for dynamic variable selection. However, their method only provides point estimates, and does so under the restrictive assumptions that (i) the regression coefficients follow identical time series models (AR(1) processes) with known parameters and (ii) the observation error variance is known and non-dynamic. These key limitations are not present in our framework.

Perhaps most comparable to the proposed methodology, Kalli and Griffin (2014) propose a class of priors which exhibit dynamic shrinkage using normal-gamma autoregressive processes. The Kalli and Griffin (2014) prior is a dynamic extension of the normal-gamma prior of Griffin and Brown (2010), and provides improvements in forecasting performance relative to non-dynamic shrinkage priors. However, the Kalli and Griffin (2014) model requires careful specification of several hyperparameters and hyperpriors, and the computation requires sophisticated adaptive MCMC techniques, which results in lengthy computation times. By comparison, our proposed class of dynamic shrinkage processes is far more general, and includes the dynamic horseshoe process as a special case—which notably does not require tuning of sensitive hyperparameters. Empirically, for time-varying parameter regression models with dynamic shrinkage, the Kalli and Griffin (2014) MCMC sampler requires several hours, while our proposed MCMC sampler runs in only a few minutes (see Section 4.1 for details and a comparison of these methods).

We apply dynamic shrinkage processes to develop a dynamic fundamental factor model for asset pricing. We build upon the five-factor Fama-French model (Fama and French, 2015), which extends the three-factor Fama-French model (Fama and French, 1993) for modeling equity returns with common risk factors. We propose a dynamic extension which allows for time-varying factor loadings, possibly with localized or irregular features, and include a sixth factor, momentum (Carhart, 1997). Despite the popularity of the three-factor Fama-French model, there is not yet consensus regarding the necessity of all five factors in Fama and French (2015) or the momentum factor. Dynamic shrinkage processes provide a mechanism for addressing this question: within a time-varying parameter regression model, dynamic shrinkage processes provide the necessary flexibility to adapt to rapidly-changing features,

while shrinking unnecessary factors to zero. Our dynamic analysis shows that with the exception of the market risk factor, no other risk factors are important except for brief periods.

We introduce the dynamic shrinkage process in Section 2 and discuss relevant properties, including the Pólya-Gamma parameter expansion for efficient computations. In Section 3, we apply the prior to develop a more adaptive Bayesian trend filtering model for irregular curve-fitting, and we compare the proposed procedure with state-of-the-art alternatives through simulations and a CPU usage application. We propose in Section 4 a time-varying parameter regression model with dynamic shrinkage processes for adaptive regularization and evaluate the model using simulations and an asset pricing example. In Section 5, we discuss the details of the Gibbs sampling algorithm, and conclude with Section 6. Proofs are in the Appendix, with additional details in the supplement.

## 2 Dynamic Shrinkage Processes

The proposed dynamic shrinkage process contains three prominent features: (i) a dependent model for the local scale parameters, $\lambda_t$; (ii) a log-scale representation of a broad class of global-local priors to propagate desirable shrinkage properties to the dynamic setting; and (iii) a Gaussian scale-mixture representation of the implied log-variance innovations to produce an efficient Gibbs sampling algorithm. In this section, we provide the relevant details regarding these features, and explore the properties of the resulting process.

### 2.1 Log-Scale Representations of Global-Local Priors

In model (2) and (4), we propose to model the (dynamic) dependence of the local scale parameters $\lambda_t$ via the log-variance $h_t = \log(\tau^2 \lambda_t^2)$ of the Gaussian prior (1). As we demonstrate below, our specification of (2) and (4) produces desirable locally adaptive shrinkage properties. However, such shrinkage behavior is not automatic: we must consider

8

appropriate distributions for $\mu$ and $\eta_t$. To illustrate this point, suppose $\eta_t \overset{iid}{\sim} N(0, \sigma_\eta^2)$ in (4), which is a common assumption in stochastic volatility (SV) modeling (Kim et al., 1998). For the likelihood $y_t \sim N(\omega_t, 1)$ and the prior (1), the posterior expectation of $\omega_t$ is $\mathbb{E}[\omega_t | \{y_s\}, \tau] = (1 - \mathbb{E}[\kappa_t | \{y_s\}, \tau]) y_t$, where

$$\kappa_t \equiv \frac{1}{1 + \mathrm{Var}[\omega_t | \tau, \lambda_t]} = \frac{1}{1 + \tau^2 \lambda_t^2} \tag{5}$$

is the *shrinkage parameter*. As noted by Carvalho et al. (2010), $\mathbb{E}[\kappa_t | \{y_s\}, \tau]$ is interpretable as the amount of shrinkage toward zero *a posteriori*: $\kappa_t \approx 0$ yields minimal shrinkage (for signals), while $\kappa_t \approx 1$ yields maximal shrinkage to zero (for noise). For the standard SV model and fixing $\phi = \mu = 0$ for simplicity, $\lambda_t = \exp(\eta_t/2)$ is log-normally distributed, and the shrinkage parameter has density $[\kappa_t] \propto \frac{1}{\kappa_t(1-\kappa_t)} \exp\left\{ -\frac{1}{2\sigma_\eta^2} \left[ \log\left( \frac{1-\kappa_t}{\kappa_t} \right) \right]^2 \right\}$. Notably, the density for $\kappa_t$ approaches zero as $\kappa_t \to 0$ and as $\kappa_t \to 1$. As a result, direct application of the Gaussian SV model may overshrink true signals and undershrink noise.

By comparison, consider the horseshoe prior of Carvalho et al. (2010). The horseshoe prior combines (1) with $[\lambda_t] \overset{iid}{\sim} C^+(0,1)$, where $C^+$ denotes the half-Cauchy distribution. For fixed $\tau = 1$, the half-Cauchy prior on $\lambda_t$ is equivalent to $\kappa_t \overset{iid}{\sim} \mathrm{Beta}(1/2, 1/2)$, which induces a "horseshoe" shape for the shrinkage parameter (see Figure 2). The horseshoe-like behavior is ideal in sparse settings, since the prior density allocates most of its mass near zero (minimal shrinkage of signals) and one (maximal shrinkage of noise). Theoretical results, simulation studies, and a variety of applications confirm the effectiveness of the horseshoe prior (Carvalho et al., 2009, 2010; Datta and Ghosh, 2013; van der Pas et al., 2014).

To emulate the robustness and sparsity properties of the horseshoe and other shrinkage priors in the dynamic setting, we represent a general class of global-local shrinkage priors on the log-scale. As a motivating example, consider the special case of (1) and (4) with $\phi = 0$: $\omega_t \overset{indep}{\sim} N(0, \tau^2 \lambda_t^2)$ with $\log(\lambda_t^2) = \eta_t$. This example is illuminating: we equivalently express the (static) horseshoe prior by letting $\eta_t \overset{D}{=} \log \lambda_t^2$, where $\overset{D}{=}$ denotes equality in distribution.
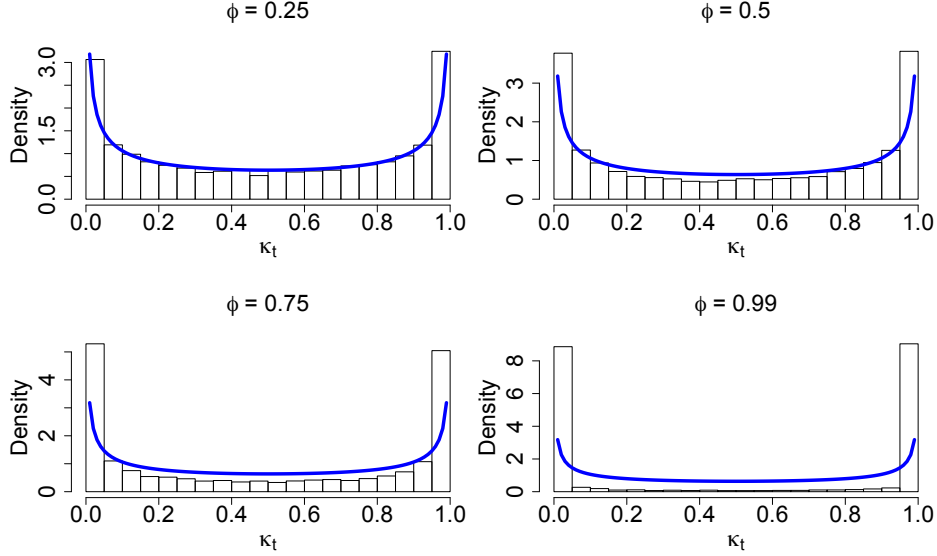
Figure 2: Simulation-based estimate of the stationary distribution of $\kappa_t$ for various AR(1) coefficients $\phi$. The blue line indicates the density of $\kappa_t$ in the static ($\phi = 0$) horseshoe, $[\kappa] \sim \text{Beta}\,(1/2, 1/2)$.

In particular, $\lambda_t \sim C^+(0,1)$ implies

$$\left[\lambda_t^2\right] \propto \left(\lambda_t^2\right)^{-1/2}\left(1 + \lambda_t^2\right)^{-1}$$

which implies

$$[\eta_t] = \pi^{-1}\exp(\eta_t/2)\left[1 + \exp(\eta_t)\right]^{-1}$$

so $\eta_t$ is $Z$-distributed with $\eta_t \sim Z(\alpha = 1/2, \beta = 1/2, \mu_z = 0, \sigma_z = 1)$. Importantly, $Z$-distributions may be written as mean-variance scale mixtures of Gaussian distributions (Barndorff-Nielsen et al., 1982), which produces a useful framework for a parameter-expanded Gibbs sampler.

More generally, consider the inverted-Beta prior, denoted $IB(\beta, \alpha)$, for $\lambda^2$ with density

$$[\lambda^2] \propto \left(\lambda^2\right)^{\alpha-1}\left(1 + \lambda^2\right)^{-(\alpha+\beta)}, \lambda > 0$$

(e.g., Armagan et al., 2011; Polson and Scott, 2012a,b). Special cases of the inverted-Beta

distribution are provided in Table 1.

This broad class of priors may be equivalently constructed via the variances $\lambda_t^2$, the shrinkage parameters $\kappa_t$, or the log-variances $\eta_t$.

**Proposition 1.** *The following distributions are equivalent:*

1. *$\lambda^2 \sim IB(\beta, \alpha)$;*

2. *$\kappa = 1/(1 + \lambda^2) \sim Beta(\beta, \alpha)$;*

3. *$\eta = \log(\lambda^2) = \log(\kappa^{-1} - 1) \sim Z(\alpha, \beta, 0, 1)$.*

Note that the ordering of the parameters $\alpha, \beta$ is identical for the inverted-Beta and Beta distributions, but reversed for the $Z$-distribution.

Now consider the dynamic setting in which $\phi \neq 0$. Model (4) implies that the conditional prior variance for $\omega_t$ in (1) is $\exp(h_t) = \exp(\mu + \phi(h_{t-1} - \mu) + \eta_t) = \tau^2 \lambda_{t-1}^{2\phi} \tilde{\lambda}_t^2$, where $\tau^2 = \exp(\mu)$, $\lambda_{t-1}^2 = \exp(h_{t-1} - \mu)$, and $\tilde{\lambda}_t^2 = \exp(\eta_t) \overset{iid}{\sim} IB(\beta, \alpha)$, as in the non-dynamic setting. This prior generalizes the $IB(\beta, \alpha)$ prior via the local variance term, $\lambda_{t-1}^{2\phi}$, which incorporates information about the shrinkage behavior at the previous time $t-1$ in the prior for $\omega_t$. We formalize the role of this local adjustment term with the following results.

**Proposition 2.** *Suppose $\eta \sim Z(\alpha, \beta, \mu_z, 1)$ for $\mu_z \in \mathbb{R}$. Then $\kappa = 1/(1 + \exp(\eta)) \sim TPB(\beta, \alpha, \exp(\mu_z))$, where $\kappa \sim TPB(\beta, \alpha, \gamma)$ denote the* three-parameter Beta distribution *with density $[\kappa] = [B(\beta, \alpha)]^{-1} \gamma^\beta \kappa^{\beta-1} (1 - \kappa)^{\alpha-1} [1 + (\gamma - 1)\kappa]^{-(\alpha+\beta)}, \kappa \in (0, 1), \gamma > 0$.*

The three-parameter Beta (TPB) distribution (Armagan et al., 2011) generalizes the Beta distribution: $\gamma = 1$ produces the $Beta(\beta, \alpha)$ distribution, while $\gamma > 1$ (respectively, $\gamma < 1$) allocates more mass near zero (respectively, one) relative to the $Beta(\beta, \alpha)$ distribution. For dynamic shrinkage processes, the TPB distribution arises as the conditional prior distribution of $\kappa_{t+1}$ given $\{\kappa_s\}_{s \leq t}$.

**Theorem 1.** *For the dynamic shrinkage process (4), the conditional prior distribution of*

11

*the shrinkage parameter* $\kappa_{t+1} = 1/(1 + \tau^2 \lambda_{t+1}^2)$ *is*

$$[\kappa_{t+1} | \{\kappa_s\}_{s \leq t}, \phi, \tau] \sim TPB\left(\beta, \alpha, \tau^{2(1-\phi)} \left[\frac{1 - \kappa_t}{\kappa_t}\right]^{\phi}\right) \tag{6}$$

*or equivalently,* $[\kappa_{t+1} | \{\lambda_s\}_{s \leq t}, \phi, \tau] \sim TPB(\beta, \alpha, \tau^2 \lambda_t^{2\phi})$.

The proof of Theorem 1 is in the Appendix. Naturally, the previous value of the shrinkage parameter, $\kappa_t$, together with the AR(1) coefficient $\phi$, inform both the magnitude and the direction of the distributional shift of $\kappa_{t+1}$.

**Theorem 2.** *For the dynamic horseshoe process of* (4) *with* $\alpha = \beta = 1/2$ *and fixed* $\tau = 1$, *the conditional prior distribution* (6) *satisfies* $\mathbb{P}(\kappa_{t+1} < \varepsilon | \{\kappa_s\}_{s \leq t}, \phi) \to 1$ *as* $\kappa_t \to 0$ *for any* $\varepsilon \in (0, 1)$ *and fixed* $\phi \neq 0$.

The proof of Theorem 2 is in the Appendix. Importantly, Theorem 2 demonstrates that the mass of the conditional prior distribution for $\kappa_{t+1}$ concentrates near zero—corresponding to minimal shrinkage of signals—when $\kappa_t$ is near zero, so the shrinkage behavior at time $t$ informs the (prior) shrinkage behavior at time $t + 1$.

We similarly characterize the posterior distribution of $\kappa_{t+1}$ given $\{\kappa_s\}_{s \leq t}$ in the following theorem, which extends the results of Datta and Ghosh (2013) to the dynamic setting.

**Theorem 3.** *Under the likelihood* $y_t \overset{indep}{\sim} N(\omega_t, 1)$, *the prior* (1), *and the dynamic horseshoe process* (4) *with* $\alpha = \beta = 1/2$ *and fixed* $\phi \neq 0$, *the posterior distribution of* $\kappa_{t+1}$ *given the history of the shrinkage process* $\{\kappa_s\}_{s \leq t}$ *satisfies the following properties:*

(a) *For any fixed* $\varepsilon \in (0, 1)$, $\mathbb{P}(\kappa_{t+1} > 1 - \varepsilon | y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau) \to 1$ *as* $\gamma_t \to 0$ *uniformly in* $y_{t+1} \in \mathbb{R}$, *where* $\gamma_t = \tau^{2(1-\phi)} [(1 - \kappa_t)/\kappa_t]^{\phi}$.

(b) *For any fixed* $\varepsilon \in (0, 1)$ *and* $\gamma_t < 1$, $\mathbb{P}(\kappa_{t+1} < \varepsilon | y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau) \to 1$ *as* $|y_{t+1}| \to \infty$.

