

# Factor MIDAS for Nowcasting and Forecasting with Ragged-Edge Data: A Model Comparison for German GDP\*

MASSIMILIANO MARCELLINO<sup>†</sup> and CHRISTIAN SCHUMACHER<sup>‡</sup>

<sup>†</sup>*Economics Department, European University Institute, Villa San Paolo, Florence and Bocconi University, Milan, Italy (e-mail: massimiliano.marcellino@eui.eu)*

<sup>‡</sup>*Deutsche Bundesbank, Economics Research Centre, Wilhelm-Epstein-Straße 14, 60431 Frankfurt/Main, Germany (e-mail: christian.schumacher@bundesbank.de)*

## Abstract

In this article, we merge two strands from the recent econometric literature. First, factor models based on large sets of macroeconomic variables for forecasting, which have generally proven useful for forecasting. However, there is some disagreement in the literature as to the appropriate method. Second, forecast methods based on mixed-frequency data sampling (MIDAS). This regression technique can take into account unbalanced datasets that emerge from publication lags of high- and low-frequency indicators, a problem practitioners have to cope with in real time. In this article, we introduce Factor MIDAS, an approach for nowcasting and forecasting low-frequency variables like gross domestic product (GDP) exploiting information in a large set of higher-frequency indicators. We consider three alternative MIDAS approaches (basic, smoothed and unrestricted) that provide harmonized projection

\*The authors are grateful for helpful comments and discussions to three anonymous referees, Riccardo Cristadoro, Sandra Eickmeier, Petra Gerlach-Kristen, Malte Knüppel, Gerhard Rünstler, Karsten Ruth, Klaus Wohlrabe and participants at the Bank of England CCBS research forum 'New Developments in Dynamic Factor Modelling' 2007, the Workshop 'Forecasting Short-Term Developments and the Role of Econometric Models' at the Bank of Canada 2007, the Pfingsttagung of the DStatG 2008, the Annual Conference of the Verein für Socialpolitik 2008, the Joint Research Workshop OeNB-SNB-Bbk 2008, the Workshop 'Forecasting Macroeconomic Variables Using Dynamic Factor Models' at the Banque de France 2008, the CFE Workshop 2008 in Neuchâtel and a seminar at the Bundesbank. The codes for this article were written in Matlab. Some functions were taken from the Econometrics Toolbox written by James P. LeSage from <http://www.spatial-econometrics.com>. Other codes were kindly provided by Mario Forni from [http://www.economia.unimore.it/forni\\_mario/matlab.htm](http://www.economia.unimore.it/forni_mario/matlab.htm), Arthur Sinko from [www.unc.edu/~sinko/midas.zip](http://www.unc.edu/~sinko/midas.zip) and Gerhard Rünstler.

JEL Classification numbers: E37, C53.

methods that allow for a comparison of the alternative factor estimation methods with respect to nowcasting and forecasting. Common to all the factor estimation methods employed here is that they can handle unbalanced datasets, as typically faced in real-time forecast applications owing to publication lags. In particular, we focus on variants of static and dynamic principal components as well as Kalman filter estimates in state-space factor models. As an empirical illustration of the technique, we use a large monthly dataset of the German economy to nowcast and forecast quarterly GDP growth. We find that the factor estimation methods do not differ substantially, whereas the most parsimonious MIDAS projection performs best overall. Finally, quarterly models are in general outperformed by the Factor MIDAS models, which confirms the usefulness of the mixed-frequency techniques that can exploit timely information from business cycle indicators.

## I. Introduction

The use of large sets of macroeconomic variables for forecasting has received increased attention in the recent literature. In particular, different types of large factor models have been widely discussed; see, for example, the comparisons and surveys by Boivin and Ng (2005), Stock and Watson (2006), D'Agostino and Giannone (2006), Eickmeier and Ziegler (2008) and Schumacher (2007). Another strand of the recent literature has considered modelling and forecasting with mixed-frequency data, in particular the mixed-data sampling (MIDAS) approach, as introduced by Ghysels, Sinko and Valkanov (2007). This approach allows for forecasting a low-frequency variable like quarterly gross domestic product (GDP) with a small number of high-frequency indicators, and this has been introduced to the macroeconomic forecast literature by Clements and Galvão (2008, 2009); see also Ghysels and Wright (2009). In this article, we bring together these two recent strands of the literature, and introduce Factor MIDAS, an approach for nowcasting and forecasting low-frequency variables exploiting information in a large set of higher-frequency indicators.

The basic MIDAS framework consists of a regression of a low-frequency variable on a set of higher-frequency indicators, where distributed lag functions are employed to specify the dynamic relationship. The Factor MIDAS approach exploits estimated factors rather than single or small groups of economic indicators as regressors. Therefore, it directly translates the well-known two-step approach of factor forecasting as introduced by Stock and Watson (2002) and Bai and Ng (2006) for the single-frequency case to the mixed-frequency case, where factors are sampled at higher frequencies than the variable to be predicted. As in the standard MIDAS case, see Clements and Galvão (2008), Factor MIDAS is a tool for direct multi-step nowcasting and forecasting; see Marcellino, Stock and Watson (2006). As a modification to MIDAS, we follow Koenig, Dolmas and Piger (2003) and also evaluate a more general regression approach, where the dynamic relationship between the low-frequency variable and the high-frequency indicators – factors in our case – is unrestricted, in contrast to the distributed lag functions as proposed by Ghysels *et al.*

(2007). This approach is labelled as unrestricted Factor MIDAS. As a third alternative, we consider a special regression scheme proposed by Altissimo *et al.* (2006) that considers only certain frequencies of correlation between variables sampled at high and low frequencies. As the regression essentially eliminates high-frequency correlations, we call it smoothed MIDAS. To motivate and compare these variants of Factor MIDAS, we discuss time aggregation in a theoretical high-frequency factor model with dynamic factors from Boivin and Ng (2005) and derive how mixed-data sampling can approximate the forecasts from the theoretical model.

Alternative factor estimators can be used for nowcasting and forecasting in the Factor MIDAS framework, and as an additional novelty of the article, we compare three different factor estimation methods that have been developed in the recent literature, but have not been compared with each other so far in a standardized framework. The factor estimation methods we discuss have in common that they can account for unbalanced datasets. In empirical real-time applications, multivariate datasets are typically unbalanced because of non-synchronous publication dates and different publication delays of the economic indicators. This leads to the so-called ragged edge of the data, see Wallis (1986). We focus on three factor estimators that are all suited for this case. First, the factor estimator by Altissimo *et al.* (2006), which builds upon the one-sided non-parametric dynamic principal component analysis (DPCA) factor estimator of Forni *et al.* (2005). To take into account the ragged edge of the data, Altissimo *et al.* (2006) simply apply a realignment of each time series to obtain a balanced dataset. Second, the expectation-maximization (EM) algorithm combined with the factor estimator-based static principal component analysis (PCA), as introduced by Stock and Watson (2002) and applied for forecasting and interpolation by Bernanke and Boivin (2003), Angelini, Henry and Marcellino (2006) and Schumacher and Breitung (2008). Third, the two-step parametric state-space factor estimator based on the Kalman smoother of Doz, Giannone and Reichlin (2006), as applied in Giannone, Reichlin and Small (2008), Angelini *et al.* (2008), Matheson (2007) and Aastveit and Trovik (2007). Which one of the estimation performs best, is *a priori* unclear. There is a large literature on comparing factor estimation methods based on large and balanced datasets, see Boivin and Ng (2005), Stock and Watson (2006), D'Agostino and Giannone (2006), Schumacher (2007), Kapetanios and Marcellino (2009) and we extend that literature by taking into account estimation methods suited for unbalanced datasets.

Combining the three alternative MIDAS regressions with the three competing factor estimators, we have a total of nine Factor MIDAS approaches. To illustrate their implementation and assess their relative merits, we conduct a detailed empirical analysis. In particular, as policy makers regularly request information on the current state of the economy in terms of GDP, which is only available on a quarterly basis and often with substantial delays, we consider Factor MIDAS nowcasting and short-term forecasting of quarterly GDP growth with a large set of timely monthly economic indicators. We focus on Germany, the largest economy in the euro area, where GDP is released about 5–6 weeks after the end of the reference quarter.

Related to the existing literature, existing nowcasts and forecasts for German GDP growth are not fully satisfying; see, for example, Schumacher and Breitung (2008), where nowcasting with a large factor model outperforms naive benchmarks only with horizons up to one quarter ahead. A more recent paper that also contains results on nowcasts and short-term forecasts for German GDP growth is Barhoumi *et al.* (2008), where large-scale factor models can outperform simple benchmarks only for the nowcast. Compared with Barhoumi *et al.* (2008), first we introduce Factor MIDAS as a standardized tool for comparing factor estimation methods, whereas Barhoumi *et al.* (2008) use different projection methods for different factor estimation methods, which makes it difficult to disentangle their relative contribution to the forecasting performance. Second, we evaluate a larger set of factor estimation methods, in particular the DPCA estimator based on unbalanced data by Altissimo *et al.* (2006), whereas Barhoumi *et al.* (2008) only employ time-aggregated quarterly data for factor estimation based on DPCA.

To relate Factor MIDAS to the methods from the existing literature, we compare two additional approaches in the empirical exercise. First, single-frequency factor models based on quarterly time-aggregated data. Quarterly data have been often used to forecast German GDP growth; see, for example, Schumacher (2007) and Barhoumi *et al.* (2008) and for other countries and datasets, for example, by Marcellino, Stock and Watson (2005) for euro area countries' GDP using disaggregated and aggregated data, Banerjee, Marcellino and Masten (2005) for euro area GDP, Banerjee and Marcellino (2006) for US GDP, Kapetanios, Labhard and Price (2008) for UK GDP and many others. Second, we compare the two-step Factor MIDAS approach with the integrated state-space approach by Banbura and Rünstler (2007) based on Doz *et al.* (2006), where nowcasting and forecasting as well as estimation of the factors is carried out simultaneously by the Kalman smoother.

Our empirical results show that the factor estimation methods do not differ substantially, whereas the most parsimonious unrestricted MIDAS projection performs best overall. Moreover, Factor MIDAS performs similarly compared with the integrated state-space approach by Banbura and Rünstler (2007). Finally, quarterly models and benchmarks are in general outperformed, highlighting the usefulness of the Factor MIDAS approach for empirical analysis.

The article proceeds as follows. Section II introduces the different Factor MIDAS approaches. Section III discusses a simple theoretical motivation of Factor MIDAS. Section IV reviews the competing factor estimation methods. Section V presents the design of the empirical forecast exercise. Section VI reports and discusses the results. Section VII summarizes and concludes.

## II. Factor MIDAS

Let us consider the variable  $y_{t_q}$ , where  $t_q$  is the quarterly time index  $t_q = 1, 2, 3, \dots, T_q$ . The same variable can also be expressed at the monthly frequency by setting  $y_{t_m} =$

$y_{t_q} \forall t_m = 3t_q$ , with  $t_m$  as the monthly time index. Thus,  $y_{t_m}$  is observed only at months  $t_m = 3, 6, 9, \dots, T_m$  with  $T_m = 3T_q$ .

We are interested in forecasting the variable  $y$ , in our case quarterly GDP growth,  $h_q$  quarters ahead or  $h_m = 3h_q$  months ahead, based on all the available monthly information. As the final GDP observation is available in period  $T_q$ , we aim at the forecast of the value  $y_{T_q+h_q} = y_{T_m+h_m}$ . For example, as GDP for the first quarter of a given year is released around mid-May, a nowcast can be produced in January, February and March of the current year, whereas a forecast can be produced in any month of the previous year.

The information set includes a large set of stationary monthly indicators, collected in the  $N$ -dimensional vector  $\mathbf{X}_{t_m}$ . The time index  $t_m$  denotes monthly frequency, and we allow for observations in the months  $t_m = 1, 2, 3, \dots, T_m + w$ . As  $T_m$  is the final month in the quarter, where GDP is available, we allow for at most  $w > 0$  monthly values of the indicators, which are earlier available than GDP, for example, from surveys or financial indicators. However, owing to publication lags, some observations for certain time series at the end of the sample can be missing, thus rendering an unbalanced sample of  $\mathbf{X}_{t_m}$ . For nowcasting and forecasting, we want to exploit all the information available, and thus condition on information up to period  $T_m + w$ , and we aim at computing the forecast  $y_{T_m+h_m|T_m+w}$  equivalent to  $y_{T_q+h_q|T_m+w}$ .

