

Virk, Nader; Javed, Farrukh; Awartani, Basel

Working Paper

A reality check on the GARCH-MIDAS volatility models

Working Paper, No. 2/2021

Provided in Cooperation with:

Örebro University School of Business

Suggested Citation: Virk, Nader; Javed, Farrukh; Awartani, Basel (2021) : A reality check on the GARCH-MIDAS volatility models, Working Paper, No. 2/2021, Örebro University School of Business, Örebro

This Version is available at:

<https://hdl.handle.net/10419/244577>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



WORKING PAPER

2/2021

A reality check on the GARCH-MIDAS volatility models

Nader Virk, Farrukh Javed and Basel Awartani

Statistics

ISSN 1403-0586

<https://www.oru.se/institutioner/handelshogskolan/forskning/working-papers/>

Örebro University School of Business
701 82 Örebro
SWEDEN

A reality check on the GARCH-MIDAS volatility models

Nader Virk[¬], Farrukh Javed^{*} and Basel Awartani[~]

Abstract

We employ a battery of model evaluation tests for a broad-set of GARCH-MIDAS models and account for data snooping bias. We document that inferences based on standard tests for GM variance components can be misleading. Our data mining free results show that the gains of macro-variables in forecasting total (long run) variance by GM models are overstated (understated). Estimation of different components of volatility is crucial for designing differentiated investing strategies, risk management plans and pricing of derivative securities. Therefore, researchers and practitioners should be wary of data mining bias, which may contaminate a forecast that may appear statistically validated using robust evaluation tests.

Key words: GARCH-MIDAS models, component variance forecasts, macro-variables, data snooping

JEL classification: C32, C52, G11, G17,

[¬] Corresponding author. Plymouth Business School, PL4 8AA Plymouth, UK. Email: nader.virk@plymouth.ac.uk

^{*} Örebro University School of Business, SE-701 82 Örebro, Sweden, Email: farrukh.javed@oru.se

[~] Westminster Business School, 35 Marylebone, London, NW15LS, Email: b.awartani@westminster.ac.uk

The new class of two component volatility models pioneered by Engle, Ghysels and Sohn (2013) incorporates a long run variance component in the modelling of volatility.¹ The novel GARCH-MIDAS (GM) models combine transitory/short-term GARCH volatility with the variability of the low frequency information channels, predominantly economic sources. This model is suitable for modelling short and long-term volatilities in financial markets.²

Engle et al. (2013) use US equity and macroeconomic data to provide the first evidence for the utility of these models in forecasting volatility. They document that the GM total variance forecasts benefit from including information from aggregate macro variables. Asgharian, Hou and Javed (2013) and Conrad and Loch (2015) provide additional tests for volatility forecasting using GM models while expanding the set of economic variables. The evidence in both studies reinforce results in Engle et al. (2013) and document that GM models, which bundle common variability of several economic variables, using principal component (PC) analysis can enhance the forecasts of the GM for both the total and the long term variance. Conrad and Loch (2015) report that term spread, housing starts and consumer sentiment are the leading predictive variables for the US stock volatility.³ In forecasting the long run US equity volatility, Lindblad (2017) shows that there gains for adding sentiment variables, whereas Conrad and Kleen (2020) reinforce the fact that housing starts is an excellent predictor for stock volatility. In sum, the evidence suggests that there are clear advantages for using macroeconomic information in predicting volatility in the context of GM models.

In the literature many researchers use GM models in volatility prediction. Their work has been motivated by the lack of appropriate models that combine the short and the long components of the latent volatility process.⁴ However, there are three issues/gaps that need scrutiny when examining the overall gains of modelling volatility within the class of GM models. First, most GM volatility modelling evidence uses US equity data. There has been little, or no, application and evidence of GM models for volatility modelling from other global markets. Second, the emerging evidence displaying the utility of GM models compares the forecasts of GM with macro information with the predictions of a GM model that is a function of lagged monthly realized variance (RV).⁵ RV, though unbiased and consistent, is still a noisy estimate of ex-post volatility. Andersen and Bollerslev (1998) illustrate that lower frequency forecasts of the latent volatility may benefit from constructing factors at higher frequencies to mitigate the idiosyncratic noisiness of the RV measure.⁶ Conjecturing from their evidence, we hypothesise that the use of cumulative monthly RV, computed from daily squared returns may provide better information container than is usually acknowledged in forecasting different variance components, especially at monthly and above lower time frequencies.

¹ The multiplicative component structure of these models combines the GARCH volatility with the long run approximations for long run variance – estimated by the innovative Mixed Data Sampling (MIDAS) approach of Ghysels, Santa-Clara and Valkanov (2004, 2006). Given the importance of long-memory dependence in financial market volatility, the long run variance smoothing by the MIDAS regressions – forecasting mismatched time frequency (e.g. monthly) variance from high frequency data points such as past daily squared returns or economic variables – has opened up new frontiers in examining the role of long memory processes in shaping financial volatility and how cross-market long run variances and correlations depict inter- and intra-market integration patterns.

² We invariably refer long-term variance to trend/secular/MIDAS component or variance to imply the same throughout the paper.

³ In addition to modelling conditional variance, the extensions build on the GM framework have been utilised in studying the relationships between oil and stock market volatilities, oil-stock correlation, stock-bond correlation, oil-macroeconomic relationships, European equity market integration patterns and the effect of investor sentiment on US stock-bond correlation patterns (Conrad, Loch and Rittler 2014; Asgharian, Christiansen and Hou 2016, Pan, Wang, Wu, and Libo 2017, and Virk and Javed 2017, Fang, Yu and Huang 2018).

⁴ There are numerous stylized facts documented about financial market volatility, such as clustering, leverage effect, mean reversion and co-movement of volatilities among assets and across markets.

⁵ Potentially, daily RV can be is a noisy estimate for forecasting volatility at high frequencies (Andersen and Bollerselv 1997, 1998). However, we argue that the same is not applicable to the low frequency RV estimates when we model low frequency variance components at monthly, quarterly and bi-annual along with the total variance forecasts.

⁶ To this extent, Andersen and Bollerslev (1998) show that a well-specified, GARCH-type volatility filter smooth unconditional realised variance/volatility (RV), from high frequency intraday data, produces precise inter-daily predictions i.e. a latent volatility factor from high frequency intraday data improves out-of-sample forecasts at inter-daily variance predictions.

Thus, merging our conjecture with the findings from the prior empirical evidence that suggest a weak relationship between stock volatility and aggregate variables (Officer 1973, Shiller 1981a&b, Christie 1982, Schwert 1989, among others), leads us to question the evolving evidence on GM models: are there actual GM specifications that perform better than the common GM benchmark?⁷ More succinctly, we inquire if (i) the evidence from GM models with macro variables is actually differentiable from the reported evidence and (ii) if the forecasts that incorporate macro information really outperform the forecasts of the GM model that condition on lagged RV to approximate the latent volatility.