The proof of Theorem 3 is in the supplementary material, and uses the observation that

marginally, $[y_{t+1}|\{\kappa_s\}] \overset{indep}{\sim} N(0, \kappa_{t+1}^{-1})$, so the posterior distribution of $\kappa_{t+1}$ is

$$[\kappa_{t+1}|y_{t+1}, \{\kappa_s\}_{s\leq t}, \phi, \tau] \propto \left\{\kappa_{t+1}^{\beta-1}(1-\kappa_{t+1})^{\alpha-1}\left[1+(\gamma_t-1)\kappa_{t+1}\right]^{-(\alpha+\beta)}\right\} \left\{\kappa_{t+1}^{1/2}\exp\left(-y_{t+1}^2\kappa_{t+1}/2\right)\right\}$$

$$\propto (1-\kappa_{t+1})^{-1/2}\left[1+(\gamma_t-1)\kappa_{t+1}\right]^{-1}\exp\left(-y_{t+1}^2\kappa_{t+1}/2\right).$$

Theorem 3(a) demonstrates that the posterior mass of $[\kappa_{t+1}|\{\kappa_s\}_{s\leq t}]$ concentrates near one as $\tau \to 0$, as in the non-dynamic horseshoe, but also as $\kappa_t \to 1$. Therefore, the dynamic horseshoe process provides an additional mechanism for shrinkage of noise, besides the global scale parameter $\tau$, via the previous shrinkage parameter $\kappa_t$. Moreover, Theorem 3(b) shows that, despite the additional shrinkage capabilities, the posterior mass of $[\kappa_{t+1}|\{\kappa_s\}_{s\leq t}]$ concentrates near zero for large absolute signals $|y_{t+1}|$, which indicates robustness of the dynamic horseshoe process to large signals analogous to the static horseshoe prior.

When $|\phi| < 1$, the log-variance process $\{h_t\}$ is stationary, which implies $\{\kappa_t\}$ is stationary. In Figure 2, we plot a simulation-based estimate of the stationary distribution of $\kappa_t$ for various values of $\phi$ under the dynamic horseshoe process. The stationary distribution of $\kappa_t$ is similar to the static horseshoe distribution ($\phi = 0$) for $\phi < 0.5$, while for large values of $\phi$ the distribution becomes more peaked at zero (less shrinkage of $\omega_t$) and one (more shrinkage of $\omega_t$). The result is intuitive: larger $|\phi|$ corresponds to greater persistence in shrinkage behavior, so marginally we expect states of aggressive shrinkage or little shrinkage.

## 2.2 Scale Mixtures via Pólya-Gamma Processes

For efficient computations, we develop a parameter expansion of the general model (2) and the dynamic shrinkage process (4) using a conditionally Gaussian representation for $\eta_t$. In doing so, we may incorporate Gaussian models—and accompanying sampling algorithms—for dependent data in (2). Given a conditionally Gaussian parameter expansion, a Gibbs sampler for (2) proceeds as follows: (i) draw the log-variances $h_t$, for which the conditional prior (2) is Gaussian, and (ii) draw the parameters in $\mu$ and $\psi_t$, for which the conditional

likelihood (2) is Gaussian. For the log-variance sampler, we represent the likelihood for $h_t$ in (1) on the log-scale and approximate the resulting distribution using a known discrete mixture of Gaussian distributions (see Section 5). This approach is popular in SV modeling (e.g., Kim et al., 1998), which is analogous to the dynamic shrinkage process in (4). Importantly, the proposed parameter expansion inherits the computational complexity of the samplers for $h_t$ and $\psi_t$: for the dynamic shrinkage processes in (4), the proposed parameter expansion implies that the log-variance $\{h_t\}_{t=1}^T$ is a Gaussian dynamic linear model, and therefore $\{h_t\}_{t=1}^T$ may be sampled jointly in $\mathcal{O}(T)$ computations. We provide the relevant details in Section 5.

The proposed algorithm requires a parameter expansion of $\eta_t \overset{iid}{\sim} Z(\alpha, \beta, 0, 1)$ in (4) as a scale mixture of Gaussian distributions. The representation of a $Z$-distribution as a mean-variance scale mixtures of Gaussian distributions is due to Barndorff-Nielsen et al. (1982). For implementation, we build on the framework of Polson et al. (2013), who propose a Pólya-Gamma scale mixture of Gaussians representation for Bayesian logistic regression. A *Pólya-Gamma* random variable $\xi$ with parameters $b > 0$ and $c \in \mathbb{R}$, denoted $\xi \sim \mathrm{PG}(b, c)$, is an infinite convolution of Gamma random variables:

$$\xi \overset{D}{=} \frac{1}{2\pi^2} \sum_{k=1}^{\infty} \frac{g_k}{(k - 1/2)^2 - c^2/(4\pi^2)} \tag{7}$$

where $g_k \overset{iid}{\sim} \mathrm{Gamma}(b, 1)$. Properties of Pólya-Gamma random variables may be found in Barndorff-Nielsen et al. (1982) and Polson et al. (2013). Our interest in Pólya-Gamma random variables derives from their role in representing the $Z$-distribution as a mean-variance scale mixture of Gaussians.

**Theorem 4.** *The random variable $\eta \sim Z(\alpha, \beta, 0, 1)$, or equivalently $\eta = \log(\lambda^2)$ with $\lambda^2 \sim$*

$IB(\beta, \alpha)$, is a mean-variance scale mixture of Gaussian distributions with

$$\begin{cases} [\eta|\xi] \sim N\left(\xi^{-1}[\alpha - \beta]/2, \xi^{-1}\right) \\ [\xi] \sim PG(\alpha + \beta, 0). \end{cases} \tag{8}$$

Moreover, the conditional distribution of $\xi$ is $[\xi|\eta] \sim PG(\alpha + \beta, \eta)$.

The proof of Theorem 4 is in the Appendix. When $\alpha = \beta$, the $Z$-distribution is symmetric, and the conditional expectation in (8) simplifies to $\mathbb{E}[\eta|\xi] = 0$. Polson et al. (2013) propose a sampling algorithm for Pólya-Gamma random variables, which is available in the R package BayesLogit, and is extremely efficient when $b = 1$. In our setting, this corresponds to $\alpha + \beta = 1$, for which the horseshoe prior is the prime example. Importantly, this representation allows us to construct an efficient sampling algorithm that combines an $\mathcal{O}(T)$ sampling algorithm for the log-volatilities $\{h_t\}_{t=1}^T$ with a Pólya-Gamma sampler for the mixing parameters.

# 3    Bayesian Trend Filtering with Dynamic Shrinkage Processes

Dynamic shrinkage processes are particularly appropriate for dynamic linear models (DLMs). DLMs combine an observation equation, which relates the observed data to latent state variables, and an evolution equation, which allows the state variables—and therefore the conditional mean of the data—to be dynamic. By construction, DLMs contain many parameters, and therefore may benefit from structured regularization. The proposed dynamic shrinkage processes offer such regularization, and unlike existing methods, do so adaptively.

Consider the following DLM with a $D$th order random walk on the state variable, $\beta_t$:

$$\begin{cases} y_t = \beta_t + \epsilon_t, \quad [\epsilon_t|\sigma_\epsilon] \overset{iid}{\sim} N(0, \sigma_\epsilon^2), \quad t = 1, \ldots, T \\ \Delta^D \beta_{t+1} = \omega_t, \quad [\omega_t|\tau, \{\lambda_s\}] \overset{indep}{\sim} N(0, \tau^2\lambda_t^2), \quad t = D, \ldots, T \end{cases} \tag{9}$$

and $\beta_{t+1} = \omega_t \sim N(0, \tau^2\lambda_t^2)$ for $t = 0, \ldots, D-1$, where $\Delta$ is the differencing operator

and $D \in \mathbb{Z}^+$ is the degree of differencing. By imposing a shrinkage prior on $\lambda_t$, model (9) may be viewed as a Bayesian adaptation of the *trend filtering* model of Kim et al. (2009) and Tibshirani (2014): model (9) features a penalty encouraging sparsity of the $D$th order differences of the conditional mean, $\beta_t$. Faulkner and Minin (2016) provide an implementation based on the (static) horseshoe prior and the Bayesian lasso, and further allow for non-Gaussian likelihoods. We refer to model (9) as a *Bayesian trend filtering* (BTF) model, with various choices available for the distribution of the innovation standard deviations, $(\tau\lambda_t)$.

We propose a dynamic horseshoe process as the prior for the innovations $\omega_t$ in model (9). The aggressive shrinkage of the horseshoe prior forces small values of $|\omega_t| = |\Delta^D \beta_{t+1}|$ toward zero, while the robustness of the horseshoe prior permits large values of $|\Delta^D \beta_{t+1}|$. When $D = 2$, model (9) will shrink the conditional mean $\beta_t$ toward a piecewise linear function with breakpoints determined adaptively, while allowing large absolute changes in the slopes. Further, using the *dynamic* horseshoe process, the shrinkage effects induced by $\lambda_t$ are time-dependent, which provides localized adaptability to regions with rapidly- or slowly-changing features. Following Carvalho et al. (2010) and Polson and Scott (2012b), we assume a half-Cauchy prior for the global scale parameter $\tau \sim C^+(0, \sigma_\epsilon/\sqrt{T})$, in which we scale by the observation error variance and the sample size (Piironen and Vehtari, 2016). Using Pólya-Gamma mixtures, the implied conditional prior on $\mu = \log(\tau^2)$ is $[\mu|\sigma_\epsilon, \xi_\mu] \sim N(\log \sigma_\epsilon^2 - \log T, \xi_\mu^{-1})$ with $\xi_\mu \sim \text{PG}(1, 0)$. We include the details of the Gibbs sampling algorithm for model (9) in Section 5, which is notably *linear* in the number of time points, $T$: the full conditional posterior precision matrices for $\boldsymbol{\beta} = (\beta_1, \ldots, \beta_T)'$ and $\boldsymbol{h} = (h_1, \ldots, h_T)'$ are $D$-banded and tridiagonal, respectively, which admit highly efficient $\mathcal{O}(T)$ back-band substitution sampling algorithms (see the supplement for empirical evidence).

## 3.1 Bayesian Trend Filtering: Simulations

To assess the performance of the Bayesian trend filtering (BTF) model (9) with dynamic horseshoe innovations (**BTF-DHS**), we compared the proposed methods to several competitive alternatives using simulated data. We considered the following variations on BTF model (9): normal-inverse-Gamma (**BTF-NIG**) innovations via $\tau^{-2} \sim \text{Gamma}(0.001, 0.001)$ with $\lambda_t = 1$; and (static) horseshoe priors for the innovations (**BTF-HS**) via $\tau, \lambda_t \stackrel{iid}{\sim} C^+(0, 1)$. In addition, we include the (non-Bayesian) trend filtering model of Tibshirani (2014) implemented using the R package `genlasso` (Arnold and Tibshirani, 2014), for which the regularization tuning parameter is chosen using cross-validation (**Trend Filtering**). For all trend filtering models, we select $D = 2$, but the relative performance is similar for $D = 1$. Among non-trend filtering models, we include a smoothing spline estimator implemented via `smooth.spline()` in R (**Smoothing Spline**); the wavelet-based estimator of Abramovich et al. (1998) (**BayesThresh**) implemented in the `wavethresh` package (Nason, 2016); and the nested Gaussian Process (**nGP**) model of Zhu and Dunson (2013), which relies on a state space model framework for efficient computations, comparable to—but empirically less efficient than—the BTF model (9).

We simulated 100 data sets from the model $y_t = y_t^* + \epsilon_t$, where $y_t^*$ is the true function and $\epsilon_t \stackrel{indep}{\sim} N(0, \sigma_*^2)$. We use the following true functions $y_t^*$ from Donoho and Johnstone (1994): *Doppler*, *Bumps*, *Blocks*, and *Heavisine*, implemented in the R package `wmtsa` (Constantine and Percival, 2016). The noise variance $\sigma_*^2$ is determined by selecting a root-signal-to-noise ratio (RSNR) and computing $\sigma_* = \sqrt{\frac{\sum_{t=1}^{T}(y_t^* - \bar{y}^*)^2}{T-1}} \Big/ \text{RSNR}$, where $\bar{y}^* = \frac{1}{T}\sum_{t=1}^{T} y_t^*$. As in Zhu and Dunson (2013), we select RSNR $= 7$ and use a moderate length time series, $T = 128$.

In Figure 3, we provide an example of each true curve $y_t^*$, together with the proposed BTF-DHS posterior expectations and credible bands. Notably, the Bayesian trend filtering model (9) with $D = 2$ and dynamic horseshoe innovations provides an exceptionally accurate fit to each data set. Importantly, the posterior expectations and the posterior credible bands adapt
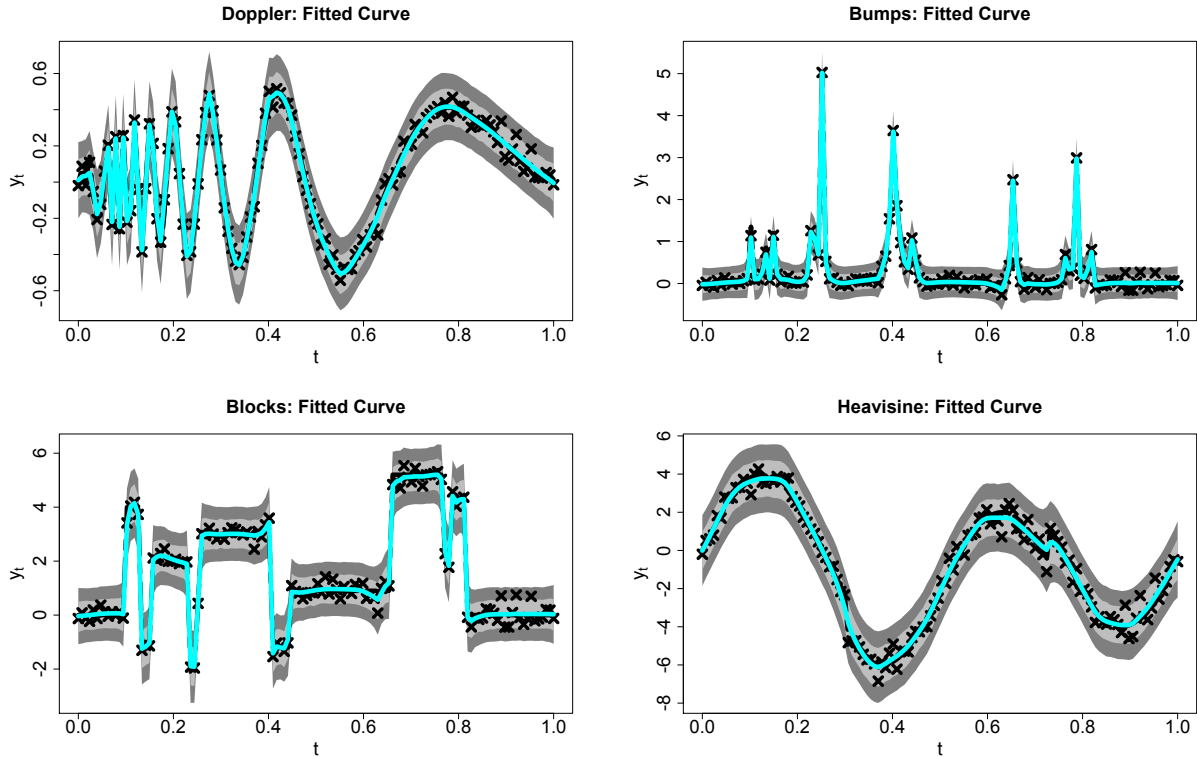
Figure 3: Fitted curves for simulated data with $T = 128$ and RSNR $= 7$. Each panel includes the simulated observations (x-marks), the posterior expectations of $\beta_t$ (cyan), and the 95% pointwise HPD credible intervals (light gray) and 95% simultaneous credible bands (dark gray) for the posterior predictive distribution of $\{y_t\}$ under BTF-DHS model (9) with $D = 2$. The proposed estimator, as well as the uncertainty bands, accurately capture both slowly- and rapidly-changing behavior in the underlying functions.

to both slowly- and rapidly-changing behavior in the underlying curves. The implementation is also efficient: the computation time for 10,000 iterations of the Gibbs sampling algorithm, implemented in R (on a MacBook Pro, 2.7 GHz Intel Core i5), is about 45 seconds.