We want to model  $\mathbf{X}_{t_m}$  using a factor representation, where  $r$  factors  $\mathbf{F}_{t_m}$  are estimated to efficiently summarize the information in the data, taking into account missing observations for parts of the dataset. After having estimated the factors,  $\{\hat{\mathbf{F}}_{t_m}\}_{t_m=1}^{T_m+w}$ , they will be used as predictors in the projection for GDP growth. In this section, we will assume that the estimated factors are just available, and that factor estimation does not create generated regressors problems.<sup>1</sup> The factor estimation methods are explained in detail in section IV.

To forecast GDP using the estimated monthly factors, we rely on the MIDAS approach as proposed by Ghysels *et al.* (2007), Ghysels and Wright (2009) and Clements and Galvão (2008, 2009). The combination of MIDAS and factor estimation methods yields the Factor MIDAS approach. We consider three alternative Factor MIDAS approaches in each of the following subsections. In each case, we adopt direct estimation for the Factor MIDAS regression, namely, we relate future values of GDP to current and lagged indicators, thus yielding different forecast models for each forecast horizon; see Marcellino *et al.* (2006) as well as Chevillon and Hendry (2005) for detailed discussions of this issue in the single-frequency case.

### Basic Factor MIDAS

In the standard MIDAS approach as in Clements and Galvão (2008, 2009), economic variables at higher frequency are used as regressors. In the Factor MIDAS approach,

<sup>1</sup> See Bai and Ng (2006) for precise conditions on the number of variables  $N$  and temporal observations  $T_m$ ; basically, it must be that  $T_m^{1/2}/N$  is  $o(1)$ .

the explanatory variables are estimated factors. Let us assume for simplicity that we have only one factor  $\hat{f}_{t_m}$  for forecasting ( $r = 1$ ). The Factor MIDAS model for forecast horizon  $h_q$  quarters with  $h_q = h_m/3$  is:

$$y_{t_q+h_q} = y_{t_m+h_m} = \beta_0 + \beta_1 b(L_m, \theta) \hat{f}_{t_m+w}^{(3)} + \varepsilon_{t_m+h_m}, \quad (1)$$

where the polynomial  $b(L_m, \theta)$  is the exponential Almon lag with

$$b(L_m, \theta) = \sum_{k=0}^K c(k, \theta) L_m^k, \quad c(k, \theta) = \frac{\exp(\theta_1 k + \theta_2 k^2)}{\sum_{k=0}^K \exp(\theta_1 k + \theta_2 k^2)}, \quad (2)$$

with  $\theta = \{\theta_1, \theta_2\}$ . In the Factor MIDAS approach, the quarterly variable  $y_{t_q+h_q} = y_{t_m+h_m}$  is directly related to the factor  $\hat{f}_{t_m+w}^{(3)}$  and its monthly lags. For estimation of the coefficients,  $\hat{f}_{t_m+w}^{(3)}$  is skip sampled from the monthly factor  $\hat{f}_{t_m+w}$ . More precisely, it is defined as  $\hat{f}_j^{(3)} = \hat{f}_j$ ,  $j = 3 + w, 6 + w, \dots, T_m + w$ , so that every third observation starting from the final one is included in the regressor  $\hat{f}_{t_m+w}^{(3)}$ . However, no monthly information is discarded for the forecast, as lagged monthly observations are included in the regressor  $b(L_m, \theta) \hat{f}_{t_m+w}^{(3)}$  given that the lag operator  $L_m$  operates at the monthly frequency. For example, for quarterly observations of GDP in periods  $t_q + h_q$  and  $t_q + h_q - 1$ , the MIDAS equation (1) with  $k = 2$  implies

$$\begin{aligned} y_{t_q+h_q} &= \beta_0 + \beta_1 c(0, \theta) \hat{f}_{t_m+w} + \beta_1 c(1, \theta) \hat{f}_{t_m+w-1} + \beta_1 c(2, \theta) \hat{f}_{t_m+w-2} + \varepsilon_{t_m+h_m}, \\ y_{t_q+h_q-1} &= \beta_0 + \beta_1 c(0, \theta) \hat{f}_{t_m+w-3} + \beta_1 c(1, \theta) \hat{f}_{t_m+w-1-3} + \beta_1 c(2, \theta) \hat{f}_{t_m+w-2-3} + \varepsilon_{t_m+h_m-3}. \end{aligned}$$

The equations always contain monthly lags of the factors on the right-hand side, so all monthly observations of  $\hat{f}_j$  are covered for  $j = t_m + w - 5, \dots, t_m + w$ ; see Andreou, Ghysels and Kourtellis (2009).

Note that month  $T_m$  corresponds to the final quarter observation of GDP  $T_q = T_m/3$ , whereas the final value of the factor is  $T_m + w$ . Thus, equation (1) exploits the timely information incorporated in the factor estimation. In this way, we take into account that a monthly indicator is typically available within the quarter for which no GDP figure is available yet; see Clements and Galvão (2008, 2009) and Giannone *et al.* (2008).

Given  $\theta_1$  and  $\theta_2$ , the exponential lag function  $b(L_m, \theta)$  provides a parsimonious way to consider monthly lags of the factors as we can allow for large  $K$  to approximate the impulse response function of GDP growth from the factors. To estimate  $\theta$ , we can employ nonlinear least squares (NLS) in a regression of  $y_{t_m}$  onto  $\hat{f}_{t_m+w-h_m}^{(3)}$  and lags, yielding coefficients  $\hat{\theta}_1$ ,  $\hat{\theta}_2$ ,  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . The forecast is given by

$$y_{T_m+h_m|T_m+w} = \hat{\beta}_0 + \hat{\beta}_1 b(L_m, \hat{\theta}) \hat{f}_{T_m+w}. \quad (3)$$

Therefore, the projection is based on the final values of the estimated factors from periods  $T_m + w, T_m + w - 1, \dots$ . For the case of  $r > 1$  with  $\mathbf{F}_{t_m} = (f_{1,t_m}, \dots, f_{r,t_m})'$ , the MIDAS regression generalizes to

$$y_{t_m+h_m} = \beta_0 + \sum_{i=1}^r \beta_{1,i} b_i(L_m, \theta_i) \hat{f}_{i,t_m+w}^{(3)} + \varepsilon_{t_m+h_m}. \quad (4)$$

Here, the parameters  $\theta_i$ , which determine the curvature of the impulse response function, can vary between the different factors. The estimation and forecast is otherwise the same.

As an extension to the previous MIDAS model, Clements and Galvão (2008) add autoregressive (AR) dynamics to the MIDAS specification. In addition to providing a more flexible dynamic specification, this approach is also useful to handle eventual serial correlation in the idiosyncratic components, see the theoretical discussion in section III. The AR Factor MIDAS can be written as:

$$y_{t_m+h_m} = \beta_0 + \lambda y_{t_m} + \sum_{i=1}^r \beta_{1,i} b_i(L_m, \theta_i) (1 - \lambda L_m^3) \hat{f}_{i,t_m+w}^{(3)} + \varepsilon_{t_m+h_m}. \quad (5)$$

To rule out discontinuities of the impulse response function of the factors on the left-hand side variable, Clements and Galvão (2008) impose a common factor restriction on the coefficients of the factors and their lags. The AR coefficient  $\lambda$  can be estimated together with the other coefficients by NLS.

In the following, we drop the prefix ‘Factor’ and denote this approach as ‘MIDAS-basic’, as all our applications are factor-based. The MIDAS specification with AR dynamics will be labelled ‘AR-MIDAS’.

### Smoothed MIDAS

Another way to formulate a mixed-frequency projection is employed in the New Eurocoin Index, see Altissimo *et al.* (2006). New Eurocoin is a composite coincident business cycle indicator of the Euro area economy and can be regarded as a projection of smoothed GDP growth on monthly factors, see Altissimo *et al.* (2006), section 4. Although the methods in that paper primarily aim at deriving a composite coincident indicator, not explicitly nowcasts or forecasts, one can directly generalize them for these purposes. In particular, the projection in Altissimo *et al.* (2006) can be modified to  $h_m$ -step forecasts according to

$$y_{T_m+h_m|T_m+w} = \hat{\mu} + \mathbf{G}\hat{\mathbf{F}}_{T_m+w}, \quad (6)$$

$$\mathbf{G} = \tilde{\Sigma}_{\mathbf{y}\mathbf{F}}(h_m - w) \times \hat{\Sigma}_{\mathbf{F}}^{-1}, \quad (7)$$

where  $\hat{\mu}$  is the sample mean of GDP growth, assuming that the factors have mean zero, and  $\mathbf{G}$  is a projection coefficient matrix.  $\hat{\Sigma}_{\mathbf{F}}$  is the estimated sample covariance of the factors, and  $\tilde{\Sigma}_{\mathbf{y}\mathbf{F}}(k)$  is a particular cross-covariance with  $k$  monthly lags between

GDP growth and the factors. The tilde denotes that  $\tilde{\Sigma}_{yF}(k)$  is not an estimate of the sample cross-covariance between factors and GDP growth, rather a cross-covariance between smoothed GDP growth and factors. The smoothing aspect is introduced into  $\tilde{\Sigma}_{yF}(k)$  as follows. Assume that both the factors and GDP growth are demeaned. Then, let the covariance between  $\hat{F}_{t_m-k}$  and  $y_{t_m}$  be estimated by

$$\hat{\Sigma}_{yF}(k) = \frac{1}{T^* - 1} \sum_{t_m=M+1}^{T_m+w} y_{t_m} \hat{F}_{t_m-k}^{(3)'} \quad (8)$$

where  $T^* = \text{floor}[(T_m + w - (M + 1))/3]$  is the number of observations available to compute the cross-covariances for  $k = -M, \dots, M$  and  $M \geq 3h_q = h_m$ . Note that skip-sampled factors  $\hat{F}_{t_m-k}^{(3)'}$  enter  $\hat{\Sigma}_{yF}(k)$ , as we have only quarterly observations of GDP. Given  $\hat{\Sigma}_{yF}(k)$ , we can estimate the cross-spectral matrix

$$\hat{S}_{yF}(\omega_j) = \sum_{k=-M}^M \left(1 - \frac{|k|}{M+1}\right) \hat{\Sigma}_{yF}(k) e^{-i\omega_j k} \quad (9)$$

at frequencies  $\omega_j = \frac{2\pi j}{2H}$  for  $j = -H, \dots, H$  using a Bartlett lead-lag window. The low-frequency relationship between  $\hat{F}_{t_m-k}$  and  $y_{t_m}$  in New Eurocoin is obtained by filtering out cross fluctuations at frequencies larger than  $\pi/6$ , using the frequency-response function  $\alpha(\omega_j)$ , which is defined as  $\alpha(\omega_j) = 1 \forall |\omega_j| < \pi/6$  and zero otherwise. By inverse Fourier transform we obtain the autocovariance matrix  $\tilde{\Sigma}_{yF}(k)$  reflecting low-frequency comovements between  $\hat{F}_{t_m-k}$  and  $y_{t_m}$

$$\tilde{\Sigma}_{yF}(k) = \frac{1}{2H+1} \sum_{j=-H}^H \alpha(\omega_j) \hat{S}_{yF}(\omega_j) e^{i\omega_j k}, \quad (10)$$

which enters equation (7) for  $k = 1, 2, \dots, h_m$ . For given  $M$  and  $H$ , we can compute the projection (6). We will denote this MIDAS approach as ‘MIDAS-smooth’.

The relationship between the basic MIDAS approach in equation (1) or equation (4) and MIDAS-smooth is immediately clear when we disregard the smoothing aspect for a moment, and consider  $\hat{\Sigma}_{yF}(k)$  instead of  $\tilde{\Sigma}_{yF}(k)$  in the projection coefficient  $\hat{\Sigma}_{yF}(h_m) \times \hat{\Sigma}_F^{-1}$  in equation (7). First note that  $\hat{\Sigma}_{yF}(k)$  is a consistent estimator of the true cross-covariance, if the sample size is sufficiently large, despite the missing values. MIDAS-basic (1) and its multivariate extension (4) are based on the same finding as the smooth projection: one regresses low-frequency GDP growth on skip-sampled high-frequency factors, but with a different functional (exponential lag) form and allows for non-zero lag orders. Thus, in terms of lags considered, the New Eurocoin projection is a restricted form of MIDAS-basic, but with a different weighting.



### Unrestricted MIDAS

The MIDAS-basic approach relies on a distributed lag function. As an alternative to these approaches, we also consider an unrestricted lag model

$$y_{t_m+h_m} = \beta_0 + \mathbf{D}(L_m)\hat{\mathbf{F}}_{t_m+w}^{(3)} + \varepsilon_{t_m+h_m}, \quad (11)$$

where  $\mathbf{D}(L_m) = \sum_{j=0}^k \mathbf{D}_j L_m^j$  is an unrestricted lag polynomial of order  $k$ . Koenig *et al.* (2003) proposed a similar model in the context of forecasting with real-time data, but not factors. A theoretical justification for the unrestricted Factor MIDAS is provided in section III, where we derive MIDAS as an approximation to a forecast equation from a high-frequency factor model in the presence of mixed sampling frequencies of the data.