We identify that the forecasting exercises in the studies using GM models typically use weak, yet informative tests, to drawing inferences about the relative performance of the GM variance forecasts against the benchmark models.⁸ These tests mostly use mean squared error (MSE) type loss function. In the spirit of Engle et al. (2013), an informative/relative MSE ratios test inform the efficiency gains of the competing model against the benchmark model or vice a versa. However, Asgharian et al. (2013) have assessed the forecasting merits of GM model using the Diebold and Mariano (DM, 1995) test – a robust statistical test for forecast evaluation.⁹ We also note that prior studies, with the definite exception of Asgharian et al, have usually generalized the evidence coming from GM models for total volatility to long run variance forecasts.

Finally, we observe that given the variety of economic sources as well as financial information variables and for that there are several candidate GM models to forecast volatility. The repetitive use of same equity data is bound to cast a shadow over the reliability and accuracy of the GM forecasting volatility inferences. Because of the data mining biases, we are sceptical about the reliability of the prior evidence on the forecasting efficiency of GM models particularly as it does not clearly differentiate between total variance and long run variance forecasts. In order to provide rigour to these forecasting comparisons using the GM models, we argue that the model comparisons of GM models should account for statistical tests that control for data snooping biases before inferring on the value of GM macro models in volatility prediction.

These identified research gaps motivate our work. To limit our study, we note that Engle et al. (2013) specify two types of GM models. The first type the one-sided filters where long-term volatility of daily stock returns is expressed as a weighted average of lagged values of lower-frequency financial/macroeconomic variables. The other type of models are the two-sided filters that use lagged and future values of the MIDAS input variables. Our work only studies the first type of GM models. We provide new evidence in two ways. First, we replicate research using the one-sided GM models across four leading global equity markets i.e. France, Germany, the UK and the US. Second, we assess the robustness of earlier evidence using tests that account for data snooping bias..

To do this we investigate the gains of several macro-variables in the GM models to forecast total and long run variance components over the use of a common GM benchmark that only uses monthly RV. We carry out informative and robust statistical tests that correspond with prior research on GM models. Furthermore, we use a far larger set of macroeconomic variables than all previous studies that assessed the importance of information coming from different aggregated channels. The study by Conrad and Loch (2015) comes closest to ours in this respect. Nonetheless, our work extends evidence on GARCH-MIDAS modelling and forecasting for a cross-section of four developed equity markets and thus provide out of sample evidence.

⁷ Although the importance of economic sources influencing market volatility is undeniable, there are limitations in the modelling the contribution of economic variables to the total and long run variance forecasts. These include identification of the aggregate variable(s) that contributes to the evolution of financial market variance and prediction of financial market volatility on the basis of an information set that is prone to measurement errors and revisions compared with other variables to proxy long run variations in the conditional variance such as RV.

⁸ This applies to tests for both total variance and long run variance components regardless benchmark model is the GM model that smooth realized variance in MIDAS regressions (Engle et al. 2013 and Conrad and Loch 2015) or the baseline GARCH (1, 1) specification (Asgharian et al. 2013 and Conrad et al. 2015).

⁹ The pairwise DM test examines the equal predictability (EPA) of the alternate model against the benchmark.

Second, we employ powerful (multiple/joint) tests that account for the data mining issues i.e. the White (2000) reality check (RC) test and Hansen (2005) superior predictive ability (SPA) test. Both tests assess the SPA of the benchmark model when contrasted against a host of alternative models and are robust against data mining biases. Specifically, these tests evaluate the amount of increase in the probability of finding SPA among the competing models when the number of competing model increases (White 2000 and Hansen 2005).¹⁰ Our evidence accounting for the data snooping bias for the class of GM models with one benchmark, is the first in this regard. Here it is important to note that, contemporary work by Lindblad (2017) and Conrad and Kleen (2020) have made use of model confidence set (MCS) of Hansen et al. (2011) for forecast evaluation of GM models.¹¹ The evidence in latter study shows that the GM RV model is outperformed by models that use macroeconomic information. Furthermore, the lag structure for the input MIDAS variables is given by the weighting scheme – unconstrained or constrained – adopted for the beta polynomial by which MIDAS variance evolves. Thus, the estimation of long run variance is dependent on the weights given to the lagged values of the input variable(s) in the MIDAS filter. Conrad and Loch (2015) confirm the evidence in Engle et al. (2013) that the optimal weighting scheme for the GM-RV (benchmark) model is the restricted one. The former study does show that some macroeconomic variables require an unrestricted scheme while RV does not.

Nonetheless, there are studies that adopt an ad-hoc restricted weighted scheme even for macroeconomic variables in the MIDAS filter such as Asgharian et al. (2013) and Conrad et al. (2014) among others. We expect that imposing this weighting scheme for other low frequency input variables may not optimally forecast long run variance given the flexible requirement for weight convergence for them as shown in Conrad and Loch (2015). To add value to our empirical analysis, we estimate all GM models using both unrestricted and restricted weighting schemes for the MIDAS variance component across cross-section of four developed equity markets – just not the US equity volatility.

In summary, we carry out a large-scale empirical exercise for the class of GM models, using one-sided filters, across four large global equity markets. Through this exercise we investigate the gains and efficacy of volatility forecasts when total and long run volatility evolves, using macro information in the MIDAS filter relative to the benchmark model, while controlling for the data snooping biases. We follow Engle et al. (2013) and make the GM model that uses rolling window (RW) monthly RV in the MIDAS filter as the benchmark model, referred as GM-RV model hereafter.

Our analysis covers daily equity and monthly macro data from 1999 to 2016. For every market, we estimate 27 GM models resulting in the same number of forecasts for the total variance and trend component. This number is doubled since we use two weighting schemes in the MIDAS regressions. Using the informative MSE ratio test of Engle et al. (2013) and the Diebold and Mariano (1995) test, our results are congruent with the findings from the US equity data. All comparisons use one step ahead forecasts i.e., pseudo out-of-sample (POS) GM variance forecasts (unless otherwise stated).

Our results show that the evidence using the MSE ratio and DM test overstates the gains of the total variance forecasts coming from GM macro models when contrasted against the benchmark model forecast: there are several competing models that outperform the benchmark model's MSE. The DM test concurs when it uses an unconstrained weighting scheme. However, with the restricted weighting gains are limited to forecasts for the MIDAS/long run variance component only –

¹⁰ Here we note that the model comparisons in our study are for long run variance and total variance as provided by the GM models. The comparability tests for volatility predictive ability tests for short run variance of GARCH type models have already been studied extensively, see Lunde and Hansen (2005) and Gonzalez, Lee and Mishra (2004) and others, and therefore we refrain from reporting that evidence to conserve space.