To compare the aforementioned procedures, we compute the root mean squared errors $\text{RMSE}(\hat{y}) = \sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^* - \hat{y}_t)^2}$ for all estimators $\hat{y}$ of the true function, $y^*$. The results are displayed in Figure 4. The proposed BTF-DHS implementation substantially outperforms all competitors, especially for rapidly-changing curves (Doppler and Bumps). The exceptional performance of BTF-DHS is paired with comparably small variability of RMSE, especially relative to non-dynamic horseshoe model (BTF-HS). Interestingly, the magnitude and variability of the RMSEs for BTF-DHS are related to the AR(1) coefficient, $\phi$: the 95%

HPD intervals (corresponding to Figure 3) are $(0.77, 0.97)$ (Doppler), $(0.81, 0.97)$ (Bumps), $(0.76, 0.96)$ (Blocks), and $(-0.04, 0.74)$ (Heavisine). For the smoothest function, Heavisine, there is less separation among the estimators. Nonetheless, BTF-DHS performs the best, even though the HPD interval for $\phi$ is wider and contains zero.

We are also interested in uncertainty quantification, and in particular how the dynamic horseshoe prior compares to the horseshoe prior. We compute the mean credible intervals widths $\text{MCIW} = \frac{1}{T} \sum_{t=1}^{T} (\hat{\beta}_t^{(97.5)} - \hat{\beta}_t^{(2.5)})$ where $\hat{\beta}_t^{(97.5)}$ and $\hat{\beta}_t^{(2.5)}$ are the 97.5% and 2.5% quantiles, respectively, of the posterior distribution of $\beta_t$ in (9) for the BTF-DHS and BTF-HS. The results are in Figure 5. The dynamic horseshoe provides massive reductions in MCIW, again in all cases except for Heavisine, for which the methods perform similarly. Therefore, in addition to more accurate point estimation (Figure 4), the BTF-DHS model produces significantly tighter credible intervals—while maintaining the approximately correct nominal (frequentist) coverage.

## 3.2 Bayesian Trend Filtering: Application to CPU Usage Data

To demonstrate the adaptability of the dynamic horseshoe process for model (9), we consider the CPU usage data in Figure 1a. The data exhibit substantial complexity: an overall smooth intraday trend but with multiple irregularly-spaced jumps, and an increase in volatility from 16:00-18:00. Our goal is to provide an accurate measure of the trend, including jumps, with appropriate uncertainty quantification. For this purpose, we employ the BTF-DHS model (9), which we extend to include stochastic volatility for the observation error: $y_t \overset{indep}{\sim} N(\beta_t, \sigma_t^2)$ with an AR(1) model on $\log \sigma_t^2$ as in (4) with $\eta_t \overset{iid}{\sim} N(0, \sigma_\eta^2)$. For the additional sampling step of the stochastic volatility parameters, we use the algorithm of Kastner and Frühwirth-Schnatter (2014) implemented in the `R` package `stochvol` (Kastner, 2016).

The resulting model fit is summarized in Figure 1. The posterior expectation and posterior credible bands accurately model both irregular jumps and smooth trends, and capture the increase in volatility from 16:00-18:00 (see Figure 1c). By examining regions of nonover-
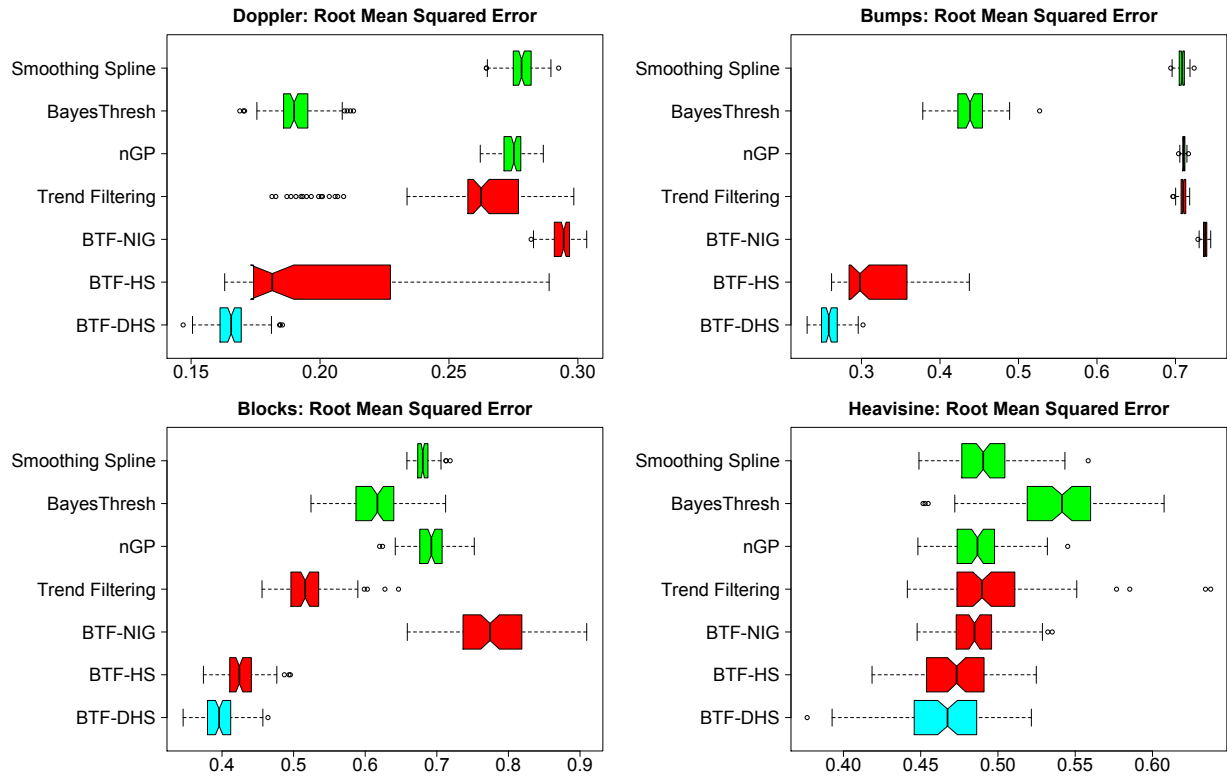
Figure 4: Root mean squared errors for simulated data with $T = 128$ and RSNR $= 7$. Non-overlapping notches indicate significant differences between medians. The Bayesian trend filtering (BTF) estimators differ in their innovation distributions, which determines the shrinkage behavior of the second order differences ($D = 2$): normal-inverse-Gamma (NIG), horseshoe (HS), and dynamic horseshoe (DHS).

lapping simultaneous posterior credible bands, we may assess change points in the level of the data. In particular, the model fit suggests that the CPU usage followed a slowly increasing trend interrupted by jumps of two distinct magnitudes prior to 16:00, after which the volatility increased and the level decreased until approximately 18:00.

We augment the simulation study of Section 3.1 with a comparison of out-of-sample estimation and inference of the CPU usage data. We fit each model using 90% ($T = 1296$) of the data selected randomly for training and the remaining 10% ($T = 144$) for testing, which was repeated 100 times. Models were compared using RMSE and MCIW.

Unlike the simulation study in Section 3.1, the subsampled data are *not* equally spaced. Taking advantage of the computational efficiency of the proposed BTF methodology, we
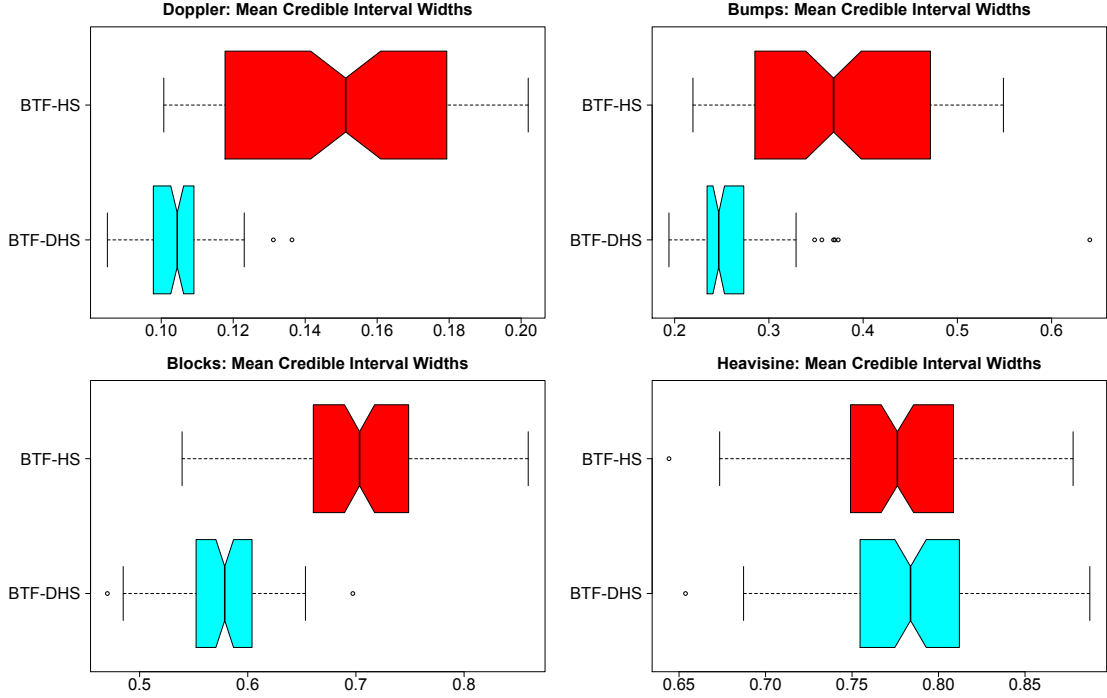
Figure 5: Mean credible interval widths for simulated data with $T = 128$ and RSNR = 7. Non-overlapping notches indicate significant differences between medians. The Bayesian trend filtering (BTF) estimators differ in their innovation distributions, which determines the shrinkage behavior of the second order differences ($D = 2$): normal-inverse-Gamma (NIG), horseshoe (HS), and dynamic horseshoe (DHS).

employ a model-based imputation scheme similar to Elerian et al. (2001), which is valid for missing observations. For unequally-spaced data $y_{t_i}, i = 1, \ldots, T$, we expand the operative data set to include missing observations along an equally-spaced grid, $t^* = 1, \ldots, T^*$, such that for each observation point $i$, $y_{t_i} = y_{t^*}$ for some $t^*$. Although $T^* \geq T$, possibly with $T^* \gg T$, all computations within the sampling algorithm, including the imputation sampling scheme for $\{y_{t^*} : t^* \neq t_i\}$, are linear in the number of (equally-spaced) time points, $T^*$. Therefore, we may apply the same Gibbs sampling algorithm as before, with the additional step of drawing $y_{t^*} \overset{indep}{\sim} N(\beta_{t^*}, \sigma^2_{t^*})$ for each unobserved $t^* \neq t_i$. Implicitly, this procedure assumes that the unobserved points are missing at random, which is satisfied by the aforementioned subsampling scheme.

The results of the out-of-sample estimation study are displayed in Figure 6. The BTF procedures are notably superior to the non-Bayesian trend filtering and smoothing spline

estimators, and, as with the simulations of Section 3.1, the proposed BTF-DHS model sub-
stantially outperforms all competitors. Importantly, the significant reduction in MCIW by
BTF-DHS indicates that the posterior credible intervals for the out-of-sample points $y_{t^*}$ are
substantially tighter for our method. By reducing uncertainty—while maintaining the ap-
proximately correct nominal (frequentist) coverage—the proposed BTF-DHS model provides
greater power to detect local features. In addition, the MCMC for the BTF-DHS is fast,
despite the imputation procedure: 10,000 iterations runs in about 80 seconds (in R on a
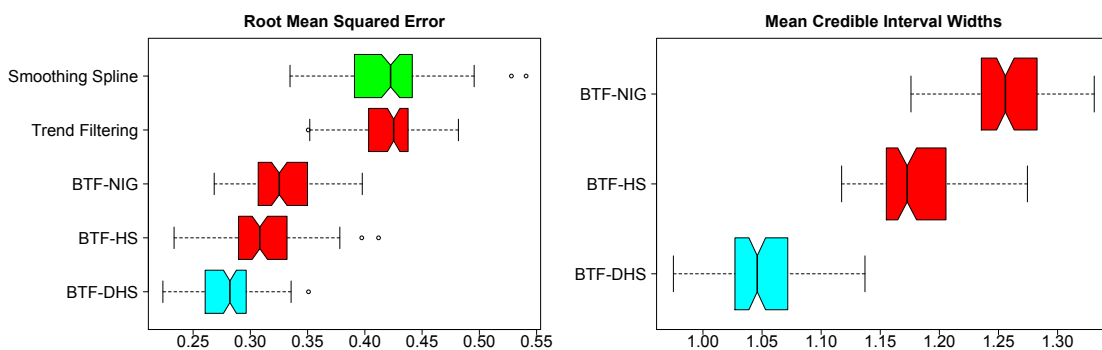MacBook Pro, 2.7 GHz Intel Core i5).



Figure 6: Root mean squared error (**left**) and mean credible interval widths (**right**) for out-of-sample minute-by-
minute CPU usage data. Non-overlapping notches indicate significant differences between medians. The Bayesian
trend filtering (BTF) estimators differ in their innovation distributions, which determines the shrinkage behavior of
the second order differences ($D = 2$): normal-inverse-Gamma (NIG), horseshoe (HS), and dynamic horseshoe (DHS).

# 4   Joint Shrinkage for Time-Varying Parameter Models

Dynamic shrinkage processes are appropriate for multivariate time series and functional data
models that may benefit from locally adaptive shrinkage properties. As outlined in Dangl and
Halling (2012), models with *time-varying parameters* are particularly important in financial
and economic applications, due to the inherent structural changes in regulations, monetary
policy, market sentiments, and macroeconomic interrelations that occur over time. Consider
the following time-varying parameter regression model with multiple dynamic predictors

$\boldsymbol{x}_t = (x_{1,t}, \ldots, x_{p,t})'$:

$$\begin{cases} y_t = \boldsymbol{x}_t' \boldsymbol{\beta}_t + \epsilon_t, & [\epsilon_t | \sigma_\epsilon] \overset{indep}{\sim} N(0, \sigma_\epsilon^2) \\ \Delta^D \boldsymbol{\beta}_{t+1} = \boldsymbol{\omega}_t, & [\omega_{j,t} | \tau_0, \{\tau_k\}, \{\lambda_{k,s}\}] \overset{indep}{\sim} N(0, \tau_0^2 \tau_j^2 \lambda_{j,t}^2) \end{cases} \tag{10}$$

where $\boldsymbol{\beta}_t = (\beta_{1,t}, \ldots, \beta_{p,t})'$ is the vector of dynamic regression coefficients and $D \in \mathbb{Z}^+$ is the degree of differencing. Model (10) is also a (discretized) *concurrent functional linear model* (e.g., Ramsay and Silverman, 2005) and a *varying-coefficient model* (Hastie and Tibshirani, 1993) in the index $t$, and therefore is broadly applicable. The prior for the innovations $\omega_{j,t}$ incorporates three levels of global-local shrinkage: a global shrinkage parameter $\tau_0$, a predictor-specific shrinkage parameter $\tau_j$, and a predictor- and time-specific local shrinkage parameter $\lambda_{j,t}$. Relative to existing time-varying parameter regression models, our approach incorporates an additional layer of dynamic dependence: not only are the parameters time-varying, but also the *relative influence* of the parameters is time-varying via the shrinkage parameters—which are dynamically dependent themselves.

We also considered a VARIMA alternative to (10): $\Delta^D \boldsymbol{\beta}_{t+1} = \boldsymbol{\Gamma} \Delta^D \boldsymbol{\beta}_t + \boldsymbol{\omega}_t$, where $\boldsymbol{\Gamma}$ is a $p \times p$ VAR coefficient matrix. While the VARIMA model allows for lagged cross-correlations between components of $\Delta^D \boldsymbol{\beta}_t$, it does not produce smooth paths for $\beta_{j,t}$, so we do not pursue it further.