We can estimate  $\mathbf{D}(L_m)$  and  $\beta_0$  by ordinary least squares (OLS). To specify the lag order in the empirical application, we consider a fixed scheme with  $k=0$  and an automatic lag length selection using the Bayesian information criterion (BIC). In general, the analysis in Inoue and Kilian (2006) shows that in-sample BIC model selection can produce good out-of-sample forecasts under mis-specification. In the empirical literature on factor forecasting, Stock and Watson (2002) also successfully employ BIC for model selection, although other criteria could in principle be employed.

Note that for  $k=0$ , we consider only  $(t_m+w)$ -dated factors for forecasting. Thus, with  $k=0$  the projection model is close to the MIDAS-smooth projection as employed in the New Eurocoin index, see the previous discussion. The difference is of course that the smoothing aspect is neglected here. The difference between MIDAS-basic and unrestricted MIDAS will depend on the trade-off between parsimony of the distributed lag function and flexibility of the unrestricted polynomial. If the true lag order and the differences in sampling frequencies are relatively small, we have to estimate a few coefficients only in unrestricted MIDAS. However, the longer the lead-lag relationship in the data is, the more unrestricted MIDAS will suffer from sampling uncertainty owing to the increase of the number of coefficients associated with increasing lag length. In the end, it is a mostly an empirical issue which variant of MIDAS will perform better.

We will denote the unrestricted MIDAS with  $k=0$  as ‘MIDAS-U0’, and with estimated lag order by BIC as ‘MIDAS-U’.

### III. A simple theoretical motivation for Factor MIDAS

Factor MIDAS can be regarded as a plug-in forecast, where factors are estimated in a first step, and then plugged into a regression model for the predictand. Thus, we follow closely the factor forecast literature based on single-frequency data like Stock and Watson (2002) and Boivin and Ng (2005). To motivate the relationship between factor models and Factor MIDAS forecasting, we now discuss a simple benchmark model from Boivin and Ng (2005) and extend this to the mixed-frequency case. In particular, let us consider the monthly one-factor ( $r=1$ ) model

$$\mathbf{X}_{t_m} = \Lambda f_{t_m} + \xi_{t_m}, \quad (12)$$

and assume that monthly unobserved GDP growth is part of  $\mathbf{X}_{t_m}$ . Thus,  $y_{t_m} \in \mathbf{X}_{t_m}$  and

$$y_{t_m} = \lambda_y f_{t_m} + \xi_{y,t_m}. \quad (13)$$

Furthermore, assume the AR model for the factor

$$f_{t_m} = \rho_f f_{t_m} + e_{f,t_m}. \quad (14)$$

We also take into account serial correlation of the idiosyncratic component such that  $\xi_{y,t_m} = \rho_{\xi_y} \xi_{y,t_m-1} + e_{\xi_y,t_m}$ . The error terms  $e_{f,t_m}$  and  $e_{\xi_y,t_m}$  are assumed to be i.i.d. white noise. Given this monthly factor specification, we can write the factor representation for  $y$  as:

$$y_{t_m+1} = \lambda_y f_{t_m+1} + \xi_{y,t_m+1} = \lambda_y (\rho_f f_{t_m} + e_{f,t_m+1}) + (\rho_{\xi_y} \xi_{y,t_m} + e_{\xi_y,t_m+1}). \quad (15)$$

Replacing  $\xi_{y,t_m} = y_{t_m} - \lambda_y f_{t_m}$  and rearranging terms yields

$$y_{t_m+1} = \rho_{\xi_y} y_{t_m} + \lambda_y (\rho_f - \rho_{\xi_y}) f_{t_m} + e_{f,t_m+1} + e_{\xi_y,t_m+1}. \quad (16)$$

Thus, monthly GDP growth 1 month ahead  $y_{t_m+1}$  is explained by lagged GDP growth as well as the lagged factor and white noise terms. In Boivin and Ng (2005), equation (16) serves as a motivation for plug-in estimators with single frequency data. In particular, for a forecast  $h_m$  months ahead and given estimates of the factors  $\hat{f}_{t_m}$ , they propose a forecast equation, where  $y_{t_m+h_m}$  is regressed on lags of  $y_{t_m}$  and  $\hat{f}_{t_m}$ , where the theoretical restrictions from the factor model are ignored. With more general lag polynomials of higher order, this unrestricted equation (16) is then equal to the well-known forecast equation by Stock and Watson (2002). Note that in the theoretical model, the existence of serial correlation in the idiosyncratic component leads to an AR term in the forecast regression, see the term  $\rho_{\xi_y} y_{t_m}$ .

Now, let us consider the mixed-frequency case and assume that  $y_{t_m}$  is a flow variable that can only be observed once in each quarter, that is, for  $t_m = 3, 6, 9, \dots$ , whereas the factor  $f_{t_m}$  can be estimated each month,  $t_m = 1, 2, 3, \dots$ . We want to evaluate the consequences of this mixed-sampling scheme. If  $y_{t_m}$  is regarded as month-on-month GDP growth and GDP is a flow variable, then time aggregation to quarterly frequency is given by  $y_{t_q} = (1 + 2L_m + 3L_m^2 + 2L_m^3 + L_m^4)y_{t_m}$  for  $t_m = 3, 6, 9, \dots$  and  $t_m = 3t_q$ , see Mariano and Murasawa (2003) and Angelini *et al.* (2006). Let us rewrite the time aggregation as  $y_{t_q} = \phi_y(L_m)y_{t_m}$ . As GDP is quarterly, we have to predict for quarterly forecast horizons. For example, assuming a quarterly horizon of one quarter or 3 months, that is,  $h_q = 1$  and  $h_m = 3$ , it turns out that, after a few calculations, we can write equation (16) as:

$$y_{t_m+3} = \rho_{\xi_y}^3 y_{t_m} + \rho_{\xi_y}^2 \lambda_y (\rho_f - \rho_{\xi_y}) (1 + \rho_{\xi_y} + \rho_{\xi_y}^2) f_{t_m} + \rho_{\xi_y} \lambda_y (\rho_f - \rho_{\xi_y}) e_{f,t_m} + \lambda_y (\rho_f - \rho_{\xi_y}) (\rho_f e_{f,t_m} + e_{f,t_m+1}) + \sum_{i=0}^2 \rho_{\xi_y}^i (e_{\xi_y,t_m+3-i} + e_{\xi_y,t_m+3-i}). \quad (17)$$

For notational convenience, we simplify this equation to

$$y_{t_m+3} = \kappa_0 y_{t_m} + \kappa_1 f_{t_m} + \kappa_{ef}(L_m) e_{f,t_m+3} + \kappa_{\xi y}(L_m) e_{\xi y,t_m+3}, \quad (18)$$

which can be regarded as an equation to be estimated by direct estimation. This is in contrast to, for example, vector autoregressive (VAR) models that are generally solved iteratively for forecasting. Thus, the long-lasting discussion of the relative merits of direct vs. iterative forecasting also applies in the present context. Marcellino *et al.* (2006) and Chevillon and Hendry (2005) are recent contributions; see also Bhansali (2002) for a survey. The literature shows that there are arguments in favour of both approaches and, generally, the direct approach seems to dominate in case of substantial mis-specifications. In our case, if the precise dynamics of the deep factor model are unknown, estimating equation (18) directly without taking into account the complicated dynamics can be useful in avoiding mis-specification of the deep model.

Let us now turn to the time aggregation implied in MIDAS. Given the time aggregation scheme  $y_{t_q} = \varphi_y(L_m)y_{t_m}$ , we can now apply time aggregation to both sides of the forecast equation (18)

$$(1 - \kappa_0 L_m^3) \varphi_y(L_m) y_{t_m+3} = \kappa_1 \varphi_y(L_m) f_{t_m} + \varphi_y(L_m) (\kappa_{ef}(L_m) e_{f,t_m+3} + \kappa_{\xi y}(L_m) e_{\xi y,t_m+3}). \quad (19)$$

As time aggregation implies  $y_{t_q+1} = \varphi_y(L_m)y_{t_m+3}$ , the left-hand side of equation (19) contains quarterly GDP growth one quarter ahead, as well as lagged quarterly GDP growth in quarter  $t_q$ . The right-hand side of equation (19) contains complicated monthly lag polynomials of the factors as well as the errors in the factor and idiosyncratic dynamics. Note that the lag polynomials are determined both from the structure of the underlying factor model as well as the aggregation scheme chosen.

Without any theoretical restrictions from the factor models imposed, equation (19) can be regarded as a variant of unrestricted MIDAS, compare it with equation (11). Note that  $\kappa_0$  in equation (19) is a function of the AR coefficient in the idiosyncratic component. Thus, in case there are strong autocorrelations in the idiosyncratic component, we should consider autoregressive terms in MIDAS-U.

MIDAS-basic imposes a certain weighting scheme on the parameters of the monthly regressors in equation (19) through the particular choice of the exponential lag polynomial. Thus, it departs further from the theoretical model that provides the parameter restrictions behind equation (18). However, note that in case the exact dynamics of the true factor model are unknown, for example, lag orders, degree of persistence in the idiosyncratic component and so on, it might be desirable to approximate them by MIDAS as proposed by Ghysels *et al.* (2007). However, depending on the sources of model uncertainty, it is *a priori* unclear which approximation works best empirically.

These theoretical considerations suggest that MIDAS regressions, and in particular Factor MIDAS, should be regarded as approximations to unknown models, rather than structural models. For example, Ghysels and Valkanov (2006) provide a theoretical justification for MIDAS, if the true data-generating process (DGP) is a

high-frequency VAR model. Our approach is similar, but extends their framework to a factor model framework.

#### IV. Estimating the factors with ragged-edge data

Factor forecasting with large, single-frequency datasets is often carried out using a two-step procedure; see, for example, Boivin and Ng (2005). First, the factors are estimated and, second, an AR model for the variable to be predicted is augmented by the estimated factors; see Bai and Ng (2006) for technical details on the properties of the resulting forecasts. A similar two-step procedure can be used in the mixed-frequency case. We have discussed the second step before, that is, the specification of the MIDAS regression. Here, we discuss the first step, factor estimation.

The starting assumption is that the monthly indicators admit a factor representation such as:

$$\mathbf{X}_{t_m} = \mathbf{A}\mathbf{F}_{t_m} + \boldsymbol{\xi}_{t_m}, \quad (20)$$

where  $\mathbf{X}_{t_m}$  is the  $N$ -dimensional vector of endogenous variables, and the  $r$ -dimensional factor vector is denoted by  $\mathbf{F}_{t_m} = (f_{1,t_m}, \dots, f_{r,t_m})'$ . The factor times the  $(N \times r)$  loadings matrix  $\mathbf{A}$  represents the common components of each variable. The idiosyncratic components  $\boldsymbol{\xi}_{t_m}$  are that part of  $\mathbf{X}_{t_m}$  not explained by the factors.

Under the assumption that  $\mathbf{X}_{t_m}$  is balanced and certain general assumptions regarding the idiosyncratic components  $\boldsymbol{\xi}_{t_m}$ , see for example Bai and Ng (2002), various ways to estimate the factors have been provided in the literature. Two of the most widely used approaches are based on PCA, as in Stock and Watson (2002), or DPCA, as in Forni *et al.* (2005). For overviews, see the surveys by Stock and Watson (2006), section 4, and Boivin and Ng (2005) and the comparisons by D'Agostino and Giannone (2006) and Schumacher (2007). However, as mentioned, we are interested in the case where  $\mathbf{X}_{t_m}$  is not balanced, owing to the presence of missing values at the end of the sample for some indicators, the so-called ragged edge. Therefore, we describe in the following subsections three factor estimation methods that can handle ragged-edge data.<sup>2</sup>

##### Vertical realignment of data and dynamic principal components factors

A very convenient way to solve the ragged-edge problem is provided by Altissimo *et al.* (2006) for estimating the New Eurocoin indicator. They propose to realign each time series in the sample to obtain a balanced dataset, see also Schneider and Spitzer (2004). Assume that variable  $i$  is released with  $k_i$  months of publication lag. Thus, given a dataset in period  $T_m + w$ , the final observation available for this time series is for period  $T_m + w - k_i$ . The realignment proposed by Altissimo *et al.* (2006) is then simply

<sup>2</sup>To focus on ragged-edge and mixed-frequency problems, we abstract from additional complications such as those resulting from seasonal adjustment and data revisions.

$$\tilde{x}_{i,t_m} = x_{i,t_m - k_i}, \quad (21)$$

for  $t_m = k_i + 1, k_i + 2, \dots, T_m + w$ . Applying this procedure for each series, and harmonizing at the beginning of the sample, yields a balanced dataset  $\tilde{\mathbf{X}}_{t_m}$  for  $t_m = \max(\{k_i\}_{i=1}^N) + 1, \dots, T_m + w$ .