¹¹ As mentioned earlier, initial work on GM modelling has used various pairwise tests e.g. Diebold-Mariano tests. Recent work on GM modelling and forecasting is transcending to evaluate its gains using stronger statistical tests. For example, Lindblad (2017) relies on the model confidence set approach of Hansen et al. (2011). Similarly, Conrad and Kleen (2020) also use MCS test for GM models. However, both the studies use only the US data.

an aspect that is typically overlooked in the GM model evaluation literature. These summary results from the DM test largely hold across all markets as we account for data snooping problems. The only exception is the fact that joint tests show that no competing model has SPA over the benchmark model total variance forecasts regardless of what weighting scheme is adopted.

We examine the consistency of our results by using alternate variance proxies and an out-of-sample forecasting scheme.¹² The generality of our results for the POS is maintained using alternate variance proxies. The results for the OS rolling forecasting comparisons – using any type of variance proxy – confirm our POS evidence. There are exceptions, however. We note that with the OS forecasting procedures the total (long run) variance forecast evaluations are sensitive to what type of weighting scheme is adopted, which endorses our scepticism on model comparison exercises using GM models with restricted weights. Our evidence shows that results from a particular weighting scheme cannot be generalized for forecasts coming from GM models with other weighting structures which is in accordance with Conrad and Loch (2015), where the authors estimated all models with a restricted and an unrestricted weighting scheme and used the likelihood ratio tests to identify the appropriate specification. This also applies to using alternative forecasting schemes such as one step ahead or multiple step ahead forecasts. Most importantly, we should scrutinise the usefulness of the GM with respect to its two distinctive forecasts i.e. total variance and long run variance – neither of the two can substitute for the other. Finally, the use of powerful tests is suggested to examine the forecasting gains of a particular model over the benchmark model, e.g. when p-values of robust statistical tests reject the superiority of the forecasting performance of the benchmark model.

Overall, we conclude that, although macroeconomic information may not improve the accuracy of total GM volatility forecasts, there is overwhelming evidence of its usefulness in improving the long-term equity volatility forecasts. Nonetheless, guided by the sum of our results from data mining free POS and OS forecasting comparisons, we broadly find that the existing literature over(under)-states the gains for forecasting GM total (long run) variance with aggregate variables. Our results have value academics, investors and practitioners alike when we know that financial volatility has different components. The emphasis is on the econometrician to be rigorous and liberal in their search for forecasts for the different components of conditional financial volatility: if GM-RV model has SPA over competing models its does not translate into superiority for the low frequency variance component and without limiting the evolution of the secular variance part. We know investors, money and risk managers with heterogeneous investment and risk management needs require volatility forecasts for different time horizons, investment mandates, and geographical presence as well as informational flow heterogeneity.

The rest of the paper is organised as follows. Sections 2 and 3 discuss the GM methodology and the choice of latent variance proxy in the scope of our work. Section 4 details the forecasting comparison testing procedures. Section 5 specifies data, sources and summary statistics. The results are summarized in section 6. And section 7 concludes the findings and implications of work.

2 GARCH-MIDAS models

In this section, we outline the GARCH-MIDAS methodology. The two-component GM volatility models break down the total variance for a financial asset into a short-term transitory component and trend/secular component. The multiplicative total variance is modelled by a unit variance GARCH (1, 1) and a secular component, which is estimated by the MIDAS filter.¹³

¹² The POS forecasting comparisons are one step ahead only and can be taken equivalent to fixed forecasting scheme, while our OS joint tests give multiple step ahead forecasts. More factually one year at time, for details see section 6.4.

¹³ The estimates for secular volatility can be obtained through several economic, financial and sentiment related variables.

To set up the model notations, we assume the compounded return for a price series on day i in month t is $r_{i,t} = \mu + \sqrt{\tau_t g_{i,t}} \varepsilon_{i,t}$ where $\varepsilon_{i,t} | \Phi_{i-1,t} \sim N(0,1)$ and $\Phi_{i-1,t}$ is the information set until day $i - 1$ in period t (month in our case). The compounded return series have mean μ as its location parameter and its scale parameter, i.e. the total conditional variance $\sigma_t^2 = \tau_t g_{i,t}$ comprises a transitory (GARCH) component $g_{i,t}$ and a long-run (MIDAS) component τ_t .

Engle et al. (2013) specifies the short-term volatility component as a GARCH (1, 1) process:

$$g_{i,t} = (1 - \alpha - \beta) + \alpha \frac{(r_{i-1,t} - \mu)^2}{\tau_t} + \beta g_{i-1,t} \quad (1),$$

where $\alpha > 0$, $\beta \geq 0$ and $\alpha + \beta < 1$.

Here, the long-term volatility process, τ_t , can evolve by a range of low frequency variables such as macroeconomic variables. Under the MIDAS setting, this component is estimated:

$$\tau_t = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_1, w_2) X_{t-k}, \quad (2)$$

where X is any variable of interest, however, just as in any another regression, it can accommodate more than one variable as well, around which the secular component should be determined:

$$E_{t-1}(r_{i,t} - \mu)^2 = \tau_t E_{t-1}(g_{i,t}) = \tau_t.$$

To estimate the long-term variance at t , $\phi_k(w)$ is a smoothing function that provides a weighting scheme for the lagged values of the variable(s) of interest in the MIDAS filter. We choose the beta smoothing function that determines optimal weights with which the lags of the input variables in the MIDAS regressions are going to shape the secular volatility through the parameter θ_1 . It is specified:

$$\phi_k(w_1, w_2) = \frac{(k/K)^{w_1-1} (1-k/K)^{w_2-1}}{\sum_{j=1}^K (j/K)^{w_1-1} (1-j/K)^{w_2-1}} \quad (3)$$

where w_1, w_2 are weights to be estimated. Using the flexible/unrestricted beta smoothing function, the long-term volatility of daily returns in equation (2) is expressed as a weighted average of lower-frequency financial and/or macroeconomic variables. This beta-polynomial is independently estimated for each MIDAS regression and for each input variable therein. Studies have used restricted version of the above weighting scheme by fixing $w_1 = 1$ (Engle et al. 2013, Asgharian et al.

2013, Conrad et al. 2014, among others): $\phi_k(w_2) = \frac{(1-k/K)^{w_2-1}}{\sum_{j=1}^K (1-j/K)^{w_2-1}}$. Ghysels, Sinko and Valkanov (2007) report that the

unrestricted smoothing scheme allows for a hump-shaped decaying pattern. For restricted weighting scheme, the fixed weight of $w_1 = 1$ ensures a decaying pattern whereas the size of w_2 determines the speed of decay: large (small) values of w_2 generate an accelerating (decelerating) decaying pattern for the lagged values of input variable(s) in the MIDAS filter. It important to note that the unrestricted case is flexible in providing/estimating large weights for distant lags of MIDAS input variable (i.e. the hump shape decline), fixing $w_1 = 1$ makes it a strictly declining case for the number of lags involved in the MIDAS smoothing.