To provide jointly localized shrinkage of the dynamic regression coefficients $\{\beta_{j,t}\}$ analogous to the Bayesian trend filtering model of Section 3, we extend (4) to allow for multivariate time dependence via a vector autoregression (VAR) on the log-variance:

$$\begin{cases} [\omega_{j,t} | \tau_0, \{\tau_k\}, \{\lambda_{k,s}\}] \overset{indep}{\sim} N(0, \tau_0^2 \tau_j^2 \lambda_{j,t}^2) \\ h_{j,t} = \log(\tau_0^2 \tau_j^2 \lambda_{j,t}^2), & j = 1, \ldots, p, t = 1, \ldots, T \\ \boldsymbol{h}_{t+1} = \boldsymbol{\mu} + \boldsymbol{\Phi}(\boldsymbol{h}_t - \boldsymbol{\mu}) + \boldsymbol{\eta}_t, & \eta_{j,t} \overset{iid}{\sim} Z(\alpha, \beta, 0, 1) \end{cases} \tag{11}$$

where $\boldsymbol{h}_t = (h_{1,t}, \ldots, h_{p,t})'$, $\boldsymbol{\mu} = (\mu_1, \ldots, \mu_p)'$, $\boldsymbol{\eta}_t = (\eta_{1,t}, \ldots, \eta_{p,t})'$, and $\boldsymbol{\Phi}$ is the $p \times p$ VAR

coefficient matrix. We assume $\mathbf{\Phi} = \text{diag}(\phi_1, \ldots, \phi_p)$ for simplicity, but non-diagonal extensions are available. Contemporaneous dependence may be introduced via a copula model for the log-variance innovations, $\boldsymbol{\eta}_t$ (Joe, 2015), but may reduce computational and MCMC efficiency. As in the univariate setting, we use Pólya-Gamma mixtures (independently) for the log-variance evolution errors, $[\eta_{j,t}|\xi_{j,t}] \overset{indep}{\sim} N(\xi_{j,t}^{-1}[\alpha - \beta]/2, \xi_{j,t}^{-1})$ with $\xi_{j,t} \overset{iid}{\sim} \text{PG}(\alpha + \beta, 0)$ and $\alpha = \beta = 1/2$. We augment model (11) with half-Cauchy priors for the predictor-specific and global parameters, $\tau_j \overset{indep}{\sim} C^+(0,1)$ and $\tau_0 \sim C^+(0, \sigma_\epsilon/\sqrt{Tp})$, in which we scale by the observation error variance and the number of innovations $\{\omega_{j,t}\}$ (Piironen and Vehtari, 2016). These priors may be equivalently represented on the log-scale using the Pólya-Gamma parameter expansion $[\mu_j|\mu, \xi_{\mu_j}] \sim N(\mu, \xi_{\mu_j}^{-1})$ and $[\mu_0|\sigma_\epsilon, \xi_{\mu_0}] \sim N(\log \sigma_\epsilon^2 - \log T, \xi_{\mu_0}^{-1})$ with $\xi_{\mu_j}, \xi_{\mu_0} \overset{iid}{\sim} \text{PG}(1,0)$ and the identification $\mu_j = \log(\tau_0^2 \tau_j^2)$ and $\mu_0 = \log(\tau_0^2)$.

## 4.1 Time-Varying Parameter Models: Simulations

We conducted a simulation study to evaluate competing variations of the time-varying parameter regression model (10), in particular relative to the proposed dynamic shrinkage process (**DHS**) in (11). Similar to the simulations of Section 3.1, we focus on the distribution of the innovations, $\omega_{j,t}$, and again include the normal-inverse-Gamma (**NIG**) and the (static) horseshoe (**HS**) as competitors, in each case selecting $D = 1$. We also include Belmonte et al. (2014), which uses the Bayesian Lasso as a prior on the innovations (**BL**). Lastly, we include Kalli and Griffin (2014), which offers an alternative approach for dynamic shrinkage (**KG**). Among models with non-dynamic regression coefficients, we include a lasso regression (Tibshirani, 1996) and an ordinary linear regression. These non-dynamic methods were non-competitive and are excluded from the figures.

We simulated 100 data sets of length $T = 200$ from the model $y_t = \boldsymbol{x}_t'\boldsymbol{\beta}_t^* + \epsilon_t$, where the $p = 20$ predictors are $x_{1,t} = 1$ and $x_{j,t} \overset{iid}{\sim} N(0,1)$ for $j > 2$, and $\epsilon_t \overset{iid}{\sim} N(0, \sigma_*^2)$. We also consider autocorrelated predictors $x_{j,t}$ in the supplement with similar results. The true regression coefficients $\boldsymbol{\beta}_t^* = (\beta_{1,t}^*, \ldots, \beta_{p,t}^*)'$ are the following: $\beta_{1,t}^* = 2$ is constant; $\beta_{2,t}^*$ is

piecewise constant with $\beta_{2,t}^* = 0$ everywhere except $\beta_{2,t}^* = 2$ for $t = 41, \ldots, 80$ and $\beta_{2,t}^* = -2$ for $t = 121, \ldots, 160$; $\beta_{3,t}^* = \frac{1}{\sqrt{100}} \sum_{s=1}^{t} Z_s$ with $Z_s \stackrel{iid}{\sim} N(0,1)$ is a scaled random walk for $t \le 100$ and $\beta_{3,t}^* = 0$ for $t > 100$; and $\beta_{j,t}^* = 0$ for $j = 4, \ldots, p = 20$. The predictor set contains a variety of functions: a constant nonzero function, a locally constant function, a slowly-varying function that thresholds to zero for $t > 100$, and 17 true zeros. The noise variance $\sigma_*^2$ is determined by selecting a root-signal-to-noise ratio (RSNR) and computing $\sigma_* = \sqrt{\frac{\sum_{t=1}^{T}(y_t^* - \bar{y}^*)^2}{T-1}} \Big/ \text{RSNR}$, where $y_t^* = \boldsymbol{x}_t' \boldsymbol{\beta}_t^*$ and $\bar{y}^* = \frac{1}{T} \sum_{t=1}^{T} y_t^*$. We select RSNR = 3.

We evaluate competing methods using RMSEs for both $y_t^*$ and $\boldsymbol{\beta}_t^*$ defined by $\text{RMSE}(\hat{y}) = \sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^* - \hat{y}_t)^2}$ and $\text{RMSE}(\hat{\boldsymbol{\beta}}) = \sqrt{\frac{1}{Tp} \sum_{t=1}^{T} \sum_{j=1}^{p} \left( \beta_{j,t}^* - \hat{\beta}_{j,t} \right)^2}$ for all estimators $\hat{\boldsymbol{\beta}}_t$ of the true regression functions, $\boldsymbol{\beta}_t^*$ with $\hat{y}_t = \boldsymbol{x}_t' \hat{\boldsymbol{\beta}}_t$. The results are displayed in Figure 7. The proposed BTF-DHS model substantially outperforms the competitors in both recovery of the true regression functions, $\beta_{j,t}^*$ and estimation of the true curves, $y_t^*$. Our closest competitor is Kalli and Griffin (2014), which also uses dynamic shrinkage, yet is less accurate in estimating the regression coefficients $\beta_{j,t}^*$ and the fitted values $y_t^*$. In addition, our MCMC algorithm is vastly more efficient: for 10,000 MCMC iterations, the Kalli and Griffin (2014) sampler ran for 3 hours and 40 minutes (using Matlab code from Professor Griffin's website), while our proposed algorithm completed in 6 minutes (on a MacBook Pro, 2.7 GHz Intel Core i5).
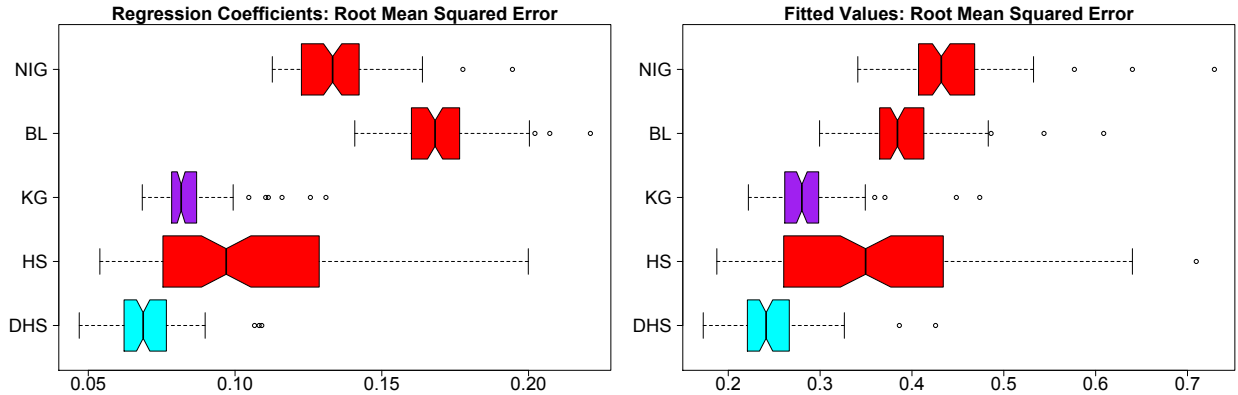


Figure 7: Root mean squared errors for the regression coefficients, $\beta_{j,t}^*$ (**left**) and the true curves, $y_t^* = \boldsymbol{x}_t' \boldsymbol{\beta}_t^*$ (**right**) for simulated data. Non-overlapping notches indicate significant differences between medians.

## 4.2 Time-Varying Parameter Models: The Fama-French Asset Pricing Model

Asset pricing models commonly feature highly structured factor models to parsimoniously model the co-movement of stock returns. Such fundamental factor models identify common risk factors among assets, which may be treated as exogenous predictors in a time series regression. Popular approaches include the one-factor Capital Asset Pricing Model (CAPM, Sharpe, 1964) and the three-factor Fama-French model (FF-3, Fama and French, 1993). Recently, the five-factor Fama-French model (FF-5, Fama and French, 2015) was proposed as an extension of FF-3 to incorporate additional common risk factors. However, outstanding questions remain regarding which, and how many, factors are necessary. Importantly, an attempt to address these questions must consider the dynamic component: the relevance of individual factors may change over time, particularly for different assets.

We apply model (10) to extend these fundamental factor models to the dynamic setting, in which the factor loadings are permitted to vary—perhaps rapidly—over time. For further generality, we append the momentum factor of Carhart (1997) to FF-5 to produce a fundamental factor model with six factors and dynamic factor loadings. Importantly, the shrinkage towards sparsity induced by the dynamic horseshoe process allows the model to effectively select out unimportant factors, which also may change over time. As in Section 3.2, we modify model (10) to include stochastic volatility for the observation error, $[\epsilon_t | \{\sigma_s\}] \overset{indep}{\sim} N(0, \sigma_t^2)$.

To study various market sectors, we use weekly industry portfolio data from the website of Kenneth R. French, which provide the value-weighted return of stocks in the given industry. We focus on manufacturing (Manuf) and healthcare (Hlth). For a given industry portfolio, the response variable is the returns in excess of the risk free rate, $y_t = R_t - R_{F,t}$, with predictors $\boldsymbol{x}_t = (1, R_{M,t} - R_{F,t}, SMB_t, HML_t, RMW_t, CMA_t, MOM_t)'$, defined as follows: the *market risk factor*, $R_{M,t} - R_{F,t}$ is the return on the market portfolio $R_{M,t}$ in excess of the risk free rate $R_{F,t}$; the *size factor*, $SMB_t$ (small minus big) is the difference in returns between

portfolios of small and large market value stocks; the *value factor*, $HML_t$ (high minus low) is the difference in returns between portfolios of high and low book-to-market value stocks; the *profitability factor*, $RMW_t$ is the difference in returns between portfolios of robust and weak profitability stocks; the *investment factor*, $CMA_t$ is the difference in returns between portfolios of stocks of low and high investment firms; and the *momentum factor*, $MOM_t$ is the difference in returns between portfolios of stocks with high and low prior returns. These data are publicly available on Kenneth R. French's website, which provides additional details on the portfolios. We standardize all predictors and the response to have unit variance.

We conduct inference on the time-varying regression coefficients $\beta_{j,t}$ using simultaneous credible bands. Unlike pointwise credible intervals, simultaneous credible bands control for multiple testing, and may be computed as in Ruppert et al. (2003). Letting $B_{j,t}(\alpha)$ denote the $(1-\alpha)\%$ simultaneous credible band for predictor $j$ at time $t$, we compute *Simultaneous Band Scores* (SimBaS; Meyer et al., 2015), $P_{j,t} = \min\{\alpha : 0 \notin B_{j,t}(\alpha)\}$. The SimBaS $P_{j,t}$ indicate the minimum level for which the simultaneous bands do not include zero, while controlling for multiple testing, and therefore may be used to detect which predictors $j$ are important at time $t$. Globally, we compute *global Bayesian p-values* (GBPVs; Meyer et al., 2015), $P_j = \min_t\{P_{j,t}\}$ for each predictor $j$, which indicate whether or not a predictor is important at *any* time $t$. SimBaS and GBPVs have proven useful in functional regression models, but also are suitable for time-varying parameter regression models to identify important predictors while controlling for multiple testing.

In Figures 8 and 9, we plot the posterior expectation and credible bands for the time-varying regression coefficients and observation error stochastic volatility for the weekly manufacturing and healthcare industry data sets, respectively, from 4/1/2007 - 4/1/2017 ($T = 522$). For the manufacturing industry, the important factors identified by the GBPVs at the 5% level are the market risk ($R_{M,t} - R_{F,t}$, GBPV $= 0.000$), investment ($CMA_t$, GBPV $= 0.024$), and momentum ($MOM_t$, GBPV $= 0.019$). However, the SimBaS $P_{j,t}$ for $CMA_t$ and $MOM_t$ are below 5% only for brief periods (red lines), which suggests that these

27

important effects are intermittent. For the healthcare industry, the GBPVs identify market risk (GBPV = 0.001) and value ($HML_t$, GBPV = 0.023) as the only important factors. Notably, the only common factor flagged by GBPVs in both the manufacturing and healthcare industries under model (10) over this time period is the market risk. This result suggests that the aggressive shrinkage behavior of the dynamic shrinkage process is important in this setting, since several factors may be effectively irrelevant for some or all time points.



Figure 8: Posterior expectations (cyan), 95% pointwise HPD credible intervals (light gray) and 95% simultaneous credible bands (dark gray) for $\beta_{j,t}$ and $\sigma_t$ (bottom right) under the BTF-DHS model given by (10) and (11) for value-weighted manufacturing industry returns. The solid black line is zero, the dashed green line is the ordinary linear regression estimate, and the solid red line indicates periods for which the 95% simultaneous credible bands do not contain zero, or, equivalently, $P_{j,t}$ (SimBaS) is less than 0.05.