Given this monthly data, Altissimo *et al.* (2006) propose DPCA to estimate the factors. As the dataset is balanced, the two-step estimation techniques by Forni *et al.* (2005) directly apply. In our applications that follow, we will denote the combination of vertical realignment and DPCA factors as ‘VA-DPCA’.

The vertical realignment solution to the ragged-edge problem is easy to use. A disadvantage is that the availability of data determines dynamic cross-correlations between variables. Furthermore, statistical release dates for data are not the same over time, for example, owing to major revisions. In this case, dynamic correlations within the data change and factors can change over time. The same holds if factors are re-estimated at a higher frequency than the frequency of the factor model. This is a very common scenario, for example, if a monthly factor model is re-estimated several times within a month when new monthly observations are released. If this is the case, the realignment of the data changes the correlation structure all the time. However, DPCA as in Forni *et al.* (2005) exploits the dynamic cross-correlations in the frequency domain and might be able to account for these changes in realignment of the data, at least in principle.

### Principal component factors and the EM algorithm

To consider missing values in the data for estimating factors, Stock and Watson (2002) propose an EM algorithm together with the standard PCA. For a discussion of the properties in the presence of ragged-edge data in real time, see Schumacher and Breitung (2008). Consider a variable  $i$  from the dataset  $\mathbf{X}_{t_m}$  as a full data column vector  $\mathbf{X}_i = (x_{i,1}, \dots, x_{i,T_m+w})'$ . Assume that not all the observations are available owing to the ragged-edge problem. The vector  $\mathbf{X}_i^{\text{obs}}$  contains the observations available for variable  $i$ , which is only a subset of  $\mathbf{X}_i$  owing to missing values. We can formulate the relationship between observed and not fully observed data by

$$\mathbf{X}_i^{\text{obs}} = \mathbf{A}_i \mathbf{X}_i, \quad (22)$$

where  $\mathbf{A}_i$  is a matrix that can tackle missing values or mixed frequencies. In case no observations are missing,  $\mathbf{A}_i$  is the identity matrix. In case an observation is missing at the end of the sample, the corresponding final row of the identity matrix is removed to ensure equation (22). The EM algorithm proceeds as follows:

1. Provide an initial (naive) guess of observations  $\hat{\mathbf{X}}_i^{(0)} \forall i$ . These guesses together with the fully observable monthly time series yields a balanced dataset  $\hat{\mathbf{X}}^{(0)}$ . Standard PCA provides initial monthly factors  $\hat{\mathbf{F}}^{(0)}$  and loadings  $\hat{\mathbf{\Lambda}}^{(0)}$ .

2. **E-step:** An update estimate of the missing observations for variable  $i$  is provided by the expectation of  $\mathbf{X}_i$  conditional on observations  $\mathbf{X}_i^{\text{obs}}$ , factors  $\hat{\mathbf{F}}^{(j-1)}$  and loadings  $\hat{\mathbf{\Lambda}}_i^{(j-1)}$  from the previous iteration is given by

$$\hat{\mathbf{X}}_i^{(j)} = \hat{\mathbf{F}}^{(j-1)} \hat{\mathbf{\Lambda}}_i^{(j-1)} + \mathbf{A}_i'(\mathbf{A}_i' \mathbf{A}_i)^{-1}(\mathbf{X}_i^{\text{obs}} - \mathbf{A}_i \hat{\mathbf{F}}^{(j-1)} \hat{\mathbf{\Lambda}}_i^{(j-1)}). \quad (23)$$

The update consists of two components: the common component from the previous iteration  $\hat{\mathbf{F}}^{(j-1)} \hat{\mathbf{\Lambda}}_i^{(j-1)}$ , plus the low-frequency idiosyncratic component  $\mathbf{X}_i^{\text{obs}} - \mathbf{A}_i \hat{\mathbf{F}}^{(j-1)} \hat{\mathbf{\Lambda}}_i^{(j-1)}$ , distributed by the projection coefficient  $\mathbf{A}_i'(\mathbf{A}_i' \mathbf{A}_i)^{-1}$  on the high-frequency periods; for details, see Schumacher and Breitung (2008). Repeat the E-step for all  $i$  yielding again a balanced dataset  $\hat{\mathbf{X}}^{(j)}$ .

3. **M-step:** Re-estimate the factors and loadings,  $\hat{\mathbf{F}}^{(j)}$  and  $\hat{\mathbf{\Lambda}}^{(j)}$  by PCA, and go to Step 2 until convergence.

After convergence, the EM algorithm provides monthly factor estimates  $\hat{\mathbf{F}}_{t_m}$  as well as estimates of the missing values of the time series. Thus, interpolation of missing values as well as factor estimation is carried out consistently in the factor framework (20) with factors estimated by PCA. Thus, an alignment of data as in VA-DPCA can be avoided with this approach. For a detailed discussion of the properties of the EM algorithm for interpolation and backcasting, see Angelini *et al.* (2006). In the applications that follow, we will denote this factor estimator as ‘EM-PCA’.

### Estimation of a large factor model in state-space form

The factor estimation approach followed by Doz *et al.* (2006) and Giannone *et al.* (2008) is based on a complete representation of the large factor model in state-space form. The model consists of the factor representation (20) and a VAR model for the factors. The full state-space model has the form

$$\mathbf{X}_{t_m} = \mathbf{A} \mathbf{F}_{t_m} + \boldsymbol{\xi}_{t_m}, \quad (24)$$

$$\boldsymbol{\Psi}(L_m) \mathbf{F}_{t_m} = \mathbf{B} \boldsymbol{\eta}_{t_m}. \quad (25)$$

Equation (24) is the static factor representation of  $\mathbf{X}_{t_m}$ . Equation (25) specifies a VAR of the factors with lag polynomial  $\boldsymbol{\Psi}(L_m) = \sum_{i=1}^p \boldsymbol{\Psi}_i L_m^i$ . The  $q$ -dimensional vector  $\boldsymbol{\eta}_{t_m}$  contains the orthogonal dynamic shocks that drive the  $r$  factors, where the matrix  $\mathbf{B}$  is  $(r \times q)$ -dimensional. The model is already in state-space form, as the factors  $\mathbf{F}_{t_m}$  are the states. If the dimension of  $\mathbf{X}_{t_m}$  is small, the model can be estimated using iterative maximum likelihood (ML). To account for large datasets, Doz *et al.* (2006) propose quasi-ML to estimate the factors, as iterative ML is infeasible in this framework. For a given number of factors  $r$  and dynamic shocks  $q$ , the estimation proceeds in the following steps:

1. Estimate  $\hat{\mathbf{F}}_{t_m}$  using PCA as an initial estimate. Here, estimation is based on the balanced part of the data. We can obtain this by removing as many values as possible at the end of the sample as long as the dataset is unbalanced.
2. Estimate  $\hat{\mathbf{A}}$  by regressing  $\mathbf{X}_{t_m}$  on the estimated factors  $\hat{\mathbf{F}}_{t_m}$ . The covariance of the idiosyncratic components  $\hat{\xi}_{t_m} = \mathbf{X}_{t_m} - \hat{\mathbf{A}}\hat{\mathbf{F}}_{t_m}$ , denoted as  $\hat{\Sigma}_{\xi}$ , is also estimated.
3. Estimate a VAR( $p$ ) on the factors  $\hat{\mathbf{F}}_{t_m}$  yielding  $\hat{\Psi}(L)$  and the residual covariance of  $\hat{\xi}_{t_m} = \hat{\Psi}(L_m)\hat{\mathbf{F}}_{t_m}$ , denoted as  $\hat{\Sigma}_{\xi}$ .
4. To obtain an estimate for  $\mathbf{B}$ , given the number of dynamic shocks  $q$ , apply an eigenvalue decomposition of  $\hat{\Sigma}_{\xi}$ . Let  $\mathbf{M}$  be the  $(r \times q)$ -dimensional matrix of the eigenvectors corresponding to the  $q$  largest eigenvalues, and let the  $(q \times q)$ -dimensional matrix  $\mathbf{P}$  contain the largest eigenvalues on the main diagonal and zero otherwise. Then, the estimate of  $\mathbf{B}$  is  $\hat{\mathbf{B}} = \mathbf{M} \times \mathbf{P}^{-1/2}$ . The coefficients and auxiliary parameters of the system of equations (24) and (25) is fully specified numerically. The model is cast into state-space form.
5. The Kalman smoother then yields new estimates of the monthly factors. The dataset used for Kalman smoother estimation is now the unbalanced dataset for  $t_m = 1, \dots, T_m + w$ , and  $T_m + w$  is the latest observation available in the entire set of monthly time series.

If missing values at the end of the sample are present, as in our setup, the Kalman smoother also yields optimal estimates and forecasts for these values conditional on the model structure and properties of the shocks. Thus, it is well suited to tackle ragged-edge problems as in the present context. Nonetheless, one has to keep in mind that in this case the coefficients in system matrices have to be estimated from a balanced sub-sample of data, as in step 1 a fully balanced dataset is needed for PCA initialization. However, although the system matrices are estimated on balanced data in the first step, the factor estimation based on the Kalman smoother applies to the unbalanced data and can tackle ragged-edge problems. The solution is to estimate coefficients outside the state-space model and avoid estimating a large number of coefficients by iterative ML. Compared with the other two approaches mentioned before, the state-space approach can take factor dynamics into account explicitly. However, a richer model structure requires the specification of more auxiliary coefficients and can be subject to mis-specification, so it is in the end an empirical question whose approach works best, see the discussion in Boivin and Ng (2005). In the applications that follow, we will denote the state-space model Kalman filter estimator of the factors as ‘KFS-PCA’.

## V. Design of the nowcast and forecast comparison exercise

So far, we have introduced nine alternative Factor MIDAS approaches, generated by the combination of the three types of MIDAS regressions in section II with the three factor estimators of section IV. As it is difficult to rank the alternative approaches

*a priori*, and their relative performance might depend on the DGP, we will illustrate their implementation and assess their relative merits by means of an empirical application: nowcasting and forecasting quarterly German GDP growth with a large set of monthly indicators. In this section we discuss, in turn, the dataset, the design of the exercise and the specifications of the Factor MIDAS models. In the next section, we present the results.

### Data and replication of the ragged edge

The dataset contains German quarterly GDP growth from 1992Q1 until 2006Q3 and 111 monthly indicators from 1992M1 until 2006M11 that cover a wide range of economic activity in Germany. The dataset is a final dataset. It is not a real-time dataset and does not contain vintages of data, as they are not available for Germany for such a broad coverage of time series. Furthermore, in Schumacher and Breitung (2008) a considerably smaller real-time dataset for Germany is used, but the results indicate that data revisions do not affect the forecast accuracy considerably. Similar results have been found by Boivin and Ng (2005) for the United States in a similar context. More information about the data can be found in the Appendix.

To consider the ragged edge of the data at the end of the sample owing to different publication lags, we follow Banbura and Rünstler (2007) and replicate the ragged edge from the one final vintage of data that is available. When downloading the data – the download date for the data used here was 6 December 2006 – we observe the ragged-edge pattern in terms of the missing values at the end of the data sample. For example, at the beginning of December 2006, we observe interest rates until November 2006, thus there is only one missing value at the end of the sample, whereas industrial production is available up to September 2006, implying three missing values. For each time series, we store the missing values at the end of the sample. Under the assumption that these patterns of data availability remain stable over time, we can impose the same missing values at each point in time of the recursive experiment. Thus, we shift the missing values back in time to mimic the availability of information as in real time.

### Nowcast and forecast design

To evaluate the performance of the models, we carry out recursive estimation and nowcasting, where the full sample is split into an evaluation sample and an estimation sample, which is recursively expanded over time. The evaluation sample is between 1998Q4 and 2006Q3. For each of these quarters, we want to compute nowcasts and forecasts depending on different monthly information sets. For example, for the initial evaluation quarter 1998Q4, we want to compute a nowcast in December 1998, one in November and October, whereas the forecasts are computed from September 1998 backwards in time accordingly. Thus, we have three nowcasts computed at the



beginning of each of the intra-quarter months. Concerning the forecasts, we present results up to two quarters ahead. Thus, again for the initial evaluation quarter 1998Q4, we have six forecasts computed based on information available in April 1998 up to information available in September 1998. Overall, we have nine projections for each GDP observation of the evaluation period, depending on the information available to make the projection.