To clarify, we note that the parameter space for the baseline GM model using RV as input variable results in $\Theta = \{\mu, \alpha, \beta, \theta_0, \theta_1, w_2\}$ when we use the restricted weighting scheme. It is understood that the parameter space for the GM-RV model with the unrestricted weighting scheme will result in $\Theta = \{\mu, \alpha, \beta, \theta_0, \theta_1, w_1, w_2\}$ i.e. one more parameter estimate than the restricted case that fixes $w_1 = 1$. The parameter space changes accordingly for GM specifications that have more than one exogenous variable to smooth secular component. For example, the unrestricted parameter space for two exogenous variables in the MIDAS regression will become $\Theta = \{\mu, \alpha, \beta, \theta_0, \theta_1, \theta_2, w_{1,\theta_1}, w_{2,\theta_1}, w_{1,\theta_2}, w_{2,\theta_2}\}$. Analogously, the restricted version will contain $2p$ less parameters, where p represents the number of input variables in the MIDAS

regression. In the scope of our work, we estimate all the models using both weighting schemes. First, we let each MIDAS regression search for the properties of exogenous variable(s) in the MIDAS regressions using the flexible weighting scheme i.e. estimate w_1 and w_2 . We then follow this by replicating all the MIDAS regressions using a constrained weighing scheme. For clarity, we specify the baseline long run component that uses rolling window (RW) monthly-realized volatility available at daily frequency:

$$\tau_{i,t} = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_1, w_2) RV_{i,t-k} \quad (4)$$

where rolling window $RV_{i,t} = \sum_{i=1}^N r_i^2$, $N = 22$ approximate monthly realized volatility, and K lags of the input variable(s), are utilized to smooth trend component.¹⁴ The relation in equation (4) can also be specified as

$$\log(\tau_t) = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_1, w_2) RV_{t-k} \quad (5).$$

Following Engle et al. (2013), we adhere to the log-specifications for all the models: a log version of GM-RW RV or simply GM-RV is directly comparable to the competing GM specifications that involve macroeconomic variables. The log transformation of equation (4) guarantees the non-negativity of the conditional variances when the input variables (e.g. term structure of interest rates, industrial production) can take negative values. The transformation of $\log(\tau_t)$ specification gives τ_t which is an exponential estimate of the right side of the equation (see Engle et al. 2013 eq. 19 on page 781).

In a similar fashion, one can construct the long-term volatility component using the levels of macroeconomic variables:

X_{t-k}^l denotes the level of X input variable in the MIDAS filter:

$$\log(\tau_t) = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_1, w_2) X_{t-k}^l \quad (6).$$

Or further extend this model by incorporating the variance of macroeconomic information as well:

$$\log(\tau_t) = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_{1,\theta_1}, w_{2,\theta_1}) X_{t-k}^l + \theta_2 \sum_{k=1}^K \phi_k(w_{1,\theta_2}, w_{2,\theta_2}) X_{t-k}^v \quad (7)$$

where X_{t-k}^v is the variance of X input variable in the MIDAS filter. For generality, a specification by adding RV together with the macroeconomic information, both at level and variance is

$$\log(\tau_t) = \theta_0 + \theta_1 \sum_{k=1}^K \phi_k(w_{1,\theta_1}, w_{2,\theta_1}) RV_{t-k} + \theta_2 \sum_{k=1}^K \phi_k(w_{1,\theta_2}, w_{2,\theta_2}) X_{t-k}^l + \theta_3 \sum_{k=1}^K \phi_k(w_{1,\theta_3}, w_{2,\theta_3}) X_{t-k}^v \quad (8)$$

Using the models implied by equations 5-8, the 26 MIDAS and GM log specification models are iterated to compute model forecasts for comparison with the benchmark GM-RV model. For simplicity and differentiation, we notate the benchmark model that involves RW-RV as GM-RV model and all remaining (competing M) models are referred by GM models.

3 Model Forecasts

In this section, we discuss our forecasting strategy and the unbiased nature of RV to proxy latent variance in predicting future increments in quadratic variation. Typically, when carrying out volatility forecast evaluations, researchers split the sample into the estimation and prediction samples where the estimation sample contains Q observations and the prediction sample contains P observations, and they both add up to the total number of observations in the sample T , i.e., $T = Q + P$. P is the number of observations that are reserved for out-of-sample comparisons.

Since volatility is a latent variable this makes the selection of an evaluation criterion arbitrary when assessing which loss function is appropriate to evaluate volatility models, for detailed discussion, see Bollerslev et al. (1994), Diebold and Lopez (1996) and Lopez (2001). In our case, we maintain consistency across the several tests carried out in our work by applying

¹⁴ In our baseline specification, we use RW-RV in the MIDAS smoothing at daily frequency: each daily realised variance is the rolling sum of 22-daily squared returns. Whereas, the long run variance smoothing for all other GM specifications including macro variables and/or principal components is at monthly frequencies only: long run variance component changes at monthly frequency and stays constant for the days in a month.

the MSE based loss function, a metric that has more often been utilized in GM forecast evaluation tests. These tests include informative comparability ratios and pairwise comparison tests. Finally, to evaluate the closeness of the model forecasts, we require a valid measure of latent volatility, so we proxy the unobservable volatility using the realized volatility. The ex-post monthly realized volatilities of the index returns are computed by adding the daily squared returns during the trading days of the month. The total conditional daily volatility of the GM models is compared against daily squared continuously compounded returns.

Assuming that the variance proxy converges to the true volatility, both at daily and monthly frequencies¹⁵, therefore, despite the noisiness of squared returns, the proxy is still suitable if a quadratic loss function is used to evaluate the forecasts.¹⁶ It can be shown that under a quadratic loss, a correct inference on the relative accuracy of the models can be obtained. For instance, suppose that σ_t^2 , $\sigma_{t,GM-RV}^2$, and $\sigma_{t,m}^2$ are, respectively, true latent variance, the total variance forecasts for the true latent volatility from the benchmark GM-RV model and competing M GM models. For the ease of the discussion, we drop the daily i subscript.¹⁷ Moreover, assume that the proxy for the unobservable volatility σ_t^2 is defined by squared returns r_t^2 . Hence, we may write for a competing m GM forecast:

$$\begin{aligned} E\left((r_t^2 - \sigma_{t,k}^2)^2\right) &= E\left(\left((r_t^2 - \sigma_t^2) - (\sigma_{t,k}^2 - \sigma_t^2)\right)^2\right) \\ &= E((r_t^2 - \sigma_t^2)^2) + E((\sigma_{t,k}^2 - \sigma_t^2)^2) \end{aligned}$$

Here the expectation of the cross product is zero due to the independence between the competing GM model's error with the squared returns proxy error. Re-writing the above by including the benchmark forecast yields,

$$\frac{1}{T} \sum_{t=1}^T \left((r_t^2 - \sigma_{t,GM-RV}^2)^2 - (r_t^2 - \sigma_{t,k}^2)^2 \right) \xrightarrow{p} E((\sigma_t^2 - \sigma_{t,GM-RV}^2)^2) - E((\sigma_t^2 - \sigma_{t,k}^2)^2)$$

where the left-hand side is only negative when

$$E((\sigma_{t,GM-RV}^2 - \sigma_t^2)^2) < E((\sigma_{t,k}^2 - \sigma_t^2)^2)$$

This proves that the squared returns proxy leads to a correct inference about the relative accuracy of the competing GM models, when economic information is incorporated, compared to the basic GM-RV benchmark model.