# 5  MCMC Sampling Algorithm and Computational Details

We design a Gibbs sampling algorithm for the dynamic shrinkage process. The sampling algorithm is both computationally and MCMC efficient, and builds upon two main components: (i) a log-variance sampling algorithm (Kastner and Frühwirth-Schnatter, 2014)
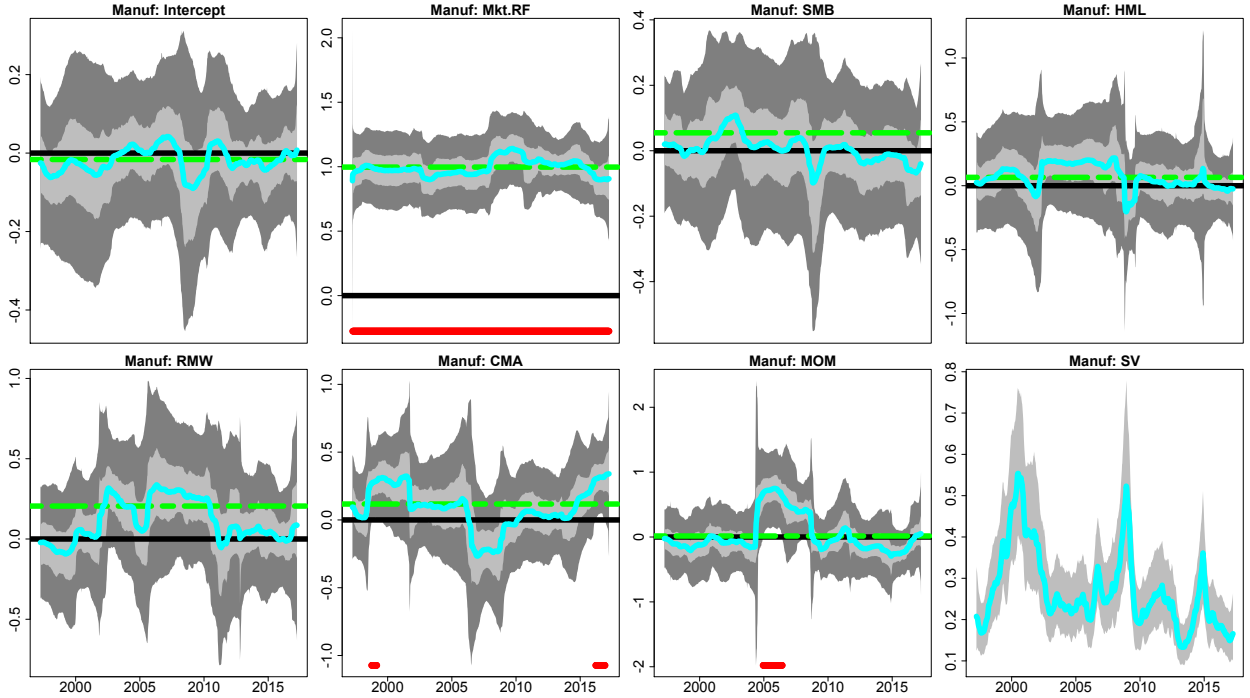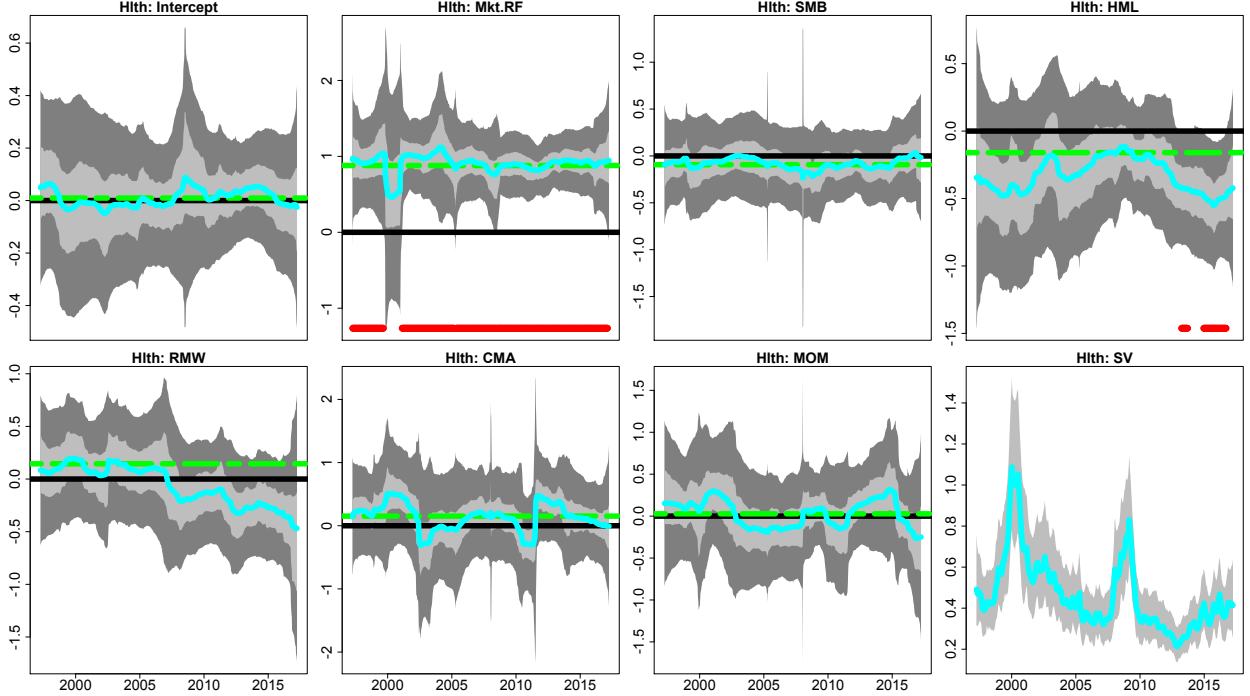
Figure 9: Posterior expectations (cyan), 95% pointwise HPD credible intervals (light gray) and 95% simultaneous credible bands (dark gray) for $\beta_{j,t}$ and $\sigma_t$ (bottom right) under the BTF-DHS model given by (10) and (11) for value-weighted healthcare industry returns. The solid black line is zero, the dashed green line is the ordinary linear regression estimate, and the solid red line indicates periods for which the 95% simultaneous credible bands do not contain zero, or, equivalently, $P_{j,t}$ (SimBaS) is less than 0.05.

augmented with a Pólya-Gamma sampler (Polson et al., 2013); and (ii) a Cholesky Factor Algorithm (CFA, Rue, 2001) for sampling the state variables in the dynamic linear model. Importantly, both components employ algorithms that are linear in the number of time points, which produces a highly efficient sampling algorithm.

The general sampling algorithm is as follows: (i) sample the dynamic shrinkage components (the log-volatilities $\{h_t\}$, the Pólya-Gamma mixing parameters $\{\xi_t\}$, the unconditional mean of log-variance $\mu$, the AR(1) coefficient of log-variance $\phi$, and the discrete mixture component indicators $\{s_t\}$); (ii) sample the state variables $\{\boldsymbol{\beta}_t\}$; and (iii) sample the observation error variance $\sigma_\epsilon^2$. We provide details of the dynamic shrinkage process sampling algorithm in Section 5.1 and include the details for sampling steps (ii) and (iii) in the supplement.

## 5.1 Efficient Sampling for the Dynamic Shrinkage Process

Consider the (univariate) dynamic shrinkage process in (4) with the Pólya-Gamma parameter expansion of Theorem 4. We provide implementation details for the dynamic horseshoe process with $\alpha = \beta = 1/2$, but extensions to other cases are straightforward. The sampling framework of Kastner and Frühwirth-Schnatter (2014) represents the likelihood for $h_t$ on the log-scale, and approximates the ensuing $\log \chi_1^2$ distribution for the errors via a known discrete mixture of Gaussian distributions. In particular, let $\tilde{y}_t = \log(\omega_t^2 + c)$, where $c$ is a small offset to avoid numerical issues. Conditional on the mixture component indicators $s_t$, the likelihood is $\tilde{y}_t \stackrel{indep}{\sim} N(h_t + m_{s_t}, v_{s_t})$ where $m_i$ and $v_i, i = 1, \ldots, 10$ are the pre-specified mean and variance components of the 10-component Gaussian mixture provided in Omori et al. (2007). Under model (4), the evolution equation is $h_{t+1} = \mu + \phi(h_t - \mu) + \eta_t$ with initialization $h_1 = \mu + \eta_0$ and innovations $[\eta_t | \xi_t] \stackrel{indep}{\sim} N(0, \xi_t^{-1})$ for $[\xi_t] \stackrel{iid}{\sim} \text{PG}(1, 0)$. Note that model (2) provides a more general setting, which similarly may be combined with the Gaussian likelihood for $\tilde{y}_t$ above.

To sample $\boldsymbol{h} = (h_1, \ldots, h_T)$ jointly, we directly compute the posterior distribution of $\boldsymbol{h}$ and exploit the tridiagonal structure of the resulting posterior precision matrix. In particular, we equivalently have $\tilde{\boldsymbol{y}} \sim N(\boldsymbol{m} + \tilde{\boldsymbol{h}} + \tilde{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_v)$ and $\boldsymbol{D}_\phi \tilde{\boldsymbol{h}} \sim N(\boldsymbol{0}, \boldsymbol{\Sigma}_\xi)$, where $\boldsymbol{m} = (m_{s_1}, \ldots, m_{s_T})'$, $\tilde{\boldsymbol{h}} = (h_1 - \mu, \ldots, h_T - \mu)'$, $\tilde{\boldsymbol{\mu}} = (\mu, (1 - \phi)\mu, \ldots, (1 - \phi)\mu)'$, $\boldsymbol{\Sigma}_v = \text{diag}\left(\{v_{s_t}\}_{t=1}^T\right)$, $\boldsymbol{\Sigma}_\xi = \text{diag}\left(\{\xi_t^{-1}\}_{t=1}^T\right)$, and $\boldsymbol{D}_\phi$ is a lower triangular matrix with ones on the diagonal, $-\phi$ on the first off-diagonal, and zeros elsewhere. We sample from the posterior distribution of $\boldsymbol{h}$ by sampling from the posterior distribution of $\tilde{\boldsymbol{h}}$ and setting $\boldsymbol{h} = \tilde{\boldsymbol{h}} + \mu\boldsymbol{1}$ for $\boldsymbol{1}$ a $T$-dimensional vector of ones. The required posterior distribution is $\tilde{\boldsymbol{h}} \sim N\left(\boldsymbol{Q}_{\tilde{h}}^{-1}\boldsymbol{\ell}_{\tilde{h}}, \boldsymbol{Q}_{\tilde{h}}^{-1}\right)$, where $\boldsymbol{Q}_{\tilde{h}} = \boldsymbol{\Sigma}_v^{-1} + \boldsymbol{D}_\phi'\boldsymbol{\Sigma}_\xi^{-1}\boldsymbol{D}_\phi$ is a tridiagonal symmetric matrix with diagonal elements

$d_0(Q_{\tilde{h}})$ and first off-diagonal elements $d_1(Q_{\tilde{h}})$ defined as

$$d_0(Q_{\tilde{h}}) = \left[(v_{s_1}^{-1} + \xi_1 + \phi^2\xi_2), (v_{s_2}^{-1} + \xi_2 + \phi^2\xi_3), \ldots, (v_{s_{T-1}}^{-1} + \xi_{T-1} + \phi^2\xi_T), (v_{s_T}^{-1} + \xi_T)\right],$$

$$d_1(Q_{\tilde{h}}) = \left[(-\phi\xi_2), (-\phi\xi_3), \ldots, (-\phi\xi_{T-1})\right], \text{ and}$$

$$\ell_{\tilde{h}} = \Sigma_v^{-1}(\tilde{y} - m - \tilde{\mu}) = \left[\frac{\tilde{y}_1 - m_{s_1} - \mu}{v_{s_1}}, \frac{\tilde{y}_2 - m_{s_2} - (1-\phi)\mu}{v_{s_2}}, \ldots, \frac{\tilde{y}_T - m_{s_T} - (1-\phi)\mu}{v_{s_T}}\right]'.$$

Drawing from this posterior distribution is straightforward and efficient, using band back-substitution described in Kastner and Frühwirth-Schnatter (2014): (i) compute the Cholesky decomposition $Q_{\tilde{h}} = LL'$, where $L$ is lower triangle; (ii) solve $La = \ell_{\tilde{h}}$ for $a$; and (iii) solve $L'\tilde{h} = a + e$ for $\tilde{h}$, where $e \sim N(0, I_T)$.

Conditional on the log-volatilities $\{h_t\}$, we sample the AR(1) evolution parameters: the log-innovation precisions $\{\xi_t\}$, the autoregressive coefficient $\phi$, and the unconditional mean $\mu$. The precisions are distributed $[\xi_t|\eta_t] \sim \text{PG}(1, \eta_t)$ for $\eta_t = h_{t+1} - \mu - \phi(h_t - \mu)$, which we sample using the `rpg()` function in the `R` package `BayesLogit` (Polson et al., 2013). The Pólya-Gamma sampler is efficient: using only exponential and inverse-Gaussian draws, Polson et al. (2013) construct an accept-reject sampler for which the probability of acceptance is uniformly bounded below at 0.99919, which does not require any tuning. Next, we assume the prior $[(\phi + 1)/2] \sim \text{Beta}(a_\phi, b_\phi)$, which restricts $|\phi| < 1$ for stationarity, and sample from the full conditional distribution of $\phi$ using the slice sampler of Neal (2003). We select $a_\phi = 10$ and $b_\phi = 2$, which places most of the mass for the density of $\phi$ in $(0,1)$ with a prior mean of $2/3$ and a prior mode of $4/5$ to reflect the likely presence of persistent volatility clustering. The prior for the global scale parameter is $\tau \sim C^+(0, \sigma_\epsilon/\sqrt{T})$, which implies $\mu = \log(\tau^2)$ is $[\mu|\sigma_\epsilon, \xi_\mu] \sim N(\log(\sigma_\epsilon^2/T), \xi_\mu^{-1})$ with $\xi_\mu \sim \text{PG}(1, 0)$. Including the initialization $h_1 \sim N(\mu, \xi_0^{-1})$ with $\xi_0 \sim \text{PG}(1, 0)$, the posterior distribution for $\mu$ is $\mu \sim N(Q_\mu^{-1}\ell_\mu, Q_\mu^{-1})$ with $Q_\mu = \xi_\mu + \xi_0 + (1-\phi)^2 \sum_{t=1}^{T-1} \xi_t$ and $\ell_\mu = \xi_\mu \log(\sigma_\epsilon^2/T) + \xi_0 h_1 + (1-\phi) \sum_{t=1}^{T-1} \xi_t(h_{t+1} - \phi h_t)$. Sampling $\xi_\mu$ and $\xi_0$ follows the Pólya-Gamma sampling scheme above.

Finally, we sample the discrete mixture component indicators $s_t$. The discrete mixture

probabilities are straightforward to compute: the prior mixture probabilities are the mixing proportions given by Omori et al. (2007) and the likelihood is $\tilde{y}_t \stackrel{indep}{\sim} N(h_t + m_{s_t}, v_{s_t})$; see Kastner and Frühwirth-Schnatter (2014) for details.

Note that the use of the discrete mixture approximation for log-variance models as a component within a larger MCMC sampling algorithm has been used widely in the literature (Clark, 2011; D'Agostino et al., 2013; Belmonte et al., 2014; Carriero et al., 2015; Kastner et al., 2017).

# 6  Discussion and Future Work

Dynamic shrinkage processes provide a computationally convenient and widely applicable mechanism for incorporating adaptive shrinkage and regularization into existing models. By extending a broad class of global-local shrinkage priors to the dynamic setting, the resulting processes inherit the desirable shrinkage behavior, but with greater time-localization. The success of dynamic shrinkage processes suggests that other priors may benefit from log-scale or other appropriate representations, with or without additional dependence modeling.

As demonstrated in Sections 3 and 4, dynamic shrinkage processes are particularly appropriate for dynamic linear models, including trend filtering and time-varying parameter regression. In both settings, the dynamic linear models with dynamic horseshoe innovations outperform all competitors in simulated data, and produce reasonable and interpretable results for real data applications. Dynamic shrinkage processes may be useful in other dynamic linear models, such as incorporating seasonality or change points with appropriately-defined (dynamic) shrinkage. Given the exceptional curve-fitting capabilities of the Bayesian trend filtering model (9) with dynamic horseshoe innovations (BTF-DHS), a natural extension would be to incorporate the BTF-DHS into more general additive, functional, or longitudinal data models in order to capture irregular or local curve features.

An important extension of the dynamic fundamental factor model of Section 4.2 is to

incorporate a large number of assets, possibly with residual correlation among stock returns beyond the common factors of FF-5. Building upon Carvalho et al. (2011), a reasonable approach may be to combine a set of known factors, such as the Fama-French factors, with a set of unknown factors to be estimated from the data, where *both* sets of factor loadings are endowed with dynamic shrinkage processes to provide greater adaptability yet sufficient capability for shrinkage of irrelevant factors.