The estimation sample depends on the information available at each period in time when computing the nowcasts and forecasts. Assume again we want to nowcast GDP growth for 1998Q4 in December 1998, then we have to identify the time-series observations available at that period in time. For this purpose, we exploit the ragged-edge structure from the end of the full sample of data, as discussed in the previous subsection. For example, for the nowcast GDP growth for 1998Q4 made in December 1998, we know from our full sample that at each period in time, we have one missing value for interest rates and three missing values of industrial production. These missing values are imposed also for the period December 1998, thus replicating the same ragged-edge pattern of data availability. We do this accordingly in every recursive subsample to determine the pseudo real-time final observation of each time series. The first observation for each time series is the same for all recursions, namely 1992M1. This implies the recursive design with increasing information over time available for estimating the factor models. To replicate the publication lags of GDP, we exploit the fact that GDP of the previous quarter is available for nowcasting and forecasting at the beginning of the third month of the next quarter. Note that we re-estimate the factors and forecast equations every recursion when new information becomes available; so, factor weights and forecast model coefficients are allowed to change over time.

For each evaluation period, we compute nine nowcasts and forecasts depending on the available information. To compare the nowcasts with the realizations of GDP growth, we use the mean-squared error (MSE). As a measure of informativeness of the nowcasts, we relate the MSE to the variance of GDP growth, where the variance is computed over the evaluation period, see Forni *et al.* (2003). A relative MSE to the variance less than one indicates that the forecast of a model for the chosen nowcast and forecast horizon is to some extent informative for current and future GDP growth; see the discussion of on forecastability in Clements and Hendry (1998), section 2.5. Note that this relative statistic can also be interpreted as a measure to compare the MSE of the factor models with the corresponding MSE of the out-of-sample mean of GDP growth as a naive forecast.

In the empirical application, we also compare the monthly Factor MIDAS forecasts with those from purely quarterly models. As a benchmark to those models, we consider a simple quarterly AR model, with the lag order specified by BIC. Furthermore, we provide the recursive in-sample mean of GDP growth as an additional benchmark. Note that this benchmark is equivalent to a random-walk model of the log level of GDP. In the recent forecasting literature, this in-sample mean of GDP growth benchmark has turned out to be a strong competitor to

more sophisticated approaches, see, for example, De Mol, Giannone and Reichlin (2008).

### Specification of models

To specify the number of factors in the applications that follow, we follow two approaches. We determine the number of static and dynamic factors,  $r$  using information criterion  $IC_{p2}$  from Bai and Ng (2002). We determine  $q$  by the criterion suggested by Bai and Ng (2007). Additionally, we compute nowcasts and forecasts for all possible combinations of  $r$  and  $q$  and evaluate them. In our application, we consider a maximum of  $r = 6$  and all combinations of  $r$  and  $q$  with  $q \leq r$ . Details can be found in appendix C of the working paper version. The key result from this exercise is that only for the case  $r = 1$  and partly for  $r = 2$ , nowcasts and forecasts have information content for current and future GDP. Apart from a few exceptions, all other combinations of numbers of factors – including those determined by information criteria – performed worse and less stable than the specifications we provide results for in the main text that follows. A reason for this result is the combination of the rather short estimation sample and the substantial likelihood of parameter changes. In this case, Banerjee, Marcellino and Masten (2008) show that there is a substantial deterioration in the performance of forecasts based on many factors, and model specification by information criteria is not helpful. Also for the United States, it was shown that only very few factors can obtain satisfactory forecast results, see Stock and Watson (2002) and Schumacher and Breitung (2008). As a result of these findings and so as to preserve space, we only present results for  $r = 1$  in the following.

For estimating the state-space factor model, a lag order determination is required to specify the factor VAR( $p$ ). For this purpose, we apply the BIC with a maximum lag order of  $p = 6$  months. The chosen lag lengths are usually very small with only one or two lags in most of the cases. To specify the dynamic PC estimator and MIDAS-smooth, we use the frequency domain parameters  $M = 24$  and  $H = 60$  for estimating the spectral density; see appendix B in Marcellino and Schumacher (2008) for details.

The EM algorithm we implement for monthly factor estimation is slightly different from what has been described before. In particular, we do not update the factor weights during the iterations. We rather exploit the fact that the covariance matrix of the monthly data can be consistently estimated despite the missing values at the end of the sample. To estimate the covariance, we simply compute pairwise covariances over the periods when both series are available. Thus, the EM algorithm is only used to interpolate the missing values and estimate the factors by the fixed weights times the data, which partly consists of estimated observations. We adopt this simplification to prevent convergence problems and to speed up the convergence process. As a stopping rule, we assume that convergence is achieved if the change in the average sum of squares of the idiosyncratic components is smaller than  $10^{-5}$ .

Concerning the specifications of MIDAS, we use a large variety of initial parameter specifications, and compute the residual sum of squares (RSS). The parameter

set with the smallest RSS then serves as the initial parameter set for NLS estimation. The parameters of the exponential lag function are restricted to  $\theta_1 < 2/5$  and  $\theta_2 < 0$ , in line with Ghysels *et al.* (2007). The maximum number of lags chosen for MIDAS is  $K = 12$  months.

## VI. Empirical nowcast and forecast comparison for German GDP

In this section, we present and discuss the empirical results for German GDP. In particular, following the previous methodological discussion, we present a comparison of factor estimation methods that can tackle ragged-edge data in the subsequent subsection, followed by a comparison of Factor MIDAS projections in the subsection on ‘Role of type of MIDAS projection’. Then, to relate our results to earlier empirical findings and conceptual discussions in the factor forecast literature, the subsection on the comparison of monthly and quarterly models compares the monthly nowcast models with quarterly factor models. Static vs. dynamic factor estimation is discussed next and the last subsection evaluates the relative merits of our two-step Factor MIDAS approach and of the integrated state-space model of Giannone *et al.* (2008).

### The role of factor estimation methods

Nowcast and forecast results for the different combinations of MIDAS projections and factor estimation methods can be found in Table 1. This table is divided into four parts. For each of the four MIDAS projections, we can compare the different factor estimation methods. The table shows MSEs relative to the GDP growth variance, and rankings based on those relative MSEs, where models with the smallest MSE rank first. The nowcast and forecast horizons are shown in Table 1 for monthly horizons  $h_m = 1, \dots, 9$ , where horizons 1–3 belong to the nowcast. Horizon  $h_m = 1$  is a nowcast made in the third month of the respective quarter, whereas horizon  $h_m = 2$  is the nowcast made in the second month of the current quarter. Thus, similar to standard forecast comparisons, increasing horizons correspond to less information available for nowcasting and forecasting, and we expect an increasing MSE for increasing horizons  $h_m$ .

The projections from the factor models have information content for the nowcast, as the MSEs of virtually all combinations of factor estimation and projection methods yield MSEs smaller than one, see Table 1. For the one-quarter ahead forecast, we find borderline results. Comparing the factor estimation methods at horizons 4–6, the results are not clear cut, as some relative MSEs are larger than one for some horizons and smaller for others. For two quarters ahead, the relative MSEs are for all factor models larger than one, thus rendering all factor models at hand uninformative for this horizon. This indicates that the methods employed here can be regarded as suited for short-term nowcasting and forecasting only.

The differences between the factor estimation methods are relatively small overall. In the rankings of nowcast performance, EM-PCA factors do best in many cases in

TABLE 1

*Comparison of nowcast and forecast results for different factor estimation methods for  $r=1$ , mean-squared error (MSE) relative to gross domestic product (GDP) variance and ranking*

	Horizon $h_m$	Nowcast			Forecast					
		Current quarter			One quarter			Two quarters		
		1	2	3	4	5	6	7	8	9
1.a. MIDAS-basic	VA-DPCA	0.71	1.01	1.05	1.02	1.19	1.05	1.16	1.24	1.26
	EM-PCA	0.62	0.69	0.78	1.01	1.09	0.94	1.21	1.09	1.05
	KFS-PCA	0.79	0.90	0.87	1.07	1.17	1.07	1.19	1.13	1.20
1.b. Ranking	VA-DPCA	2	3	3	2	3	2	1	3	3
	EM-PCA	1	1	1	1	1	1	3	1	1
	KFS-PCA	3	2	2	3	2	3	2	2	2
2.a. MIDAS-U	VA-DPCA	0.90	1.05	1.02	1.04	1.15	1.11	1.19	1.13	1.17
	EM-PCA	0.92	0.65	0.72	1.08	1.05	0.90	1.19	1.42	1.40
	KFS-PCA	0.89	0.90	0.81	0.97	1.03	1.02	1.31	1.49	1.36
2.b. Ranking	VA-DPCA	2	3	3	2	3	3	2	1	1
	EM-PCA	3	1	1	3	2	1	1	2	3
	KFS-PCA	1	2	2	1	1	2	3	3	2
3.a. MIDAS-smooth	VA-DPCA	0.69	0.92	0.87	0.95	1.10	1.20	1.18	1.12	1.19
	EM-PCA	0.70	0.73	0.84	0.94	0.95	1.00	1.05	1.09	1.13
	KFS-PCA	0.76	0.85	0.89	0.98	1.06	1.08	1.10	1.16	1.19
3.b. Ranking	VA-DPCA	1	3	2	2	3	3	3	2	3
	EM-PCA	2	1	1	1	1	1	1	1	1
	KFS-PCA	3	2	3	3	2	2	2	3	2
4.a. MIDAS-U0	VA-DPCA	0.71	0.86	0.89	0.90	1.05	0.98	1.05	1.09	1.12
	EM-PCA	0.58	0.65	0.72	0.92	0.93	0.79	1.10	1.10	1.05
	KFS-PCA	0.68	0.85	0.80	0.95	1.01	0.93	1.08	1.09	1.06
4.b. Ranking	VA-DPCA	3	3	3	1	3	3	1	2	3
	EM-PCA	1	1	1	2	1	1	3	3	1
	KFS-PCA	2	2	2	3	2	2	2	1	2

*Notes:* The variance of GDP growth in the evaluation sample is 0.246. In the rankings, models with the smallest MSE rank first. MIDAS indicates mixed-data sampling. VA-DPCA refers to the vertical realignment and dynamic principal components analysis (PCA) used in Altissimo *et al.* (2006), EM-PCA is the expectation-maximization algorithm together with PCA as in Stock and Watson (2002) and KFS-PCA is the Kalman smoother of state-space factors according to Doz *et al.* (2006). The projection MIDAS-basic is the projection from Ghysels and Valkanov (2006), MIDAS-U is unrestricted MIDAS without exponential lag polynomial and lag specification using Bayesian information criterion. MIDAS-smooth is the projection as employed in Altissimo *et al.* (2006), and MIDAS-U0 is the MIDAS projection with unrestricted lag polynomials of order zero.

terms of ranking. However, for  $h_m = 1$  and using the MIDAS-U and MIDAS-smooth projections together with factors VA-DPCA and KFS-PCA, respectively, do better than EM-PCA. Across projection methods, there are no systematic differences in

nowcasting performance between factor estimation by VA-DPCA and KFS-PCA, as the relative MSE rankings change depending on the nowcast and forecast horizons.

Note that in Table 1 increasing the nowcast or forecast horizon month by month does not always leads to an increase of relative MSE, although this happens in most of the cases. This can be observed across all models under comparison. Thus, as new monthly information becomes available, the methods employed here cannot always improve the nowcasts and forecasts with this information. This could be a result of the relatively short sample under consideration that induces high sampling uncertainty of the estimates and nowcasts.

The relative comparison of the factor estimation methods was based on the MSE as a performance measure so far. However, as the MSE averages over observations in the evaluation period, this statistic can be dominated by differences in performance in only a few periods. Therefore, we additionally investigate the factor nowcasts over recursions. In Figure 1, the time series of nowcasts for  $h_m = 1, 2, 3$  are shown together with GDP growth observations and the in-sample mean as a benchmark nowcast for different factor estimation methods. Concerning the type of projection, Figure 1 includes results for MIDAS-U0 only. As the results are very similar for the other types of MIDAS projections, we leave them out of this comparison here. The same holds for the forecast horizons  $h_m \geq 3$ .

The results in Figure 1 show that the three-factor models perform clearly better than the simple benchmark. However, the erratic movements of GDP growth at the beginning of the sample, for example, in 2000Q2 and 2000Q3, are not predicted well by all three-factor models. Increasing the nowcast horizon from  $h_m = 1$  to 3 shows the decline in variance of the nowcasts and, thus, a decline in nowcast ability. A common finding of the figures is the high correlation between the forecasts of the three-factor models, as periods of good and bad performances are similar. Therefore, in line with the similar MSE findings before, we find no clear indications of dramatic differences between the nowcast accuracy of the three-factor models over time.