In reality, the parameters and volatility forecasts of the models are not known and have to be estimated using, as noted earlier, a sample of Q data. If the estimated parameters are \sqrt{Q} consistent and $\frac{P}{Q}$ converges to $\pi < \infty$, then

$$\frac{1}{P} \sum_{t=Q}^{Q+P-1} \left((r_{t+1}^2 - \hat{\sigma}_{t+1,GM-RV}^2)^2 - (r_{t+1}^2 - \hat{\sigma}_{t+1,k}^2)^2 \right) = \frac{1}{P} \left((r_{t+1}^2 - \sigma_{t+1,GM-RV}^2)^2 - (r_{t+1}^2 - \sigma_{t+1,k}^2)^2 \right) + O_P\left(P^{-\frac{1}{2}}\right)$$

and hence, the comparison of the models is still valid as the ranking is preserved. For simplicity, we annotate the total and long run variance forecast comparisons using MSE based quadratic loss function:

$$RV_{MSE}^{\sigma^2} \equiv \frac{1}{P} (r_{i,t}^2 - \sigma_{i,t,m}^2)^2 \quad \text{and} \quad RV_{MSE}^{\tau} \equiv \frac{1}{P} (RV_t - \tau_{t,m})^2.$$

¹⁵ Note that the monthly measure here is only valid if the time series of returns is identically and independently distributed with zero mean. In that case, the monthly volatility is the mere sum of daily volatilities.

¹⁶ A valid proxy that is widely used to measure the accuracy of volatility models is the daily realized volatility computed by summing intraday squared returns (See, Andersen et al., (2001a, 2001b); Barndorff-Nielsen and Shephard (2001, 2002); Meddahi, 2002 among others). It has been shown in these studies that as we compute over smaller and smaller intra-day intervals, it converges to the true integrated volatility process in continuous semi-martingale process. The measure has been also shown to be unbiased by Hansen and Lund (2005). However, the GM is a discrete time process and realized volatility computed over intraday data may not be consistent and hence, there isn't much to choose between realized volatility and squared returns as both are unbiased and neither is consistent.

¹⁷ Here, we assume that the volatility process is covariance stationary and that we know the parameters of the MIDAS models.

where $RV_{MSE}^{\sigma^2}$ and RV_{MSE}^{τ} are the loss functions computed for total and long run variance, respectively for the competing GM models when latent variance is approximated by RV.

4 Model Comparisons

In our work, the conditional total and long run variance forecasts from the competing M GM models are compared against the GM-RV forecasts, respectively $\sigma_{t,GM-RV}^2$ and $\tau_{t,GM-RV}$. The parameters of the volatility models are estimated using the Q sample observations. These estimates are then used to make one step ahead pseudo out-of-sample forecasts, in our case: $P = 1$, which then are taken to different model comparisons tests.¹⁸ The model comparison tests are undertaken for both the total variance forecasts at daily frequency and the long run variance forecasts at monthly frequency. The forecasting errors are calculated with respect to monthly RV, from daily stock return series for each market for the secular components and using squared returns for the total variances. So, when evaluating model forecasts, we proxy latent conditional variances by squared returns for daily data, and use the sum of daily squared returns in month t i.e. $\hat{\sigma}_t^2 = RV$ for monthly data.

Consider that there are M competing models: $M = 1, 2, \dots, m$ and GM-RV is the benchmark model. Each model m provides a forecast or series of forecasts $\{g_{m,t}, \tau_{m,t}\}_{t=1}^P$ which are compared to $\{\hat{\sigma}_t^2\}_{t=1}^P$, yielding a MSE based loss function L . Each model leads to a sequence of losses/forecast errors i.e. the losses for benchmark model are defined as, $L_{GM-RV,t} \equiv L(\hat{\sigma}_t^2, g_{GM-RV,t}, \tau_{GM-RV,t})$ and losses for a competing k model are defined as, $L_{m,t} \equiv L(\hat{\sigma}_t^2, g_{m,t}, \tau_{m,t})$. Using these quadratic losses, we conduct numerous model comparisons tests that are briefly explained below.

4.1 MSE ratios

This informative forecasting comparison metric is presented in Engle et al. (2013) to assess model performance of the GARCH-MIDAS models. The statistic, which is the ratio of the MSE of conditional variance forecasts of the m -competing model e.g. $RV_{m,MSE}^{\sigma^2}$ relative to MSE of the total variance forecast of the benchmark model i.e. GM-RV model, is defined as,

$$MSE \text{ ratio} \equiv RV_{m,MSE} / RV_{GM-RV,MSE}, \text{ for } M = 1, 2, \dots, m.$$

These ratios are computed for both long-term and the total variances of all the competing models with reference to corresponding MSE of the benchmark model. Given that the benchmark model is in the denominator, a ratio of less than one implies an improvement over the benchmark model. . In other words, a ratio lower than 1 indicates a lower loss of accuracy in the sample of competing forecasts compared to the GM-RV benchmark model.

4.2 Pairwise tests of equal predictive ability

Tests for equal predictive ability (EPA), in a general setting, were proposed by Diebold and Mariano (1995) and West (1996). The DM test is used to compare the prediction accuracy of two competing models with the null hypothesis of no difference between the accuracy of two competing forecasts.

The relative performance (RP) measure between the loss function, which in our case is quadratic, of two competing models is:

$$RP_{k,t} \equiv L_{m,t} - L_{GM-RV,t}, \text{ for } M = 1, 2, \dots, m \quad (9)$$

While the test proposed by Diebold and Mariano (1995) focuses on EPA and is pairwise, testing whether a particular forecasting procedure is superior to the alternative forecasts requires a test of superior predictive ability (SPA). Therefore,

¹⁸ The POS, which could also be regarded in-sample fitting, is a terminology that follows from Engle et al. (2013) and is adopted for ease of comparison.

we carry out robust tests of superior predictive ability proposed by White (2000) (hereafter, RC test) and Hansen (2005) (hereafter, SPA test). These tests are explained below.