Another promising area for applications of the proposed methodology is compressive sensing and signal processing, which commonly rely on approximations for estimation and prediction (e.g., Ziniel and Schniter, 2013; Wang et al., 2016). The linear time complexity of our MCMC algorithm for Bayesian trend filtering with dynamic shrinkage may offer the computational scalability to provide full Bayesian inference, and perhaps improved prediction accompanied by adequate uncertainty quantification, which is notably absent from the papers cited above.

# References

Abramovich, F., Sapatinas, T., and Silverman, B. W. (1998). Wavelet thresholding via a Bayesian approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(4):725–749.

Armagan, A., Clyde, M., and Dunson, D. B. (2011). Generalized Beta mixtures of Gaussians. In *Advances in neural information processing systems*, pages 523–531.

Arnold, T. B. and Tibshirani, R. J. (2014). `genlasso`: *Path algorithm for generalized lasso problems*. `R` package version 1.3.

Bae, K. and Mallick, B. K. (2004). Gene selection using a two-level hierarchical Bayesian model. *Bioinformatics*, 20(18):3423–3430.

Barndorff-Nielsen, O., Kent, J., and Sørensen, M. (1982). Normal variance-mean mixtures

and z distributions. *International Statistical Review/Revue Internationale de Statistique*, pages 145–159.

Belmonte, M. A., Koop, G., and Korobilis, D. (2014). Hierarchical shrinkage in time-varying parameter models. *Journal of Forecasting*, 33(1):80–94.

Berger, J. (1980). A robust generalized Bayes estimator and confidence region for a multivariate normal mean. *The Annals of Statistics*, pages 716–761.

Bitto, A. and Frühwirth-Schnatter, S. (2016). Achieving shrinkage in a time-varying parameter model framework. *arXiv preprint arXiv:1611.01310*.

Carhart, M. M. (1997). On persistence in mutual fund performance. *The Journal of finance*, 52(1):57–82.

Carriero, A., Clark, T. E., and Marcellino, M. (2015). Realtime nowcasting with a Bayesian mixed frequency model with stochastic volatility. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 178(4):837–862.

Carvalho, C. M., Lopes, H. F., and Aguilar, O. (2011). Dynamic stock selection strategies: A structured factor model framework. *Bayesian Statistics*, 9:1–21.

Carvalho, C. M., Polson, N. G., and Scott, J. G. (2009). Handling sparsity via the horseshoe. In *AISTATS*, volume 5, pages 73–80.

Carvalho, C. M., Polson, N. G., and Scott, J. G. (2010). The horseshoe estimator for sparse signals. *Biometrika*, pages 465–480.

Chan, J. C. (2013). Moving average stochastic volatility models with application to inflation forecast. *Journal of Econometrics*, 176(2):162–172.

Chan, J. C. and Jeliazkov, I. (2009). Efficient simulation and integrated likelihood estimation in state space models. *International Journal of Mathematical Modelling and Numerical Optimisation*, 1(1-2):101–120.

Chan, J. C., Koop, G., Leon-Gonzalez, R., and Strachan, R. W. (2012). Time varying dimension models. *Journal of Business & Economic Statistics*, 30(3):358–367.

Clark, T. E. (2011). Real-time density forecasts from Bayesian vector autoregressions with stochastic volatility. *Journal of Business & Economic Statistics*, 29(3):327–341.

Constantine, W. and Percival, D. (2016). *wmtsa: Wavelet Methods for Time Series Analysis*. R package version 2.0-2.

D'Agostino, A., Gambetti, L., and Giannone, D. (2013). Macroeconomic forecasting and structural change. *Journal of Applied Econometrics*, 28(1):82–101.

Dangl, T. and Halling, M. (2012). Predictive regressions with time-varying coefficients. *Journal of Financial Economics*, 106(1):157–181.

Datta, J. and Ghosh, J. K. (2013). Asymptotic properties of Bayes risk for the horseshoe prior. *Bayesian Analysis*, 8(1):111–132.

Donoho, D. L. and Johnstone, J. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3):425–455.

Durbin, J. and Koopman, S. J. (2002). A simple and efficient simulation smoother for state space time series analysis. *Biometrika*, 89(3):603–616.

Elerian, O., Chib, S., and Shephard, N. (2001). Likelihood inference for discretely observed nonlinear diffusions. *Econometrica*, 69(4):959–993.

Fama, E. F. and French, K. R. (1993). Common risk factors in the returns on stocks and bonds. *Journal of Financial Economics*, 33(1):3–56.

Fama, E. F. and French, K. R. (2015). A five-factor asset pricing model. *Journal of Financial Economics*, 116(1):1–22.

Faulkner, J. R. and Minin, V. N. (2016). Locally adaptive smoothing with Markov random fields and shrinkage priors. *Bayesian Analysis*.

Figueiredo, M. A. (2003). Adaptive sparseness for supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(9):1150–1159.

Frühwirth-Schnatter, S. and Wagner, H. (2010). Stochastic model specification search for Gaussian and partial non-Gaussian state space models. *Journal of Econometrics*, 154(1):85–100.

Griffin, J. E. and Brown, P. J. (2005). Alternative prior distributions for variable selection with very many more variables than observations. Technical report, University of Warwick, Centre for Research in Statistical Methodology.

Griffin, J. E. and Brown, P. J. (2010). Inference with normal-gamma prior distributions in regression problems. *Bayesian Analysis*, 5(1):171–188.

Hastie, T. and Tibshirani, R. (1993). Varying-coefficient models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 757–796.

James, N. A., Kejariwal, A., and Matteson, D. S. (2016). Leveraging cloud data to mitigate user experience from 'Breaking Bad'. In *2016 IEEE International Conference on Big Data*, pages 3499–3508. IEEE.

Joe, H. (2015). *Dependence modeling with copulas.* CRC Press.

Kalli, M. and Griffin, J. E. (2014). Time-varying sparsity in dynamic regression models. *Journal of Econometrics*, 178(2):779–793.

Kastner, G. (2016). Dealing with stochastic volatility in time series using the R package stochvol. *Journal of Statistical Software*, 69(5):1–30.

Kastner, G. and Frühwirth-Schnatter, S. (2014). Ancillarity-sufficiency interweaving strategy (ASIS) for boosting MCMC estimation of stochastic volatility models. *Computational Statistics & Data Analysis*, 76:408–423.

Kastner, G., Frühwirth-Schnatter, S., and Lopes, H. F. (2017). Efficient Bayesian inference for multivariate factor stochastic volatility models. *Journal of Computational and Graphical Statistics*, 26(4):905–917.

Kim, S., Shephard, N., and Chib, S. (1998). Stochastic volatility: likelihood inference and comparison with ARCH models. *The Review of Economic Studies*, 65(3):361–393.

Kim, S.-J., Koh, K., Boyd, S., and Gorinevsky, D. (2009). $\ell_1$ Trend Filtering. *SIAM review*, 51(2):339–360.

Korobilis, D. (2013). Hierarchical shrinkage priors for dynamic regressions with many predictors. *International Journal of Forecasting*, 29(1):43–59.

Kyung, M., Gill, J., Ghosh, M., and Casella, G. (2010). Penalized regression, standard errors, and Bayesian lassos. *Bayesian Analysis*, 5(2):369–411.

McCausland, W. J., Miller, S., and Pelletier, D. (2011). Simulation smoothing for state–space models: A computational efficiency analysis. *Computational Statistics & Data Analysis*, 55(1):199–212.

Meyer, M. J., Coull, B. A., Versace, F., Cinciripini, P., and Morris, J. S. (2015). Bayesian function-on-function regression for multilevel functional data. *Biometrics*, 71(3):563–574.

Nakajima, J. and West, M. (2013). Bayesian analysis of latent threshold dynamic models. *Journal of Business & Economic Statistics*, 31(2):151–164.

Nason, G. (2016). `wavethresh`: *Wavelets Statistics and Transforms*. `R` package version 4.6.8.

Neal, R. M. (2003). Slice sampling. *Annals of Statistics*, pages 705–741.

Omori, Y., Chib, S., Shephard, N., and Nakajima, J. (2007). Stochastic volatility with leverage: Fast and efficient likelihood inference. *Journal of Econometrics*, 140(2):425–449.

Piironen, J. and Vehtari, A. (2016). On the hyperprior choice for the global shrinkage parameter in the horseshoe prior. *arXiv preprint arXiv:1610.05559*.

Polson, N. G. and Scott, J. G. (2010). Shrink globally, act locally: sparse Bayesian regularization and prediction. *Bayesian Statistics*, 9:501–538.

Polson, N. G. and Scott, J. G. (2012a). Local shrinkage rules, Lévy processes and regularized regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(2):287–311.

Polson, N. G. and Scott, J. G. (2012b). On the half-Cauchy prior for a global scale parameter. *Bayesian Analysis*, 7(4):887–902.

Polson, N. G., Scott, J. G., and Windle, J. (2013). Bayesian inference for logistic models using Pólya–Gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349.

Ramsay, J. and Silverman, B. (2005). *Functional Data Analysis*. Springer.

Rockova, V. and McAlinn, K. (2017). Dynamic variable selection with spike-and-slab process priors. *arXiv preprint arXiv:1708.00085*.

Rue, H. (2001). Fast sampling of Gaussian Markov random fields. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2):325–338.

Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric regression*. Number 12. Cambridge University Press.

Sharpe, W. F. (1964). Capital asset prices: A theory of market equilibrium under conditions of risk. *The Journal of Finance*, 19(3):425–442.

Strawderman, W. E. (1971). Proper bayes minimax estimators of the multivariate normal mean. *The Annals of Mathematical Statistics*, 42(1):385–388.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288.

Tibshirani, R. J. (2014). Adaptive piecewise polynomial estimation via trend filtering. *The Annals of Statistics*, 42(1):285–323.

van der Pas, S., Kleijn, B., and van der Vaart, A. (2014). The horseshoe estimator: Posterior concentration around nearly black vectors. *Electronic Journal of Statistics*, 8(2):2585–2618.

Wang, H., Yu, H., Hoy, M., Dauwels, J., and Wang, H. (2016). Variational Bayesian dynamic compressive sensing. In *2016 IEEE International Symposium on Information Theory (ISIT)*, pages 1421–1425. IEEE.

Zhu, B. and Dunson, D. B. (2013). Locally adaptive Bayes nonparametric regression via nested Gaussian processes. *Journal of the American Statistical Association*, 108(504):1445–1456.

Ziniel, J. and Schniter, P. (2013). Dynamic compressive sensing of time-varying signals via approximate message passing. *IEEE transactions on signal processing*, 61(21):5270–5284.

# A    Proofs

*Proof.* (Proposition 1) Proposition 1 follows from Proposition 2 with $\mu_z = 0$.    □

*Proof.* (Proposition 2) Let $\eta \sim Z(\alpha, \beta, \mu_z, 1)$ with density (3), i.e.,

$$[z] = \big[\sigma B(\alpha, \beta)\big]^{-1}\big\{\exp\big[(z - \mu_z)/\sigma_z\big]\big\}^\alpha\big\{1 + \exp\big[(z - \mu_z)/\sigma_z\big]\big\}^{-(\alpha+\beta)}.$$

The density of $\lambda^2 = \exp(\eta)$ is

$$\begin{aligned}[\lambda^2] &\propto \big(\lambda^2\big)^{-1}\big\{\exp\big[\log(\lambda^2) - \mu_z\big]\big\}^\alpha\big\{1 + \exp\big[\log(\lambda^2) - \mu_z\big]\big\}^{-(\alpha+\beta)}\\ &\propto \big(\lambda^2\big)^{\alpha-1}\big[1 + \lambda^2/\exp(\mu_z)\big]^{-(\alpha+\beta)}\end{aligned}$$

and therefore the density of $\kappa = 1/(1 + \lambda^2)$ is

$$[\kappa] \propto \kappa^{-2}\left[\kappa^{-1} - 1\right]^{\alpha-1}\left[1 + (\kappa^{-1} - 1)/\exp(\mu_z)\right]^{-(\alpha+\beta)}$$

$$\propto \kappa^{-2-(\alpha-1)}(1 - \kappa)^{\alpha-1}\left\{\kappa^{-1}\left[\kappa\exp(\mu_z) + (1 - \kappa)\right]\right\}^{-(\alpha+\beta)}$$

$$\propto (1 - \kappa)^{\alpha-1}\kappa^{\beta-1}\left[\kappa\exp(\mu_z) + (1 - \kappa)\right]^{-(\alpha+\beta)}$$

i.e., $\kappa \sim \mathrm{TPB}(\beta, \alpha, \exp(\mu_z))$. $\qquad\square$

*Proof.* (Theorem 1) Under model (4), i.e.,

$$h_{t+1} = \mu + \phi(h_t - \mu) + \eta_t, \quad \eta_t \overset{iid}{\sim} Z(\alpha, \beta, 0, 1),$$

we have $[h_{t+1}|h_t, \phi, \mu] \sim Z(\alpha, \beta, \mu + \phi(h_t - \mu), 1)$. Using Proposition 2, the conditional distribution for $\kappa_{t+1}$ is $[\kappa_{t+1}|h_t, \phi, \mu] \sim \mathrm{TPB}(\beta, \alpha, \exp(\mu + \phi(h_t - \mu)))$. By substituting $\tau = \exp(\mu)$ and $\lambda_t = \exp(h_t - \mu)$, we equivalently have $[\kappa_{t+1}|\lambda_t, \phi, \tau] \sim \mathrm{TPB}(\beta, \alpha, \tau^2\lambda_t^{2\phi})$. Noting $\tau^2\lambda_t^{2\phi} = \tau^{2(1-\phi)}\left[\frac{1-\kappa_t}{\kappa_t}\right]^{\phi}$ completes the proof. $\qquad\square$

*Proof.* (Theorem 2) Let $\gamma_t = \left[(1-\kappa_t)/\kappa_t\right]^{\phi}$ and note that $\kappa \mapsto \kappa^{-1/2}$ and $\kappa \mapsto \left[1+(\gamma_t-1)\kappa\right]^{-1}$ are decreasing in $\kappa$ for $\gamma_t > 1$. It follows that, for $\gamma_t > 1$,

$$\mathbb{P}\left(\kappa_{t+1} > \varepsilon \big| \{\kappa_s\}_{s\leq t}, \phi\right) = \int_{\varepsilon}^{1} \pi^{-1}\gamma_t^{1/2}\kappa_{t+1}^{-1/2}(1 - \kappa_{t+1})^{-1/2}\left[1 + (\gamma_t - 1)\kappa_{t+1}\right]^{-1} d\kappa_{t+1}$$

$$\leq \pi^{-1}\gamma_t^{1/2}\varepsilon^{-1/2}\left[1 + (\gamma_t - 1)\varepsilon\right]^{-1}\int_{\varepsilon}^{1}(1 - \kappa_{t+1})^{-1/2} d\kappa_{t+1}$$

$$\leq 2\pi^{-1}\varepsilon^{-1/2}(1 - \varepsilon)^{1/2}\frac{\gamma_t^{1/2}}{1 + (\gamma_t - 1)\varepsilon}$$

converges to zero as $\kappa_t \to 0$, since $\kappa_t \to 0$ implies $\gamma_t \to \infty$. $\qquad\square$

*Proof.* (Theorem 3) Marginalizing over $\omega_t$, the likelihood is $[y_{t+1}|\{\kappa_s\}] \overset{indep}{\sim} N(0, \kappa_{t+1}^{-1})$. From