### **The role of the type of MIDAS projection**

Next, we discuss the different types of MIDAS projections. The nowcast results can be found in Table 2. This table contains three groups for each of the factor estimation methods. For each factor estimation method, we will compare the different MIDAS projections. In Table 2, a general finding is that the differences between the MIDAS approaches are not big as all approaches lead to nowcasts that have information content for current GDP growth, and only a few combinations of factor estimation and MIDAS projection also have predictive ability for the next quarter. Comparing the methods, we see that the difference between MIDAS-basic based on exponential lags and MIDAS-U is not clear-cut, as none of them outperforms the other across all factor estimation methods and horizons. The most simple MIDAS projections without lags of the factors, MIDAS-U0 and MIDAS-smooth, provide often better nowcasts than MIDAS based on exponential lag functions, MIDAS-basic or MIDAS-U. MIDAS-

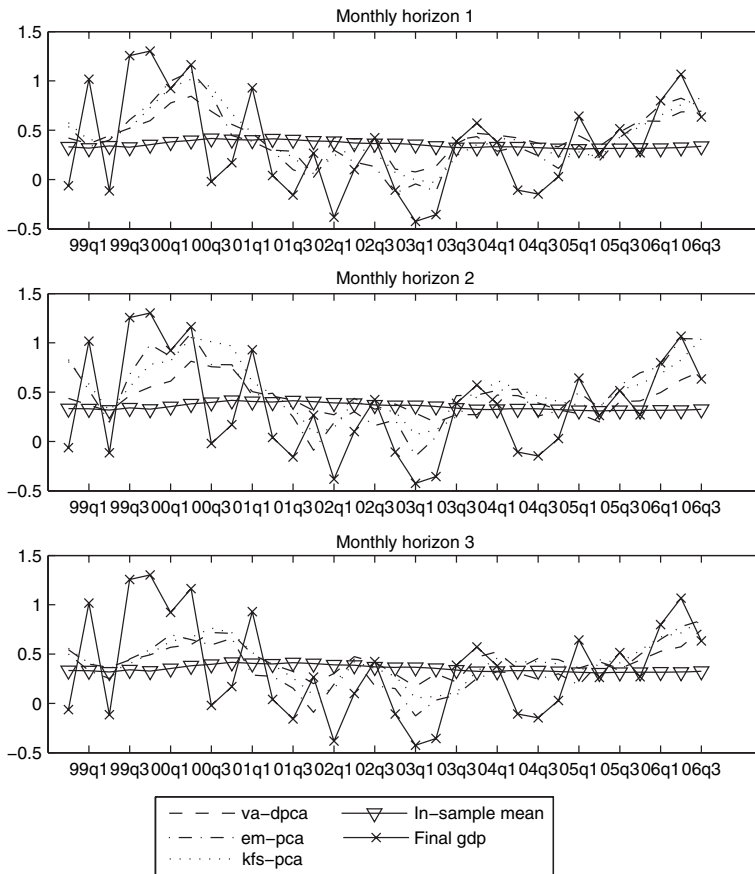


Figure 1. Nowcasts with MIDAS-U0 and different factor estimation methods for horizon  $h_m = 1, 2, 3$  and GDP observations, quarter on quarter growth, number of factors  $r = 1$  and  $q = 1$

*Note:* The figure shows nowcasts for the different factor estimation methods and the in-sample mean as a benchmark. For the model descriptions and abbreviations, see Table 1.

smooth can outperform MIDAS-U0 only for factors obtained by VA-DPCA for a few horizons. However, based on EM-PCA and KFS-PCA factors, the projection MIDAS-U0 outperforms MIDAS-smooth at all horizons. Thus, the simplest projection method MIDAS-U0 seems to work more stable than the other methods overall, as it ranks first or second in most of the cases. The good performance of MIDAS-U0 is likely because of the adoption of a very parsimonious specification, which often pays when forecasting, as it reduces estimation uncertainty and frequently only slightly biases the forecasts; see, for example, Clements and Hendry (1998, ch. 12).

As an AR model is often supposed to be an appropriate benchmark specification for GDP growth, the MIDAS-AR discussed in section II might give additional insights in which direction the other MIDAS approaches considered so far might be improved. In Table 3, the MIDAS-AR is compared with the MIDAS-basic without

TABLE 2

*Comparison of nowcast and forecast results from different MIDAS projections for  $r=1$ , mean-squared error (MSE) relative to gross domestic product (GDP) variance and ranking*

		Nowcast			Forecast					
		Current quarter			One quarter			Two quarters		
	Horizon $h_m$	1	2	3	4	5	6	7	8	9
1.a. VA-DPCA	MIDAS-basic	0.71	1.01	1.05	1.02	1.19	1.05	1.16	1.24	1.26
	MIDAS-U	0.90	1.05	1.02	1.04	1.15	1.11	1.19	1.13	1.17
	MIDAS-smooth	0.69	0.92	0.87	0.95	1.10	1.20	1.18	1.12	1.19
	MIDAS-U0	0.71	0.86	0.89	0.90	1.05	0.98	1.05	1.09	1.12
1.b. Ranking	MIDAS-basic	3	3	4	3	4	2	2	4	4
	MIDAS-U	4	4	3	4	3	3	4	3	2
	MIDAS-smooth	1	2	1	2	2	4	3	2	3
	MIDAS-U0	2	1	2	1	1	1	1	1	1
2.a. EM-PCA	MIDAS-basic	0.62	0.69	0.78	1.01	1.09	0.94	1.21	1.09	1.05
	MIDAS-U	0.92	0.65	0.72	1.08	1.05	0.90	1.19	1.42	1.40
	MIDAS-smooth	0.70	0.73	0.84	0.94	0.95	1.00	1.05	1.09	1.13
	MIDAS-U0	0.58	0.65	0.72	0.92	0.93	0.79	1.10	1.10	1.05
2.b. Ranking	MIDAS-basic	2	3	3	3	4	3	4	2	1
	MIDAS-U	4	1	1	4	3	2	3	4	4
	MIDAS-smooth	3	4	4	2	2	4	1	1	3
	MIDAS-U0	1	2	2	1	1	1	2	3	2
3.a. KFS-PCA	MIDAS-basic	0.79	0.90	0.87	1.07	1.17	1.07	1.19	1.13	1.20
	MIDAS-U	0.89	0.90	0.81	0.97	1.03	1.02	1.31	1.49	1.36
	MIDAS-smooth	0.76	0.85	0.89	0.98	1.06	1.08	1.10	1.16	1.19
	MIDAS-U0	0.68	0.85	0.80	0.95	1.01	0.93	1.08	1.09	1.06
3.b. Ranking	MIDAS-basic	3	4	3	4	4	3	3	2	3
	MIDAS-U	4	3	2	2	2	2	4	4	4
	MIDAS-smooth	2	2	4	3	3	4	2	3	2
	MIDAS-U0	1	1	1	1	1	1	1	1	1

*Note:* For model abbreviations, see Table 1.

AR terms. The results in Table 3 show that considering AR terms does not improve the nowcast and forecast performances systematically. For different horizons and different factor estimation methods, the ranking between MIDAS-AR and MIDAS-basic changes. MIDAS-AR is not generally better than MIDAS-basic, which might also indicate problems with estimating AR dynamics in German GDP growth. Note that we also tried to augment the unrestricted MIDAS with AR terms. However, also this experiment did not lead to clear-cut improvements in forecast performance.

### Comparing monthly with quarterly models

We now investigate the relative advantages of the nowcast factor models with earlier factor approaches in the literature. A widely followed approach in the previous

TABLE 3

*MIDAS-AR vs. MIDAS-basic, comparison of relative mean-squared errors (MSEs) for  $r = 1$* 

		Nowcast			Forecast					
		Current quarter			One quarter			Two quarters		
	Horizon $h_m$	1	2	3	4	5	6	7	8	9
1.a. VA-DPCA	MIDAS-AR	0.76	0.90	0.91	1.09	1.12	1.04	1.24	1.21	1.29
	MIDAS-basic	0.71	1.01	1.05	1.02	1.19	1.05	1.16	1.24	1.26
1.b. Ranking	MIDAS-AR	2	1	1	2	1	1	2	1	2
	MIDAS-basic	1	2	2	1	2	2	1	2	1
2.a. EM-PCA	MIDAS-AR	0.64	0.63	0.75	1.04	1.32	0.99	1.18	1.05	1.35
	MIDAS-basic	0.62	0.69	0.78	1.01	1.09	0.94	1.21	1.09	1.05
2.b. Ranking	MIDAS-AR	2	1	1	2	2	2	1	1	2
	MIDAS-basic	1	2	2	1	1	1	2	2	1
3.a. KFS-PCA	MIDAS-AR	0.90	0.93	0.84	1.07	1.18	1.13	1.31	1.19	1.26
	MIDAS-basic	0.79	0.90	0.87	1.07	1.17	1.07	1.19	1.13	1.20
3.b. Ranking	MIDAS-AR	2	2	1	2	2	2	2	2	2
	MIDAS-basic	1	1	2	1	1	1	1	1	1

*Notes:* The projection MIDAS-AR contains one autoregressive term as in Clements and Galvão (2008) and MIDAS basic is without AR terms. For factor model abbreviations, see Table 1.

literature on factor forecasting is time aggregation. To obtain a balanced sample of data, one can simply aggregate the monthly data to quarterly data and ignore the most recent observations of high-frequency indicators. Then, the standard techniques of factor forecasting with single-frequency data can be employed. Note that previously most of the studies for forecasting of German GDP growth were based on quarterly, partly time-aggregated data; see, for example, Schumacher (2007). As quarterly data are widely used in the empirical literature for GDP forecasting, we will also compare the mixed-frequency nowcast models with the quarterly factor models.

In particular, we employ the standard model for factor forecasting following Stock and Watson (2002). The forecast equation is essentially a quarterly factor-augmented AR model according to

$$y_{t_q+h_q} = \beta_0 + \lambda(L_q)y_{t_q} + \mathbf{E}(L_q)\hat{\mathbf{F}}_{t_q}^Q + \varepsilon_{t_q+h_q}, \quad (26)$$

where  $\mathbf{E}(L_q) = \sum_{p=0}^P \mathbf{E}_p L_q^p$  is an unrestricted lag polynomial of order  $P$ , and  $L_q$  is the quarterly lag operator now.  $\lambda(L_q)$  is now a lag polynomial of order  $R$  for AR terms. We consider fixed specifications with  $P=R=0$ , as well as BIC lag order selection with up to two lags. The factors  $\hat{\mathbf{F}}_{t_q}^Q$  are estimated by PCA, which is applied to the quarterly indicators. These time-series indicators are the same as for the nowcast models as discussed before, but aggregated over time to quarterly frequency. Note that model (26) with static factors  $\hat{\mathbf{F}}_{t_q}^Q$  works quite well for single-frequency data compared with dynamic factor estimates, see Boivin and Ng (2005), D'Agostino and



Giannone (2006) and Schumacher (2007) for German GDP. Thus, it might serve as an interesting alternative to the MIDAS nowcast models. As a benchmark for the factor nowcast models, we employ a univariate quarterly AR model for GDP growth, specified using the BIC as before. In the application, it turns out that in almost all of the recursions only one lag is chosen. Furthermore, we present the in-sample mean of GDP growth as an additional benchmark.

Table 4 contains results for the nowcasts using quarterly factor models as well as the simple benchmarks. As representatives of the nowcast models, we present results based on MIDAS-U0 and the three different ragged-edge factor estimation methods. In the empirical nowcast comparison, the simple benchmarks do not perform well, as can be seen from the bottom rows of Table 4. Both the AR model and the in-sample mean have relative MSEs larger than one. Note that, whereas the nowcast factor models employ monthly information which is updated every month and, thus, can lead to changes in nowcast and forecast MSEs, the benchmark models and the quarterly factor models change only every third month (when a new observation of GDP is available), implying a constant MSE for 3 months.