4.3 Tests for superior predictive ability

Utilizing the terminologies set earlier, the null hypothesis under both White's (2000) RC test and Hansen's (2005) test states that among competing models, the one with the smallest loss i.e. $L_{m,t}$ is not any better than the losses given by benchmark model i.e. $L_{GM-RV,t}$. However, rejecting the null hypothesis means that at least one model produces smaller forecasting errors compared to the benchmark. The loss functions, $L_{m,t}$ and $L_{GM-RV,t}$ for all the models are taken together via the relative measure as shown in equation (9), i.e. multiple testing is carried out to examine the superiority of the benchmark relative to the best performing model in the competing space or vice versa.

Using numerical results, Hansen (2005) shows that White's reality check test is conservative at detecting the SPA of competing models relative to the benchmark model forecasts and aggressively favours the null hypothesis (the benchmark forecasts) when the alternate model space contains poor and/or irrelevant forecasts. Hansen articulates this drawback and states that the p-value can be spuriously increased when inferior models are present in the spectrum of models. This is due to the variance of the loss function of a poorly specified m -model $\bar{L}_{m,t}$, which may remain large even after the inclusion of better models. He further proposes two modifications to the original RC test, the first uses studentized test statistic and the second specifies a sample-dependent null distribution for the test statistic.

Hansen (2005) document that the approximate distribution of the Hansen's SPA test-statistic is obtained using the stationary bootstrap of Politis and Romano (1994), similar to RC SPA test-statistic. Nonetheless, the resultant test-statistic has greater power and is less sensitive to poor or irrelevant specifications resulting in the non-rejection of the null hypothesis. Hansen (2005) considered different adjustments to equation (9) to get the three versions of the studentized test statistics that depend on the variance of $\bar{L}_{m,t}$.

5 Data

Our data set contains the market MSCI equity indices (obtained from DataStream) of four major global markets: France, Germany, the UK and the US. The daily and monthly dollar prices of these indices are retrieved from Thomson Reuters DataStream for the period from January 1994 to December 2016. We compute continuously compounded returns: $r_i = \ln(P_i) - \ln(P_{i-1})$ for all four equity indices, where P represents the daily closing index price levels and i is the day index. The equity markets considered in our study provide consistency and a coherent developed market perspective and therefore, are crucial for global investors, portfolio managers and institutional investors considering the criteria of investibility, replicability and cost efficiency. This is displayed by the fact that MSCI USA alone makes 54% of MSCI all country world index and France, Germany and the UK make up in aggregate 60.72% (17.34%, 15.02% and 28.36%, respectively) of the MSCI Europe index that covers 85 percent of free float capitalization of European equity markets. These indices underline many financial derivative products, exchange traded funds etc. and therefore forecasting its volatility is crucial for pricing, risk management and asset allocation decisions.

[Insert Table 1]

The macroeconomic variables, retrieved at the monthly frequency, for comparable periods are taken from three different sources. The data for France, Germany and the UK are taken from the Eurostat database, whereas the US data come from the Federal Reserve Bank of St. Louis i.e. the FRED database. If the macro data were unavailable from either Eurostat or the FRED, we used the Organization for Economic Co-operation and Development (OECD) database. From these monthly values the growth rate is computed as the natural log difference, except for the term structure of interest rates where we

take the simple difference of the yields 10-year maturity bond and monthly T-bills or any monthly duration interest bearing security. To compute the macro-volatility of undertaken aggregate variables in our work, we use innovations from autoregressive-moving average (ARMA) models applied on each macro variable growth series (Schwert 1989).¹⁹

In total, we estimate GM models using eight macroeconomic variables: the term structure of interest rates (TermStr), the exchange rate (EXR), the narrow money (NM), the broad money (BM), the consumer price index (CPI), the industrial production (IP), crude oil prices (Oil) and the unemployment rate (UnEmp).²⁰ These macro variables are strongly interdependent, especially for developed countries and using them simultaneously, can cause issues pertaining to multicollinearity, over parameterisation and model convergence.²¹ To deal with these aspects, we employ dynamic principal component analysis (PCA) on these eight macro series, for each country, to assimilate information that explains the variance of these variables through meaningful and integrating components.²² Furthermore, Asgharian et al. (2013) and Virk and Javed (2017) have shown that integrated changes in the macro environment, as captured by the first two components coming from PCA analysis, have clear information benefits relative to GM model that uses a singular macro variable. Therefore, in addition to the eight macro variables, we take the first two components from the PCA analysis for each market when we estimate competing and benchmark models in this study. The summary statistics for the daily returns and squared returns on all four equity indices are presented in Appendix Table AI.

[Insert Figures 1 and 2]

[Insert Tables 2 and 3]

6 Empirical results and discussions

6.1 sections 5.2, 5.3. Tables 5 and 7 from IRFA. Plus Table 9. Put Tables 6, 8, 10 in the Appendix.

We start our estimation period from January 1999 so we can capture recent trends in market volatility.²³ To maximise on changing trends in the period prior to year 1999, we employ MIDAS lags of 5-year²⁴ duration starting from January 1994. Hence, the one step ahead volatility is forecasted with 5 years of lagged RV and/or macroeconomic data depending on model specification.

Table 1 shows the set of GM models that are competing against the benchmark GM-RV model. Furthermore, for all models, we estimate two GM models: the first in levels using the growth in the macro variables or the principal components, see equation (6) and the second is when the lagged volatilities of each of the macro variables or the principal components are taken together with the lagged levels, see equation (7). Finally, given the noted gains of principal components in improving the GM forecasting ability in Asgharian et al (2013) and Conrad and Loch (2015), we augment our benchmark model GM-RV with PC1 or PC2, see models 12 and 13 in Table 1 and equation (8) for reference.

¹⁹ We use innovations after fitting a best fit ARMA model as given by Bayesian information criterion.

²⁰ The exchange rate for France, Germany and the UK are taken against USD, whereas for the US we take it against a basket of currencies, as provided by FRED.

²¹ The issues pertaining to convergence for MIDAS regressions was frequented quite often due to non-convex objective function when we employed flexible weighting in search of optimal weight structure to exploit the information content in each aggregate variable.

²² PCA analysis has the benefits of over parametrization in model estimations when conditional variance can be influenced by a range of factors and effectively removes noise from the signal. We apply the dynamic PCA such that the stationary macro series are transformed to have standard normal distribution with zero mean and unit variance. We note that the first two components from dynamic PCA invariably explain 70-90 % of the variability in the total factor variance of the macro environment of economies investigated in our work.

²³ This choice is to alleviate any potential structural breaks that are potentially possible due to the changing patterns in the run to introduction of the Euro that may affect the estimation of conditional volatilities, both total and long term, when their evolution may depend on different macro variables. This is especially applicable to the European markets.

²⁴ Different lag structured are implemented in the estimation process, and based on optimal convergence of the model together with minimum use of data for smoothing, 5-year lag (K=5) is selected.