Theorem 1, the posterior distribution of $\kappa_{t+1}$ may be computed as

$$[\kappa_{t+1}|y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau] \propto \left\{ \kappa_{t+1}^{\beta-1}(1 - \kappa_{t+1})^{\alpha-1} \left[1 + (\gamma_t - 1)\kappa_{t+1}\right]^{-(\alpha+\beta)} \right\}$$

$$\times \left\{ \kappa_{t+1}^{1/2} \exp\left(-y_{t+1}^2 \kappa_{t+1}/2\right) \right\}$$

$$\propto (1 - \kappa_{t+1})^{-1/2} \left[1 + (\gamma_t - 1)\kappa_{t+1}\right]^{-1} \exp\left(-y_{t+1}^2 \kappa_{t+1}/2\right)$$

for $\alpha = \beta = 1/2$, where $\gamma_t = \tau^{2(1-\phi)}\left[(1 - \kappa_t)/\kappa_t\right]^\phi$. Defining $p_1(\kappa) = (1 - \kappa)^{-1/2}$, $p_2(\kappa|\gamma_t) = \left[1 + (\gamma_t - 1)\kappa\right]^{-1}$, and $p_3(\kappa|y_{t+1}) = \exp\left(-y_{t+1}^2 \kappa/2\right)$ for $\kappa \in (0, 1)$, observe that $p_1(\cdot)$ is increasing in $\kappa$, $p_2(\kappa|\gamma_t) \leq \left[p_1(\kappa)\right]^2$ for all $\gamma_t \geq 0$, and $p_3(\cdot)$ is decreasing in $\kappa$. Similar to Datta and Ghosh (2013), the following inequalities hold for $\varepsilon \in (0, 1)$ with $\varepsilon' = 1 - \varepsilon$:

$$\mathbb{P}\left(\kappa_{t+1} < \varepsilon'|y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau\right) \leq \frac{\mathbb{P}\left(\kappa_{t+1} < \varepsilon'|y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau\right)}{\mathbb{P}\left(\kappa_{t+1} > \varepsilon'|y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau\right)}$$

$$\leq \frac{\int_0^{\varepsilon'} (1 - \kappa_{t+1})^{-3/2} \exp\left(-y_{t+1}^2 \kappa_{t+1}/2\right) d\kappa_{t+1}}{\int_{\varepsilon'}^1 \left[1 + (\gamma_t - 1)\kappa_{t+1}\right]^{-3/2} \exp\left(-y_{t+1}^2 \kappa_{t+1}/2\right) d\kappa_{t+1}}$$

$$\leq \frac{\int_0^{\varepsilon'} (1 - \kappa_{t+1})^{-3/2} d\kappa_{t+1}}{\exp\left(-y_{t+1}^2/2\right) \int_{\varepsilon'}^1 \left[1 + (\gamma_t - 1)\kappa_{t+1}\right]^{-3/2} d\kappa_{t+1}}$$

$$\leq \frac{2\left[(1 - \varepsilon')^{-1/2} - 1\right]}{\exp\left(-y_{t+1}^2/2\right) 2(\gamma_t - 1)^{-1}\left\{\left[1 + (\gamma_t - 1)\varepsilon'\right]^{-1/2} - \gamma_t^{-1/2}\right\}}$$

$$\leq \left[(1 - \varepsilon')^{-1/2} - 1\right] \exp\left(y_{t+1}^2/2\right) \gamma_t^{1/2}$$

$$\times \left\{ \frac{1 - \gamma_t}{1 - \gamma_t^{1/2}/[1 + (\gamma_t - 1)\varepsilon']^{1/2}} \right\}.$$

Noting the final term in curly braces converges to 1 as $\gamma_t \to 0$, we obtain $\mathbb{P}\left(\kappa_{t+1} < 1 - \varepsilon|y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau\right) \to 0$ as $\gamma_t \to 0$. The result for (a) follows immediately.

For $\varepsilon \in (0, 1)$ and $\gamma_t < 1$, and observing that $p_2(\kappa|\gamma_t)$ is increasing in $\kappa$ for $\gamma_t < 1$, then

for any $\delta \in (0,1)$,

$$
\begin{aligned}
\mathbb{P}\big(\kappa_{t+1} > \varepsilon \big| y_{t+1}, \{\kappa_s\}_{s \leq t}, \phi, \tau\big) &\leq \frac{\gamma_t^{-1} \exp\left(-y_{t+1}^2 \varepsilon/2\right) \int_\varepsilon^1 (1-\kappa_{t+1})^{-1/2} \, d\kappa_{t+1}}{\int_0^{\delta\epsilon'} \exp\left(-y_{t+1}^2 \kappa_{t+1}/2\right) \, d\kappa_{t+1}}, \\
&\leq \frac{\gamma_t^{-1} \exp\left(-y_{t+1}^2 \varepsilon/2\right) 2(1-\varepsilon)^{1/2}}{\exp\left(-y_{t+1}^2 \delta\varepsilon/2\right) \delta\varepsilon} \\
&= \exp\left(-y_{t+1}^2 \varepsilon[1-\delta]/2\right) \gamma_t^{-1} 2(1-\varepsilon)^{1/2}(\delta\varepsilon)^{-1}
\end{aligned}
$$

which converges to zero as $|y_{t+1}| \to \infty$, proving (b). $\qquad\square$

*Proof.* (Theorem 4) The density of $\eta \sim Z(\alpha, \beta, 0, 1)$ may be written

$$
\begin{aligned}
[\eta] &= \frac{1}{B(\alpha, \beta)} \frac{[\exp(\eta)]^\alpha}{[1 + \exp(\eta)]^{\alpha+\beta}} \\
&= \frac{1}{B(\alpha, \beta)} 2^{-(\alpha+\beta)} \exp\{\eta[\alpha - (\alpha+\beta)/2]\} \int_0^\infty \exp(-\eta^2 \xi/2) p_{\alpha+\beta}(\xi) \, d\xi
\end{aligned}
$$

using Theorem 1 of Polson et al. (2013), where $p_b(\xi)$ is the density of the random variable $\xi \sim \text{PG}(b, 0), b > 0$. It follows that

$$
[\eta] \propto \int_0^\infty \exp\left\{-\frac{1}{2}\big[\eta^2 \xi - \eta(\alpha - \beta)\big]\right\} p_{\alpha+\beta}(\xi) \, d\xi \propto \int_0^\infty f_N\big(\eta; \xi^{-1}[\alpha - \beta]/2, \xi^{-1}\big) p_{\alpha+\beta}(\xi) \, d\xi
$$

where $f_N(\eta; \mu_N, \sigma_N^2)$ is the density of the random variable $\eta \sim N(\mu_N, \sigma_N^2)$.

The conditional distribution $[\xi|\eta] \sim \text{PG}(\alpha + \beta, \eta)$ is a direct result of Polson et al. (2013).

$\qquad\square$

# B   MCMC Sampling Algorithm and Computational Details

We design a Gibbs sampling algorithm for the dynamic shrinkage process. The sampling algorithm is both computationally and MCMC efficient, and builds upon two main components: (1) a stochastic volatility sampling algorithm (Kastner and Frühwirth-Schnatter, 2014) augmented with a Pólya-Gamma sampler (Polson et al., 2013); and (2) a Cholesky

Factor Algorithm (CFA, Rue, 2001) for sampling the state variables in the dynamic linear model. Alternative sampling algorithms exist for more general DLMs, such as the simulation smoothing algorithm of Durbin and Koopman (2002). However, as demonstrated by McCausland et al. (2011) and explored in Chan and Jeliazkov (2009) and Chan (2013), the CFA sampler is often more efficient. Importantly, both components employ algorithms that are linear in the number of time points, which produces a highly efficient sampling algorithm.

The general sampling algorithm is as follows, with the details provided in the subsequent sections:

1. Sample the dynamic shrinkage components (Section 5.1)

   (a) Log-volatilities, $\{h_t\}$

   (b) Pólya-Gamma mixing parameters, $\{\xi_t\}$

   (c) Unconditional mean of log-volatility, $\mu$

   (d) AR(1) coefficient of log-volatility, $\phi$

   (e) Discrete mixture component indicators, $\{s_t\}$

2. Sample the state variables, $\{\boldsymbol{\beta}_t\}$ (Section B.2)

3. Sample the observation error variance, $\sigma_\epsilon^2$.

For the observation error variance, we follow Carvalho et al. (2010) and assume the Jeffreys' prior $[\sigma_\epsilon^2] \propto 1/\sigma_\epsilon^2$. The full conditional distribution is $[\sigma_\epsilon | \{y_t\}_{t=1}^T, \{\beta_t\}_{t=1}^T, \tau^2] \propto \sigma_\epsilon^{-1} \times \sigma_\epsilon^{-T} \exp\left\{ -\frac{1}{2\sigma_\epsilon^2} \sum_{t=1}^T (y_t - \beta_t)^2 \right\} \times \frac{\sqrt{T}}{\sigma_\epsilon(1 + T\tau^2/\sigma_\epsilon^2)}$, where the last term arises from $\tau \sim C^+(0, \sigma_\epsilon/\sqrt{T})$. We sample from this distribution using the slice sampler of Neal (2003).

If we instead use a stochastic volatility model for the observation error variance as in Sections 3.2 and 4.2, we replace this step with a stochastic volatility sampling algorithm (e.g., Kastner and Frühwirth-Schnatter, 2014), which requires additional sampling steps for the corresponding log-volatility and the unconditional mean, AR(1) coefficient, and evolution error variance of log-volatility. An efficient implementation of such a sampler is available in

43

the R package `stochvol` (Kastner, 2016). In this setting, we do not scale $\tau$ by the standard deviation, and instead assume $\tau \sim C^+(0, 1/\sqrt{T})$.

In Figure 10, we provide empirical evidence for the linear time $\mathcal{O}(T)$ computations of the Bayesian trend filtering model with dynamic horseshoe innovations. The runtime per 1000 MCMC iterations is less than 6 minutes (on a MacBook Pro, 2.7 GHz Intel Core i5) for samples sizes up to $T = 10^5$, so the Gibbs sampling algorithm is scalable.

**Computation Time for BTF-DHS (per 1000 MCMC iterations)**



Figure 10: Computation time per 1000 MCMC iterations for the Bayesian trend filtering model with dynamic horseshoe innovations (BTF-DHS).

## B.1 Efficient Sampling for the Dynamic Shrinkage Process

Consider the (univariate) dynamic shrinkage process in (4) with the Pólya-Gamma parameter expansion of Theorem 4. We provide implementation details for the dynamic horseshoe prior with $\alpha = \beta = 1/2$, but extensions to other cases are straightforward. The SV sampling

framework of Kastner and Frühwirth-Schnatter (2014) represents the likelihood for $h_t$ on the log-scale, and approximates the ensuing $\log \chi_1^2$ distribution for the errors via a known discrete mixture of Gaussian distributions. In particular, let $\tilde{y}_t = \log(\omega_t^2 + c)$, where $c$ is a small offset to avoid numerical issues. Conditional on the mixture component indicators $s_t$, the likelihood is $\tilde{y}_t \overset{indep}{\sim} N(h_t + m_{s_t}, v_{s_t})$ where $m_i$ and $v_i, i = 1, \ldots, 10$ are the pre-specified mean and variance components of the 10-component Gaussian mixture provided in Omori et al. (2007). The evolution equation is $h_{t+1} = \mu + \phi(h_t - \mu) + \eta_t$ with initialization $h_1 = \mu + \eta_0$ and innovations $[\eta_t | \xi_t] \overset{indep}{\sim} N(0, \xi_t^{-1})$ for $[\xi_t] \overset{iid}{\sim} \text{PG}(1, 0)$.

To sample $\boldsymbol{h} = (h_1, \ldots, h_T)$ jointly, we directly compute the posterior distribution of $\boldsymbol{h}$ and exploit the tridiagonal structure of the resulting posterior precision matrix. In particular, we equivalently have $\tilde{\boldsymbol{y}} \sim N(\boldsymbol{m} + \tilde{\boldsymbol{h}} + \tilde{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_v)$ and $\boldsymbol{D}_\phi \tilde{\boldsymbol{h}} \sim N(\boldsymbol{0}, \boldsymbol{\Sigma}_\xi)$, where $\boldsymbol{m} = (m_{s_1}, \ldots, m_{s_T})'$, $\tilde{\boldsymbol{h}} = (h_1 - \mu, \ldots, h_T - \mu)'$, $\tilde{\boldsymbol{\mu}} = (\mu, (1-\phi)\mu, \ldots, (1-\phi)\mu)'$, $\boldsymbol{\Sigma}_v = \text{diag}\left(\{v_{s_t}\}_{t=1}^T\right)$, $\boldsymbol{\Sigma}_\xi = \text{diag}\left(\{\xi_t^{-1}\}_{t=1}^T\right)$, and $\boldsymbol{D}_\phi$ is a lower triangular matrix with ones on the diagonal, $-\phi$ on the first off-diagonal, and zeros elsewhere. We sample from the posterior distribution of $\boldsymbol{h}$ by sampling from the posterior distribution of $\tilde{\boldsymbol{h}}$ and setting $\boldsymbol{h} = \tilde{\boldsymbol{h}} + \mu \boldsymbol{1}$ for $\boldsymbol{1}$ a $T$-dimensional vector of ones. The required posterior distribution is $\tilde{\boldsymbol{h}} \sim N\left(\boldsymbol{Q}_{\tilde{h}}^{-1} \boldsymbol{\ell}_{\tilde{h}}, \boldsymbol{Q}_{\tilde{h}}^{-1}\right)$, where $\boldsymbol{Q}_{\tilde{h}} = \boldsymbol{\Sigma}_v^{-1} + \boldsymbol{D}_\phi' \boldsymbol{\Sigma}_\xi^{-1} \boldsymbol{D}_\phi$ is a tridiagonal symmetric matrix with diagonal elements $\boldsymbol{d}_0(\boldsymbol{Q}_{\tilde{h}})$ and first off-diagonal elements $\boldsymbol{d}_1(\boldsymbol{Q}_{\tilde{h}})$ defined as

$$\boldsymbol{d}_0(\boldsymbol{Q}_{\tilde{h}}) = \left[(v_{s_1}^{-1} + \xi_1 + \phi^2 \xi_2), (v_{s_2}^{-1} + \xi_2 + \phi^2 \xi_3), \ldots, (v_{s_{T-1}}^{-1} + \xi_{T-1} + \phi^2 \xi_T), (v_{s_T}^{-1} + \xi_T)\right],$$

$$\boldsymbol{d}_1(\boldsymbol{Q}_{\tilde{h}}) = \left[(-\phi\xi_2), (-\phi\xi_3), \ldots, (-\phi\xi_{T-1})\right], \text{ and}$$

$$\boldsymbol{\ell}_{\tilde{h}} = \boldsymbol{\Sigma}_v^{-1}\left(\tilde{\boldsymbol{y}} - \boldsymbol{m} - \tilde{\boldsymbol{\mu}}\right)$$
$$= \left[\frac{\tilde{y}_1 - m_{s_1} - \mu}{v_{s_1}}, \frac{\tilde{y}_2 - m_{s_2} - (1-\phi)\mu}{v_{s_2}}, \ldots, \frac{\tilde{y}_T - m_{s_T} - (1-\phi)\mu}{v_{s_T}}\right]'.$$

Drawing from this posterior distribution is straightforward and efficient, using band back-substitution described in Kastner and Frühwirth-Schnatter (2014): (1) compute the Cholesky decomposition $\boldsymbol{Q}_{\tilde{h}} = \boldsymbol{L}\boldsymbol{L}'$, where $\boldsymbol{L}$ is lower triangle; (2) solve $\boldsymbol{L}\boldsymbol{a} = \boldsymbol{\ell}_{\tilde{h}}$ for $\boldsymbol{a}$; and (3) solve

$L'\tilde{h} = a + e$ for $\tilde{h}$, where $e \sim N(0, I_T)$.