The quarterly factor model performs better than the naive benchmarks, and has some information content for GDP growth for horizons up to 3 months. For longer horizons, there is almost no information content in the forecasts. Compared with the monthly nowcast models, the quarterly factor model is generally outperformed for the nowcast for  $h_m \leq 3$ . In many cases, this also holds for the one-quarter ahead forecast, although the differences are smaller at these horizons. Thus, according to these results,

TABLE 4

*Comparison of mixed-frequency nowcast models with MIDAS-U0 and quarterly factor and benchmark models, mean square error (MSE) relative to gross domestic product (GDP) variance and ranking*

		<i>Nowcast</i>			<i>Forecast</i>					
		<i>Current quarter</i>			<i>One quarter</i>			<i>Two quarters</i>		
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>	<i>7</i>	<i>8</i>	<i>9</i>
<i>Horizon <math>h_m</math></i>										
1. MIDAS-U0	VA-DPCA	0.71	0.86	0.89	0.90	1.05	0.98	1.05	1.09	1.12
	EM-PCA	0.58	0.65	0.72	0.92	0.93	0.79	1.10	1.10	1.05
	KFS-PCA	0.68	0.85	0.80	0.95	1.01	0.93	1.08	1.09	1.06
2. Quarterly factor models	PCA, $P=R=0$ , $r=1$	0.98	1.05	1.05	1.05	1.16	1.16	1.16	1.14	1.14
	PCA, $P=R=0$ , $r=2$	1.03	0.94	0.94	0.94	1.31	1.31	1.31	1.25	1.25
	PCA-BIC, $r=1$	0.91	1.02	1.02	1.02	1.13	1.13	1.13	1.08	1.08
	PCA-BIC, $r=2$	0.94	0.89	0.89	0.89	1.23	1.23	1.23	1.23	1.23
	PCA-BIC	0.91	1.03	1.03	1.03	1.13	1.13	1.13	1.08	1.08
3. Benchmarks	AR	1.02	1.17	1.17	1.17	1.08	1.08	1.08	1.08	1.08
	In-sample mean	1.03	1.04	1.04	1.04	1.05	1.05	1.05	1.06	1.06

*Notes:* The quarterly factor model contains factors estimated from quarterly, time-aggregated data using static PCA. The first two specifications are based on a fixed number of factors and a fixed number of lags, whereas the third and fourth are based on a fixed number of factors and the number of lags is chosen by BIC. PCA-BIC selects the number of factors as well as the lag orders using BIC as in Stock and Watson (2002). Concerning the nowcast models, the abbreviations are explained in Table 1.

taking into account ragged-edge information as in the nowcast models with monthly indicators can improve the current estimate of GDP growth. As the use of time-aggregated data implies a loss of information at the end of the sample, the results imply that the nowcast methods employed here can to some extent exploit this information. In summary, we can confirm that in general it is advisable to employ the ragged-edge data together with the different factor estimation techniques for nowcasting.

### **Static vs. dynamic factors**

Following the discussion in Boivin and Ng (2005), there is some disagreement in the literature concerning the appropriate factor estimation method to be employed for forecasting. In particular, it is unclear whether DPCA or PCA are favourable for predictive purposes. In general, there is no consensus as to the appropriate estimation method; see also the discussion in Schneider and Spitzer (2004), Den Reijer (2005), D'Agostino and Giannone (2006) and Boivin and Ng (2005) for different datasets. The problem is partly related to the fact that a static factor model such as equation (20) can approximate a dynamic specification when a sufficiently large number of factors is included, see, for example, Stock and Watson (2006) on this point. In a dataset for the German economy with balanced recursive samples, DPCA does not generally work better than PCA, and the differences between the methods are small, see Schumacher (2007).

Against the background of this discussion, we will address this issue also in the present context. In our previous applications, DPCA was employed to estimate the factors in combination with vertical realignment of the data. To assess the sensitivity of the results, we compare the existing figures using VA-DPCA with those resulting from static PCA and vertical realignment of the data, denoted as VA-PCA. Table 5 shows relative MSEs to GDP growth variance for the different factor estimates and different projection techniques. The results show that the information content of the nowcasts and forecasts does hardly change if the factors are estimated by PCA instead of DPCA. MSEs relative to GDP growth variance are in most of the cases above or below one for both factor estimators. The bottom part of Table 5 shows another relative MSE defined as the MSE obtained from using DPCA factors divided by the MSE obtained from using static PCA factors for forecasting. The results show no systematic advantages over the horizons between the two methods. Thus, the way the factors are estimated seems to be of limited importance in this application.

### **Integrated state-space model approach vs. two-step nowcasting**

The results obtained so far are entirely based on a two-step procedure: the factors are estimated first, and then forecasting is carried out using the MIDAS approaches. However, among the models, the state-space approach allows in general for joint estimation of the factors and nowcasting GDP. Banbura and Rünstler (2007)

TABLE 5

*Static PCA vs. dynamic PCA nowcasts for  $r=1$ , MSE relative to gross domestic product (GDP) variance in Parts 1 to 5, Part 6 DPCA mean square error (MSE) divided by PCA MSE*

		Nowcast			Forecast					
		Current quarter			One quarter			Two quarters		
	Horizon $h_m$	1	2	3	4	5	6	7	8	9
1. MIDAS-basic	VA-DPCA	0.71	1.01	1.05	1.02	1.19	1.05	1.16	1.24	1.26
	VA-PCA	0.69	1.05	1.02	1.01	1.21	1.03	1.07	1.24	1.33
2. MIDAS-U	VA-DPCA	0.90	1.05	1.02	1.04	1.15	1.11	1.19	1.13	1.17
	VA-PCA	0.76	1.14	1.00	1.00	1.12	1.08	1.15	1.12	1.18
3. MIDAS-smooth	VA-DPCA	0.69	0.92	0.87	0.95	1.10	1.20	1.18	1.12	1.19
	VA-PCA	0.70	0.98	0.88	0.93	1.07	1.12	1.14	1.07	1.17
4. MIDAS-U0	VA-DPCA	0.71	0.86	0.89	0.90	1.05	0.98	1.05	1.09	1.12
	VA-PCA	0.69	0.93	0.93	0.85	1.08	0.91	1.04	1.07	1.13
5. Relative MSE: DPCA/PCA	MIDAS-basic	1.04	0.96	1.04	1.02	0.98	1.01	1.08	1.00	0.95
	MIDAS-U	1.19	0.92	1.02	1.04	1.02	1.03	1.03	1.01	0.99
	MIDAS-smooth	0.99	0.93	0.98	1.02	1.02	1.08	1.04	1.05	1.02
	MIDAS-U0	1.03	0.93	0.96	1.06	0.96	1.08	1.00	1.02	0.99

*Notes:* Parts 1–4 show relative MSEs to variance of GDP. Part 5 shows another relative MSE defined as the MSE of the VA-DPCA factor model divided by the MSE of the model using static factors, denoted as VA-PCA. For model abbreviations, see Table 1.

follow this direction and propose to augment the state-space model by a simple static relationship between monthly GDP growth and the factors. This follows the seminal work by Mariano and Murasawa (2003), where combining monthly and quarterly data in a small factor state-space model has been introduced. Technically, they augment the previous state-space system, see equations (24) and (25), with further relationships that interpolate GDP growth and relate monthly GDP growth to the monthly factors. All in all, they add three equations (see Banbura and Rünstler, 2007, p. 5): equation (1) is  $y_{t_q} = \tilde{y}_{t_q} + \varepsilon_{t_q}$ , with  $\varepsilon_{t_q}$  as a measurement error, which is normally distributed with mean zero and variance  $\Sigma_\varepsilon$ ; (2) an equation for time aggregation  $\tilde{y}_{t_q} = \tilde{y}_{t_m} = (\frac{1}{3} + \frac{2}{3}L_m + L_m^2 + \frac{2}{3}L_m^3 + \frac{1}{3}L_m^4)y_{t_m}^m$  for  $t_m = 3, 6, \dots, T_m$ ; and (3) the static factor representation at the monthly frequency  $y_{t_m}^m = \Lambda_y F_{t_m}$ . Equations in (2) and (3) add to the vector state equation, whereas (1) adds to the vector observation equation of the state-space model. In line with the estimation procedure for the factor-only state-space models (24) and (25), Banbura and Rünstler (2007) estimate the coefficients  $\Lambda_y$ ,  $\Sigma_\varepsilon$  outside the state-space model by estimating a reduced form of (1) to (3), which is a regression model for quarterly GDP growth dependent on time-aggregated quarterly factors. They plug the resulting estimates of  $\Lambda_y$  and  $\Sigma_\varepsilon$  in the state-space model for Kalman smoothing, which now also provides the nowcasts and forecasts for GDP growth, as  $y_{t_q}$  is part of the observation vector in this integrated approach.

The key difference between the two-step factor estimation MIDAS approach chosen in the previous applications and the ones followed by Banbura and Rünstler

(2007) and Mariano and Murasawa (2003) is that MIDAS directly relates time series of different frequencies, whereas the state-space approaches allow for specifying relationships consistently at the higher frequency. Furthermore, MIDAS is a direct forecast device, whereas the Kalman smoother is based on a VAR model that yields iterative forecasts in the terminology of Marcellino *et al.* (2006). This approach is fully integrated as it interpolates missing values of the indicators, estimates factors and yields nowcasts of GDP in one coherent framework. To check whether this strategy can improve over the two-step approach followed here so far in terms of nowcasting and forecasting, we also provide nowcast results for the model proposed by Banbura and Rünstler (2007). Table 6 shows relative MSEs to GDP growth variance and rankings for the different state-space model nowcasts and forecasts. In the table, ‘KFS-PCA full’ denotes the fully integrated approach, whereas all the other forecasts are based on the two-step procedure, where the Kalman smoother is used to estimate the monthly factors only. Note that the coefficients of the state-space model are re-estimated for each recursion in the exercise. Therefore, factor estimates can change owing to parameter changes as well as the addition of new information at the end of the sample. The results show that the integrated approach also does well in nowcasting and forecasting. It performs better than the two-step MIDAS-basic and MIDAS-U projection, and very similar to the simple MIDAS-U0 projection. For horizons two, four and five, it performs best among all the different approaches. For horizons, one and three, the MIDAS-U0 performs best.

The similar performance of the fully integrated state-space model to the very simple MIDAS projections confirms the previous findings that simple and very parsimonious projection models seem to work better than more complicated models. Note that the equation  $y_{t_m}^m = \Lambda_y \mathbf{F}_{t_m}$  in the previous state-space model, which relates

TABLE 6

*Two-step KFS-PCA vs. fully integrated nowcast and forecast results from the state-space model for  $r = 1$ , mean-squared error (MSE) relative to gross domestic product (GDP) variance and ranking*

	Horizon $h_m$	Nowcast			Forecast					
		Current quarter			One quarter			Two quarters		
		1	2	3	4	5	6	7	8	9
1.a. Relative MSE	KFS-PCA full	0.70	0.81	0.84	0.88	1.00	0.95	1.10	1.12	1.09
	MIDAS-basic	0.79	0.90	0.87	1.07	1.17	1.07	1.19	1.13	1.20
	MIDAS-U	0.89	0.90	0.81	0.97	1.03	1.02	1.31	1.49	1.36
	MIDAS-smooth	0.76	0.85	0.89	0.98	1.06	1.08	1.10	1.16	1.19
	MIDAS-U0	0.68	0.85	0.80	0.95	1.01	0.93	1.08	1.09	1.06
1.b. Ranking	KFS-PCA full	2	1	3	1	1	2	3	2	2
	MIDAS-basic	4	5	4	5	5	4	4	3	4
	MIDAS-U	5	4	2	3	3	3	5	5	5
	MIDAS-smooth	3	3	5	4	4	5	2	4	3
	MIDAS-U0	1	2	1	2	2	1	1	1	1

*Note:* For model abbreviations, see Table 1.

monthly GDP growth and the factors, is very parsimonious and does not contain lags of the factors as is the case of the MIDAS-U0 forecast. Whether the approach is integrated within one coherent state-space model or split into two steps is, however, of second-order importance according to our findings. Therefore, we do not seem to lose much if we rely on the two-step procedure, which allows us to compare the different factor estimation methods.

## VII. Conclusions

In this article, we propose mixed-data sampling based on factors, Factor MIDAS, as a nowcasting and forecasting tool that combines methods from the recent literature on large factor models and on mixed-data sampling. Factor MIDAS serves as a projection method that allows for a harmonized comparison of alternative factor estimation methods that can exploit information from a large set of indicators subject to different publication lags that lead to missing values at the end of the multivariate sample, the so-called ragged edge. With Factor MIDAS, high-frequency factors can be used to forecast low-frequency variables; in our case, we use monthly indicators to nowcast and forecast quarterly GDP growth.

We discuss three MIDAS projections – basic, smoothed and unrestricted – and three factor estimation methods that can handle ragged-edge data: DPCAs with data realignment, static PCA with the EM algorithm and the Kalman smoother-based method of Doz *et al.* (2006). It is difficult to rank the set of resulting Factor MIDAS specifications on a purely *a priori* basis, as their performance can differ with specific features of the DGP. Therefore, we have considered their behaviour in an empirical context, where a large set of monthly indicators are used to nowcast and forecast German GDP growth, a key macroeconomic variable for the largest economy in the euro area.