In sum, we compare the performance of 26 competing models against our benchmark GM-RV while using two different weighting schemes.²⁵ That is, all GM models are estimated with unconstrained weights (both the w_1 and w_2 are estimated) or constrained weights (when only w_2 is estimated and w_1 is kept fixed, i.e., $w_1 = 1$) for the beta smoothing function $\phi_k(w_1, w_2)$ for each lagged input variable in the MIDAS filter.²⁶

6.1 Forecasting errors of the competing models relative to benchmark model

In this section, we compare the relative performance of the quadratic loss functions of all competing models relative to total and secular component forecasts coming from the benchmark GM-RV model. Table 2 presents these ratios using the unconstrained/flexible weighting scheme in the MIDAS filter for all four markets. As noted in Engle et al. (2013) the measure of forecasting, or relative performance covers POS forecasts: the parameters for all GM models are estimated using the full sample, and the forecasts for next month are computed using month end price data.

The MSE ratios under the headings σ_{w_1, w_2}^2 and τ_{w_1, w_2} are, respectively, for the total variance forecasting errors and the secular components. As described earlier, we estimate two GM models for each MIDAS input variable. Therefore, columns under the heading of X_l describe the ratios of the models that use the level of first order change in the input variables, while the columns under the heading X_{l+v} describe the ratios for models that use both the levels and the volatility of the input variables in the MIDAS filter.

The vast majority of the ratios reported in the Table 2 across all four markets are less than one which indicates that several competing models outperform the conditional predictions given by the GM-RV model. This result highlights the benefit of incorporating macro information in forecasting total and long run variance component.

However, there are exceptions – a higher quadratic loss is noted in some instances. For France, for both total and long run variance predictions, the MSE ratio is above one when the competing GM model’s MIDAS input variables are level and variance of TermStr or EXR, uses CPI related information variables, level of UEmp, PC2 or when PC1 and PC2 are taken together. The same is observed for German data for competing GM models that include level of BM, level and volatility of PC1, either combination of input variables that use PC1 and PC2 in the MIDAS smoothing. The best chance for the benchmark model, as shown by the MSE ratios, is reserved for GM estimations carried out for the UK: in 22 out of 52 MSE ratios the benchmark predictions brought smaller forecasting errors. In the US data, the benchmark model GM-RV is also outperformed by most of the competing GM models. Exceptions apply for GM models with CPI (both total conditional and long run variance forecasts), UEmp or PC1 (in either, only the total variance forecasts).

The best model across the three European markets is the GM-RV+PC1 model using both level and volatility of the PC variable in the MIDAS smoothing. Using this model, the quadratic loss of total and long run variance forecasts of the benchmark can be approximately reduced by 10% for France, by 7% for Germany and 10% for the UK. In the US, the model that brings largest efficiency in reducing forecasting errors, and in the prediction of conditional total and long run variances, is the GM-PC1+PC2 model that brings an efficiency of up to 15% and 17%, respectively. This performance in the US market is followed by GM-RV+PC1 and GM-TermStr.

²⁵ We also compare the MSE forecasting errors of competing errors using another benchmark where we fix the RV in a month i.e. a fixed span GM-RV model. The results of these comparisons are available upon request and qualitatively resemble what we report using the benchmark GM-RV model in this study.

²⁶ The GM estimations show that, across all models and markets, the estimates in the GARCH model are usually significant at 5% or below critical t-values and estimate values are in line to the vastly available evidence for GARCH (1,1) models. The regression coefficients, in all models and across markets, on RV in the MIDAS regressions are positive and super significant regardless of the choice of weighting scheme. However, we note that MIDAS input variables in the competing GM models are more often significant with unrestricted weighting scheme at conventional 5% critical t-values. The significance of these estimates with the restricted weighting scheme reduces drastically. These results are not reported to conserve space for the large number of regressions, with even larger number of regression estimates, carried in our work and available upon request.

We replicate the results in Table 2 using restricted beta smoothing function: $w_1 = 1$, where only w_2 is estimated as a free parameter for each input variable in the MIDAS filter. Table 3 provides the MSE ratios with constrained beta smoothing and the evidence against the benchmark model forecasts is starker than that observed with the flexible weighting scheme in (Table 2). These results are almost entirely uniform across the four markets. There are two exceptions: one for France, where the GM-RV+PC2 (either using level or level and volatility together in the MIDAS filter) has a larger MSE than the benchmark model and the other for the UK, where the RP ratio of the GM-BM using the level and volatility of the input variable is above 1.

Within these general gains of using macro variables (on their own or through PCs), the improvements in predicting the volatility of the US equity market are relatively higher than improvements found in other markets. The ratios of the macro models in the US are relatively lower, and this applies to both the total and long run variance forecasts across all models. The largest improvement of almost 19% is observed, relative to the benchmark GM-RV model forecasting error, when the level of the term structure of interest rates is used in the GM model. Other notable GM models that bring large improvements (approximately 14-16%) include GM-Oil, GM-UEmp and GM-CPI against the common benchmarks for the total and secular variances coming from GM-RV model benchmarks.

The gains in accuracy are lower for the UK market than gains seen in the conditional variance predictions for the US market, but they are higher than the accuracy gains observed for France and Germany. For instance, in the UK the maximum gain in predicting total volatility is around 10% and it is achieved when we use the information related to changes in TermStr or crude oil prices. The quadratic loss of the monthly secular trend in the conditional variance can also be reduced by a maximum of 10% by using growth in crude oil prices or the unemployment rates in the GM models. In Germany and France, the gains in forecasting total and long run volatilities by the competing GM models with IP or the UEmp (with levels only) bring meaningful differences. For France, other notable models outperforming the benchmark forecasts are GM-CPI (level and volatility) for the secular component and GM-PC2 and GM-PC1+PC2 (in either case, level only) for the total variance forecasts. The MSE efficiency gains for using these variables in the MIDAS regression are in the range of 8-9%.

[Insert Tables 4 and 5]

We also note that including the volatility of the input variables together with the levels does not necessarily result in a more accurate forecast while employing the restricted weighting scheme. For instance, in the US the model that uses the term structure level is preferable to all other models. However, the accuracy of this model is significantly reduced from 19% to 13% when the volatility of the term structure is used as an additional predictor.

If the imposed variance evolution structure is unrelated to the relatively unconvincing performance of GM models that add volatility of input variables, this drop in accuracy might be attributed to over-parametrisation, small sample bias and parameter estimation errors.²⁷ Obviously, in our case, for most variables the levels seems to have more explanatory power than the volatilities.

6.2 Pairwise model comparisons: tests of equal forecasting ability

The quadratic losses reported in Tables 2 and 3 are only sample estimates and hence, we need to infer the accuracy of the models in the population. We proceed by conducting a pairwise comparison using the DM test, defined in equation (9), to compare the forecasts of the GM-RV benchmark for both total and long run variances in the set of competing models. These results are reported in Tables 4 and 5.