Conditional on the log-volatilities $\{h_t\}$, we sample the AR(1) evolution parameters: the log-innovation precisions $\{\xi_t\}$, the autoregressive coefficient $\phi$, and the unconditional mean $\mu$. The precisions are distributed $[\xi_t|\eta_t] \sim \text{PG}(1, \eta_t)$ for $\eta_t = h_{t+1} - \mu - \phi(h_t - \mu)$, which we sample using the `rpg()` function in the `R` package `BayesLogit` (Polson et al., 2013). The Pólya-Gamma sampler is efficient: using only exponential and inverse-Gaussian draws, Polson et al. (2013) construct an accept-reject sampler for which the probability of acceptance is uniformly bounded below at 0.99919, which does not require any tuning. Next, we assume the prior $[(\phi + 1)/2] \sim \text{Beta}(a_\phi, b_\phi)$, which restricts $|\phi| < 1$ for stationarity, and sample from the full conditional distribution of $\phi$ using the slice sampler of Neal (2003). We select $a_\phi = 10$ and $b_\phi = 2$, which places most of the mass for the density of $\phi$ in $(0, 1)$ with a prior mean of $2/3$ and a prior mode of $4/5$ to reflect the likely presence of persistent volatility clustering. The prior for the global scale parameter is $\tau \sim C^+(0, \sigma_\epsilon/\sqrt{T})$, which implies $\mu = \log(\tau^2)$ is $[\mu|\sigma_\epsilon, \xi_\mu] \sim N(\log(\sigma_\epsilon^2/T), \xi_\mu^{-1})$ with $\xi_\mu \sim \text{PG}(1, 0)$. Including the initialization $h_1 \sim N(\mu, \xi_0^{-1})$ with $\xi_0 \sim \text{PG}(1, 0)$, the posterior distribution for $\mu$ is $\mu \sim N(Q_\mu^{-1}\ell_\mu, Q_\mu^{-1})$ with $Q_\mu = \xi_\mu + \xi_0 + (1-\phi)^2 \sum_{t=1}^{T-1} \xi_t$ and $\ell_\mu = \xi_\mu \log(\sigma_\epsilon^2/T) + \xi_0 h_1 + (1-\phi)\sum_{t=1}^{T-1}\xi_t(h_{t+1} - \phi h_t)$. Sampling $\xi_\mu$ and $\xi_0$ follows the Pólya-Gamma sampling scheme above.

Finally, we sample the discrete mixture component indicators $s_t$. The discrete mixture probabilities are straightforward to compute: the prior mixture probabilities are the mixing proportions given by Omori et al. (2007) and the likelihood is $\tilde{y}_t \overset{indep}{\sim} N(h_t + m_{s_t}, v_{s_t})$; see Kastner and Frühwirth-Schnatter (2014) for details.

In the multivariate setting $p > 1$ of (11) with $\mathbf{\Phi} = \text{diag}(\phi_1, \ldots, \phi_p)$, we may modify the log-volatility sampler of $\{h_{j,t}\}$ by redefining relevant quantities using the ordering $\boldsymbol{h} = (h_{1,1}, \ldots, h_{1,T}, h_{2,1}, \ldots, h_{p,T})'$. In particular, the posterior precision matrix is again tridiagonal, but with diagonal elements $d_0(\boldsymbol{Q}_{\tilde{h}}) = [d_{0,1}(\boldsymbol{Q}_{\tilde{h}}), \ldots, d_{0,p}(\boldsymbol{Q}_{\tilde{h}})]$ and first off-diagonal elements $\boldsymbol{d}_1(\boldsymbol{Q}_{\tilde{h}}) = [d_{1,1}(\boldsymbol{Q}_{\tilde{h}}), 0, d_{1,2}(\boldsymbol{Q}_{\tilde{h}}), 0, \ldots, 0, d_{1,p}(\boldsymbol{Q}_{\tilde{h}})]$, where $d_{0,j}(\boldsymbol{Q}_{\tilde{h}})$ and $d_{1,j}(\boldsymbol{Q}_{\tilde{h}})$ are the diagonal elements and first off-diagonal elements, respectively, for predictor $j$ as com-

puted in the univariate case above. Similarly, the linear term $\boldsymbol{\ell}_{\tilde{h}} = \left[\boldsymbol{\ell}'_{\tilde{h},1}, \ldots, \boldsymbol{\ell}'_{\tilde{h},p}\right]'$ where $\boldsymbol{\ell}_{\tilde{h},j}$ is the linear term for predictor $j$ as computed in the univariate case. The parameters $\xi_{j,t}$, $\phi_j$, and $s_{j,t}$ may be sampled independently as in the univariate case, while samplers for $\{\mu_j\}$ and $\mu_0$ proceed as in a standard hierarchical Gaussian model. For the more general case of non-diagonal $\boldsymbol{\Phi}$, we may use a simulation smoothing algorithm (e.g., Durbin and Koopman, 2002) for the log-volatilities $\{h_{j,t}\}$, while the sampler for $\boldsymbol{\Phi}$ will depend on the chosen prior.

## B.2 Efficient Sampling for the State Variables

In the univariate setting of (9), the sampler for $\boldsymbol{\beta} = (\beta_1, \ldots, \beta_T)$ is similar to the log-volatility sample in Section 5.1. We provide the details for $D = 2$; the $D = 1$ case is similar to Section 5.1 with $\phi = 1$, $\mu = 0$, and $m_{s_t} = 0$. Model (9) may be written $\boldsymbol{y} \sim N(\boldsymbol{\beta}, \boldsymbol{\Sigma}_\epsilon)$ and $\boldsymbol{D}_2\boldsymbol{\beta} \sim N(\boldsymbol{0}, \boldsymbol{\Sigma}_\omega)$, where $\boldsymbol{y} = (y_1, \ldots, y_T)'$, $\boldsymbol{\Sigma}_\epsilon = \text{diag}\left(\{\sigma_t^2\}_{t=1}^T\right)$, $\boldsymbol{\Sigma}_\omega = \text{diag}\left(\{\sigma_{\omega_t}^2\}_{t=1}^T\right)$ for $\sigma_{\omega_t}^2 = \tau^2\lambda_t^2$, and $\boldsymbol{D}_2$ is a lower triangular matrix with ones on the diagonal, $(0, -2, \ldots, -2)$ on the first off-diagonal, ones on the second off-diagonal, and zeros elsewhere. Note that we allow the observation error variance $\sigma_t^2$ to be time-dependent for full generality, as in Section 4.2. The posterior for $\boldsymbol{\beta}$ is $\boldsymbol{\beta} \sim N\left(\boldsymbol{Q}_{\boldsymbol{\beta}}^{-1}\boldsymbol{\ell}_{\boldsymbol{\beta}}, \boldsymbol{Q}_{\boldsymbol{\beta}}^{-1}\right)$, where $\boldsymbol{Q}_{\boldsymbol{\beta}} = \boldsymbol{\Sigma}_\epsilon^{-1} + \boldsymbol{D}_2'\boldsymbol{\Sigma}_\omega^{-1}\boldsymbol{D}_2$ is a pentadiagonal symmetric matrix with diagonal elements $\boldsymbol{d}_0(\boldsymbol{Q}_{\boldsymbol{\beta}})$, first off-diagonal elements $\boldsymbol{d}_1(\boldsymbol{Q}_{\boldsymbol{\beta}})$, and second-off diagonal elements $\boldsymbol{d}_2(\boldsymbol{Q}_{\boldsymbol{\beta}})$ defined as

$$
\begin{aligned}
\boldsymbol{d}_0(\boldsymbol{Q}_{\boldsymbol{\beta}}) = \Big[ &\left(\sigma_1^{-2} + \sigma_{\omega_1}^{-2} + \sigma_{\omega_3}^{-2}\right), \left(\sigma_2^{-2} + \sigma_{\omega_2}^{-2} + 4\sigma_{\omega_3}^{-2} + \sigma_{\omega_4}^{-2}\right), \ldots, \\
&\left(\sigma_t^{-2} + \sigma_{\omega_t}^{-2} + 4\sigma_{\omega_{t+1}}^{-2} + \sigma_{\omega_{t+2}}^{-2}\right), \ldots, \\
&\left(\sigma_{T-2}^{-2} + \sigma_{\omega_{T-2}}^{-2} + 4\sigma_{\omega_{T-1}}^{-2} + \sigma_{\omega_T}^{-2}\right), \left(\sigma_{T-1}^{-2} + \sigma_{\omega_{T-1}}^{-2} + 4\sigma_{\omega_T}^{-2}\right), \left(\sigma_T^{-2} + \sigma_{\omega_T}^{-2}\right) \Big], \\
\boldsymbol{d}_1(\boldsymbol{Q}_{\boldsymbol{\beta}}) = \big[ &-2\sigma_{\omega_3}^{-2}, -2\left(\sigma_{\omega_3}^{-2} + \sigma_{\omega_4}^{-2}\right), \ldots, -2\left(\sigma_{\omega_t}^{-2} + \sigma_{\omega_{t+1}}^{-2}\right), \ldots, -2\left(\sigma_{\omega_{T-1}}^{-2} + \sigma_{\omega_T}^{-2}\right), -2\sigma_{\omega_T}^{-2}\big], \\
\boldsymbol{d}_2(\boldsymbol{Q}_{\boldsymbol{\beta}}) = \big[ &\sigma_{\omega_3}^{-2}, \ldots, \sigma_{\omega_t}^{-2}, \ldots, \sigma_{\omega_T}^{-2}\big],
\end{aligned}
$$

and $\boldsymbol{\ell_\beta} = \boldsymbol{\Sigma}_\epsilon^{-1} \boldsymbol{y} = \left[ y_1/\sigma_1^2, \ldots, y_t/\sigma_t^2, \ldots, y_T/\sigma_T^2 \right]'$. Drawing from the posterior distribution is straightforward and efficient using the back-band substitution algorithm described in Section 5.1.

In the multivariate setting of (10), we similarly derive the posterior distribution for $\boldsymbol{\beta} = (\boldsymbol{\beta}_1', \ldots, \boldsymbol{\beta}_T')' = (\beta_{1,1}, \beta_{2,1}, \ldots, \beta_{p,1}, \beta_{1,2}, \ldots, \beta_{p,T})'$. Let $\boldsymbol{X} = \text{blockdiag}\left( \{\boldsymbol{x}_t'\}_{t=1}^T \right)$ denote the $T \times Tp$ block-diagonal matrix of predictors and $\boldsymbol{\Sigma}_\omega = \text{diag}\left( \{\sigma_{\omega_{j,t}}^2\}_{j,t} \right)$ for $\sigma_{\omega_{j,t}}^2 = \tau_0^2 \tau_j^2 \lambda_{j,t}^2$. The posterior distribution is $\boldsymbol{\beta} \sim N\left( \boldsymbol{Q}_\beta^{-1} \boldsymbol{\ell}_\beta, \boldsymbol{Q}_\beta^{-1} \right)$, where

$$\boldsymbol{Q}_\beta = \boldsymbol{X}' \boldsymbol{\Sigma}_\epsilon^{-1} \boldsymbol{X} + (\boldsymbol{D}_2' \otimes \boldsymbol{I}_p) \boldsymbol{\Sigma}_\omega^{-1} (\boldsymbol{D}_2 \otimes \boldsymbol{I}_p)$$

and

$$\boldsymbol{\ell}_\beta = \boldsymbol{X}' \boldsymbol{\Sigma}_\epsilon^{-1} \boldsymbol{y} = \left[ \boldsymbol{x}_1' y_1/\sigma_1^2, \ldots, \boldsymbol{x}_t' y_t/\sigma_t^2, \ldots, \boldsymbol{x}_T' y_T/\sigma_T^2 \right]'.$$

Note that $\boldsymbol{Q}_\beta$ may be constructed directly as above, but is now $2p$-banded. Alternatively, the regression coefficients $\{\beta_{j,t}\}$ may be sampled jointly using the simulation smoothing algorithm of Durbin and Koopman (2002).

# C   Additional Simulation Results

We augment the simulation study of Section 4 by considering autocorrelated predictors. Following the simulation design from Section 4, we simulated 100 data sets of length $T = 200$ from the model $y_t = \boldsymbol{x}_t' \boldsymbol{\beta}_t^* + \epsilon_t$ with $\epsilon_t \overset{iid}{\sim} N(0, \sigma_*^2)$. The $p = 20$ predictors are $x_{1,t} = 1$ and, for $j = 2, \ldots, p$, the time series $\{x_{j,t}\}_{t=1}^T$ is simulated from an AR(1) process with an autoregressive coefficient of 0.8, Gaussian innovations, unconditional mean zero, and unconditional standard deviation one. The true regression coefficients $\boldsymbol{\beta}_t^* = (\beta_{1,t}^*, \ldots, \beta_{p,t}^*)'$ are the following: $\beta_{1,t}^* = 2$ is constant; $\beta_{2,t}^*$ is piecewise constant with $\beta_{2,t}^* = 0$ everywhere except $\beta_{2,t}^* = 2$ for $t = 41, \ldots, 80$ and $\beta_{2,t}^* = -2$ for $t = 121, \ldots, 160$; $\beta_{3,t}^* = \frac{1}{\sqrt{100}} \sum_{s=1}^t Z_s$ with $Z_s \overset{iid}{\sim} N(0,1)$ is a scaled random walk for $t \le 100$ and $\beta_{3,t}^* = 0$ for $t > 100$; and $\beta_{j,t}^* = 0$

for $j = 4, \ldots, p = 20$. The predictor set contains a variety of functions: a constant nonzero function, a locally constant function, a slowly-varying function that thresholds to zero for $t > 100$, and 17 true zeros. The noise variance $\sigma_*^2$ is determined by selecting a root-signal-to-noise ratio (RSNR) and computing $\sigma_* = \sqrt{\frac{\sum_{t=1}^{T}(y_t^* - \bar{y}^*)^2}{T-1}} \big/ \mathrm{RSNR}$, where $y_t^* = \boldsymbol{x}_t' \boldsymbol{\beta}_t^*$ and $\bar{y}^* = \frac{1}{T} \sum_{t=1}^{T} y_t^*$. We select RSNR = 3.

We evaluate competing methods using RMSEs for both $y_t^*$ and $\boldsymbol{\beta}_t^*$ defined by $\mathrm{RMSE}(\hat{y}) = \sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^* - \hat{y}_t)^2}$ and $\mathrm{RMSE}(\hat{\boldsymbol{\beta}}) = \sqrt{\frac{1}{Tp} \sum_{t=1}^{T} \sum_{j=1}^{p} \left( \beta_{j,t}^* - \hat{\beta}_{j,t} \right)^2}$ for all estimators $\hat{\boldsymbol{\beta}}_t$ of the true regression functions, $\boldsymbol{\beta}_t^*$ with $\hat{y}_t = \boldsymbol{x}_t' \hat{\boldsymbol{\beta}}_t$. The results are displayed in Figure 11. The proposed BTF-DHS model outperforms the competitors in both recovery of the true regression functions, $\beta_{j,t}^*$ and estimation of the true curves, $y_t^*$.
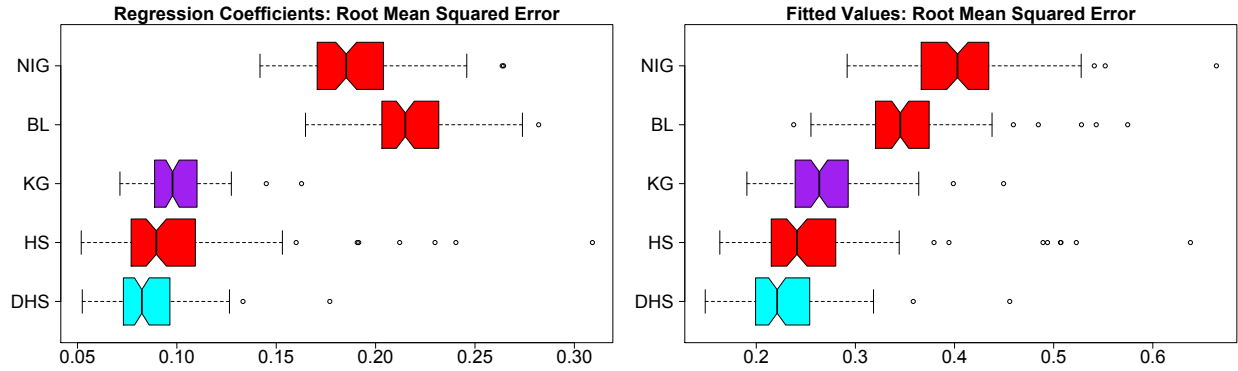


Figure 11: Root mean squared errors for the regression coefficients, $\beta_{j,t}^*$ (**left**) and the true curves, $y_t^* = \boldsymbol{x}_t' \boldsymbol{\beta}_t^*$ (**right**) for simulated data. Non-overlapping notches indicate significant differences between medians.