Concerning the differences between the MIDAS projection methods, the results indicate that MIDAS with exponentially distributed lag functions performs similarly to MIDAS with unrestricted lag polynomials. The best performing projection is in many cases a very simple MIDAS with just one lag of the factors. AR dynamics also plays only a minor role in the MIDAS projections.

The choice of the factor estimation technique has no substantial impact on the nowcast performance. All the three methods provide informative nowcasts and to a lesser extent informative forecasts one quarter ahead, and over time the three forecasts are highly correlated. There are also no systematic differences between static and dynamic PCA for nowcasting. Interestingly, factor-forecast applications with single-frequency data, see for example, D'Agostino and Giannone (2006) and Schumacher (2007), have recently obtained similar findings. Moreover, choosing an integrated state-space model as in Banbura and Rünstler (2007) rather than the two-step procedure we propose, cannot improve the nowcast performance.

The methods employed in this article can only be regarded as short-term forecasting devices, as the forecasts of all the factor models we have considered are hardly

informative for forecast horizons longer than one quarter. This feature can often be observed in the recent literature. For example, related to the debate on the ‘Great Moderation’, there is evidence of a decline in forecastability of real and nominal variables for many sophisticated forecasting procedures; see D’Agostino, Giannone and Surico (2006) and Campbell (2007), for example. However, we should also emphasize that virtually all Factor MIDAS nowcasts can improve over quarterly factor forecasts based on time-aggregated data. Thus, taking into account higher-frequency information and exploiting the most recent observations pays off for nowcasting and short-term forecasting. We expect this finding to be valid not only in our application but more generally in the presence of mixed-frequency data, and in this context the Factor MIDAS approach represents a useful tool for empirical analysis.

*Final Manuscript Received: December 2009*

## References

- Aastveit, K. A. and Trovik, T. G. (2007). *Nowcasting Norwegian GDP: The Role of Asset Prices in a Small Open Economy*, Norges Bank Working Paper 2007/9.
- Altissimo, F., Cristadoro, R., Forni, M., Lippi, M. and Veronese, G. (2006). *New Eurocoin: Tracking Economic Growth in Real Time*, CEPR Working Paper 5633.
- Andreou, E., Ghysels, E. and Kourtellis, A. (2010). ‘Regression models with mixed sampling frequencies’, *Journal of Econometrics*, in press.
- Angelini, E., Henry, J. and Marcellino, M. (2006). ‘Interpolation and backdating with a large information set’, *Journal of Economic Dynamics & Control*, Vol. 30, pp. 2693–2724.
- Angelini, E., Camba-Méndez, G., Giannone, D., Rünstler, G. and Reichlin, L. (2008). *Short-Term Forecasts of Euro Area GDP Growth*, ECB Working Paper 949.
- Bai, J. and Ng, S. (2002). ‘Determining the number of factors in approximate factor models’, *Econometrica*, Vol. 70, pp. 191–221.
- Bai, J. and Ng, S. (2006). ‘Confidence intervals for diffusion index forecasts and inference for factor-augmented regressions’, *Econometrica*, Vol. 74, pp. 1133–1150.
- Bai, J. and Ng, S. (2007). ‘Determining the number of primitive shocks in factor models’, *Journal of Business and Economic Statistics*, Vol. 25, pp. 52–60.
- Banbura, M. and Rünstler, G. (2007). *A Look into the Factor Model Black Box: Publication Lags and the Role of Hard and Soft Data in Forecasting GDP*, ECB Working Paper 751.
- Banerjee, A. and Marcellino, M. (2006). ‘Are there any reliable leading indicators for US inflation and GDP growth?’ *International Journal of Forecasting*, Vol. 22, pp. 137–151.
- Banerjee, A., Marcellino, M. and Masten, I. (2005). ‘Leading indicators for euro-area inflation and GDP growth’, *Oxford Bulletin of Economics and Statistics*, Vol. 67, pp. 785–813.
- Banerjee, A., Marcellino, M. and Masten, I. (2008). ‘Forecasting macroeconomic variables using diffusion indexes in short samples with structural change’, chapter 4 in Rapach D. and Wohar M. (eds), *Forecasting in the Presence of Structural Breaks and Model Uncertainty*, Elsevier, Amsterdam, pp. 149–194.
- Barhoumi, K., Benk, S., Cristadoro, R., Den Reijer, A., Jakaitiene, A., Jelonek, P., Rua, A., Rünstler, G., Ruth, K. and Van Nieuwenhuyze, C. (2008). *Short-term Forecasting of GDP Using Large Monthly Datasets: A Pseudo Real-Time Forecast Evaluation Exercise*, ECB Occasional Paper 84.
- Bernanke, B. and Boivin, J. (2003). ‘Monetary policy in a data-rich environment’, *Journal of Monetary Economics*, Vol. 50, pp. 525–546.

- Bhansali, R. (2002). 'Multi-step forecasting', in Clements M. and Hendry D. (eds), *A Companion to Economic Forecasting*, Blackwell, Oxford, pp. 206–221.
- Boivin, J. and Ng, S. (2005). 'Understanding and comparing factor-based forecasts', *International Journal of Central Banking*, Vol. 1, pp. 117–151.
- Campbell, S. (2007). 'Macroeconomic volatility, predictability, and uncertainty in the great moderation: evidence from the survey of professional forecasters', *Journal of Business & Economic Statistics*, Vol. 25, pp. 191–200.
- Chevillon, G. and Hendry, D. F. (2005). 'Non-parametric direct multi-step estimation for forecasting economic processes', *International Journal of Forecasting*, Vol. 21, pp. 201–218.
- Clements, M. and Galvão, A. (2008). 'Macroeconomic forecasting with mixed-frequency data: forecasting output growth in the United States', *Journal of Business & Economic Statistics*, Vol. 26, pp. 546–554.
- Clements, M. and Galvão, A. (2009). 'Forecasting US output growth using leading indicators: an appraisal using MIDAS models', *Journal of Applied Econometrics*, Vol. 24, pp. 1187–1206.
- Clements, M. and Hendry, D. F. (1998). *Forecasting Economic Time Series*, Cambridge University Press, Cambridge.
- D'Agostino, A. and Giannone, D. (2006). *Comparing Alternative Predictors Based on Large-Panel Factor Models*, ECB Working Paper 680.
- D'Agostino, A., Giannone, D. and Surico, P. (2006). *(Un)Predictability and Macroeconomic Stability*, ECB Working Paper 605.
- De Mol, C., Giannone, D. and Reichlin, L. (2008). 'Forecasting using a large number of predictors: is Bayesian regression a valid alternative to principal components?' *Journal of Econometrics*, Vol. 146, pp. 318–328.
- Den Reijer, A. (2005). *Forecasting Dutch GDP Using Large Scale Factor Models*, DNB Working Paper 28.
- Doz, C., Giannone, D. and Reichlin, L. (2006). *A Quasi Maximum Likelihood Approach for Large Approximate Dynamic Factor Models*, ECB Working Paper 674.
- Eickmeier, S. and Ziegler, C. (2008). 'How successful are dynamic factor models at forecasting output and inflation? A meta-analytic approach', *Journal of Forecasting*, Vol. 27, pp. 237–265.
- Forni, M., Hallin, M., Lippi, M. and Reichlin, L. (2003). 'Do financial variables help forecasting inflation and real activity in the euro area?' *Journal of Monetary Economics*, Vol. 50, pp. 1243–1255.
- Forni, M., Hallin, M., Lippi, M. and Reichlin, L. (2005). 'The generalized dynamic factor model: one-sided estimation and forecasting', *Journal of the American Statistical Association*, Vol. 100, pp. 830–840.
- Ghysels, E. and Valkanov, R. (2006). *Linear Time Series Processes with Mixed Data Sampling and MIDAS Regression Models*, University of North Carolina, mimeo.
- Ghysels, E. and Wright, J. (2009). 'Forecasting professional forecasters', *Journal of Business & Economic Statistics*, Vol. 27, pp. 504–516.
- Ghysels, E., Sinko, A. and Valkanov, R. (2007). 'MIDAS regressions: further results and new directions', *Econometric Reviews*, Vol. 26, pp. 53–90.
- Giannone, D., Reichlin, L. and Small, D. (2008). 'Nowcasting GDP and inflation: the real-time informational content of macroeconomic data releases', *Journal of Monetary Economics*, Vol. 55, pp. 665–676.
- Inoue, A. and Kilian, L. (2006). 'On the selection of forecasting models', *Journal of Econometrics*, Vol. 130, pp. 273–306.
- Kapetanios, G. and Marcellino, M. (2009). 'A parametric estimation method for dynamic factor models of large dimensions', *Journal of Time Series Analysis*, Vol. 30, pp. 208–238.
- Kapetanios, G., Labhard, V. and Price, S. (2008). 'Forecast combination and the Bank of England's suite of statistical forecasting models', *Economic Modelling*, Vol. 25, pp. 772–792.

- Koenig, E. F., Dolmas, D. and Piger, J. (2003). 'The use and abuse of real-time data in economic forecasting', *The Review of Economics and Statistics*, Vol. 85, pp. 618–628.
- Marcellino, M. and Schumacher, C. (2008). *Factor-MIDAS for Now- and Forecasting with Ragged-Edge Data: A Model Comparison for German GDP*, CEPR Discussion Paper No. 6708.
- Marcellino, M., Stock, J. and Watson, M. (2005). 'Macroeconomic forecasting in the euro area: country specific versus euro wide information', *European Economic Review*, Vol. 47, pp. 1–18.
- Marcellino, M., Stock, J. and Watson, M. (2006). 'A comparison of direct and iterated multistep AR methods for forecasting macroeconomic time series', *Journal of Econometrics*, Vol. 135, pp. 499–526.
- Mariano, R. and Murasawa, Y. (2003). 'A new coincident index of business cycles based on monthly and quarterly series', *Journal of Applied Econometrics*, Vol. 18, pp. 427–443.
- Matheson, T. (2007). *An Analysis of the Informational Content of New Zealand Data Releases: The Importance of Business Opinion Surveys*, Reserve Bank of New Zealand Discussion Paper Series DP2007/13.
- Schneider, M. and Spitzer, M. (2004). *Forecasting Austrian GDP Using the Generalized Dynamic Factor Model*, OeNB Working Paper 89.
- Schumacher, C. (2007). 'Forecasting German GDP using alternative factor models based on large datasets', *Journal of Forecasting*, Vol. 26, pp. 271–302.
- Schumacher, C. and Breitung, J. (2008). 'Real-time forecasting of German GDP based on a large factor model with monthly and quarterly data', *International Journal of Forecasting*, Vol. 24, pp. 368–398.
- Stock, J. and Watson, M. (2002). 'Macroeconomic forecasting using diffusion indexes', *Journal of Business & Economic Statistics*, Vol. 20, pp. 147–162.
- Stock, J. and Watson, M. (2006). 'Forecasting with many predictors', in Elliot G., Granger C. and Timmermann A. (eds), *Handbook of Economic Forecasting*, Vol. 1, Elsevier, Amsterdam, pp. 515–554.
- Wallis, K. (1986). 'Forecasting with an econometric model: the "ragged edge" problem', *Journal of Forecasting*, Vol. 5, pp. 1–13.
- Watson, M. (2003). 'Macroeconomic forecasting using many predictors', in Dewatripont M., Hansen L. and Turnovsky S. (eds), *Advances in Economics and Econometrics, Theory and Applications, Eight World Congress of the Econometric Society*, Vol. 3, Cambridge University Press, Cambridge, pp. 87–115.

## Appendix: Monthly dataset

The whole dataset for Germany contains 111 monthly time series over the sample period from 1992M1 until 2006M11. The time series cover broadly the following groups of data: prices, labour market data, financial data (interest rates, stock market indices), industry statistics, construction statistics, surveys and miscellaneous indicators.

The source of the time series is the Bundesbank database. The download date of the dataset is 6 December 2006. In this dataset, there are differing missing values at the end of the sample. For example, whereas financial time series are available up to 2006M11, industrial time series like production, orders and so on are only available up to 2006M09. This leads to a ragged-edge structure at the end of the sample, which serves as a template to replicate the ragged edges in past pseudo real-time periods as described in the main text.

Natural logarithms were taken for all time series except interest rates. Stationarity was obtained by appropriately differencing the time series. Most of the time series



taken from the aforementioned source are already seasonally adjusted. Remaining time series with seasonal fluctuations were adjusted using Census-X12 prior to the forecast simulations. Correction of extreme outliers was performed using a modification of the procedure proposed by Watson (2003). Large outliers are defined as observations that differ from the sample median by more than six times the sample interquartile range (Watson, 2003). The identified observation is set equal to the respective outside boundary of the interquartile. The entire list of time series can be found in appendix A of Marcellino and Schumacher (2008).