²⁷ Asymptotically, the larger model may not produce less accurate forecasts than the smaller model. But in finite samples this may occur due to small sample bias or parameter estimation errors.

The bold (negative) statistic values for the DM test in Table 4 imply that the competing model outperforms the benchmark model forecasts at conventional 5% or below significance levels. Thus, considering the bold and negative DM test statistic values in Table 4, there is overwhelming evidence that the competing models' POS forecasts, for both total and long variances, bring efficiency gains in reducing MSE relative to the benchmark model. These results are consistent across the four countries in the sample. Furthermore, we note that the competing GM models that perform well in the MSE ratio tests have the largest DM test statistic values i.e. their ability to suppress MSE relative to benchmark model is super significant. Table 5 shows the results with the restricted weighting scheme. The EPA (beneath the heading σ_{w2}^2) for the majority of the competing GM models' total variance forecasts is rejected – Germany is only exception where 13 models reject GM-RV model's total variance forecast.

This implies that the DM test statistic values are not statistically significant at conventional levels. There are scattered instances in which the DM test statistic values are significant e.g. for the long run variance in France GM models with changes in BM and CPI. In the UK, the GM models with the EXR changes and together with its volatility do the same. In the US, no competing GM models beats the benchmark in forecasting total volatility. Only for Germany do the total variance EPA tests show that a larger number of the competing GM models match to the benchmark model forecast: out of 26 there are 15 GM models in which the DM test statistic values are statistically significant. With respect to long run forecasts (under the columns τ_{w2}) the results show the opposite picture. All competing models, across all markets, have EPA as of the benchmark model secular forecasts: the DM test values reject the null at 5% or below significance values.

We note that the DM tests with the restricted weighting scheme casts doubt on the proclaimed usefulness of macro information in the GM models: when it comes to total variance, the evidence against the benchmark model forecast weakens drastically. That is, as our results show that even the model that uses RV and PCs with restricted weighing does not outperform the benchmark forecast for the total variance these results question the inference coming from the so-called informative, MSE, ratios. Furthermore, given the statistical rigor of the DM test compared to the MSE ratios, we note that different weighting schemes generate different test values and consequent inferences. We make two interlinked observations. First, the non-EPA of the competing GM models for the total variance forecasts with the restricted weighting scheme is consistent with the results in Engle et al. (2013) and Conrad and Loch (2015): they report that the restricted weighting scheme is optimal for the GM-RV model. We show superimposing restricted weighting on competing GM models is potentially inapt. Second, this stresses the importance of estimating the best weighting combinations for GM models that incorporate macro variables to let these specifications have a full chance computing their variance forecasts relative to the benchmark forecasts. This is particularly applicable in the evaluation of the total-variance forecasts.

Nonetheless, this does not appear to influence the evidence against the long run variance forecasts: macro information improves forecasting efficiency of MIDAS component using either weighting scheme. This aspect was previously noted with restricted weighing in Asgharian et al. (2013).

[Insert Table 6]

6.3 Multiple model comparisons: tests of superior predictive performance

The previous sections have collected evidence that shows how (relatively) weak the benchmark model is: many models outperform the GM-RV model. Invariably, the DM test shows that the inferences from MSE ratios are sensitive to which type of weighting scheme is adopted when evaluating total variance forecasts. Although our specification search for a good forecasting model is extensive and has samples from four countries, it is possible they are the results of coincidence and

data mining. To deal with this, we carry out the reality check (RC) (White, 2000) and superior predictive ability (SPA) tests (Hansen, 2005).

Following the null hypothesis that no model is better than the benchmark model, the RC and SPA tests are not only robust to data snooping issues but are also more powerful in their ability to jointly compare the performance of multiple models. A rejection of null implies that there is at least one model in the competing set of models which is better than the benchmark model.²⁸

Before examining the results of multiple testing using RC and Hansen SPA tests, we assess how many of the 26 competing models for each market outperform the benchmark model using the naïve p-values. Note that, the naïve p-values are the bi-model – a variant of DM test – bootstrapped RC test values that are computed as if the best model is the only model in the competing model space. These p-values are an important signal to check for the data snooping bias. White (2000) shows that the naïve p-values are the lower threshold for the RC multiple test p-values: therefore, he suggests it is only meaningful to test for data snooping bias when the p-values are small.

[Insert Tables 7 and 8]

As the RC test with naïve p-value below 0.05 implies rejecting null, we count all the instances for variance forecasts (both total and long run) and divide them by the total number of competing models (i.e., 26 in this case). We report these proportions in Table 6 for the POS loss functions. Results show that no model has efficiency gains in forecasting total volatility using either of the two weighting schemes. This result is observed for all the four markets.

However, there are quite a few competing GM models that outperform the benchmark model's long run variance forecasts. The results for the US results are one exception. Further, we note that using naïve p-values results, there are more competing GM models with SPA when unrestricted weighting is used relative to restricted weighting scheme. This is again endorses our expectation that using the restricted weighting scheme can undermine the utility of macro information in GM models. These results, together with the findings reported in sections 6.1 and 6.2, show that there are competing GM models that can be more informative or can reduce MSE more than the benchmark model. However, this inference is not applicable for total variance forecasts when examined by DM test with restricted weighting scheme or by naïve p-values using either weighting scheme. Thus, this observation questions the practice of drawing inferences from of simplistic tests combined with a restricted weighting scheme, at least for the total variance comparisons. But it also suggests increased caution in evaluating GM forecasts coming from competing models.

Therefore, robust SPA tests that account for multiple testing are critical when carrying out model comparisons to improve inferential reliability. Table 7 shows these results for the RC test and Hansen test. Following Hansen (2005), we compute three test statistics for the consistent (centre, c), upper (u) and lower (l) bounds for both RC test and Hansen test. A substantial difference between these upper and lower bounds indicate the presence of poor alternative specifications (for details, see Hansen 2005). Each test-statistic p-value is computed from 1000 bootstrap resamples and bootstrap parameter $q = 0.25$.²⁹ Nonetheless, for ease of discussion, we only refer to the consistent p-values in this section. The least naïve p-value from the pairwise RC tests is also reported to see the full effect of the joint testing to capture the data mining bias.³⁰

²⁸ Both are tests for superior predictive ability as ruled by the null hypotheses of RC/White test and SPA/Hansen test

²⁹ Our results remain unchanged if we increase the number of bootstrap samples to 10,000. Furthermore, we assess the sensitivity of the SPA results with different values for bootstrap parameter 'q' that controls for time dependence: a $q=1$ completely ignores the time dependence. The results using lower time dependence than $q=0.25$ do not alter our results in any manner. We also note that Hansen (2005) and Gonzalez et al. (2004) also use time dependence values of $q=0.25$ that accounts for substantial time dependence corresponding to the empirical observations in the time series modelling of the equity volatility.

³⁰ White (2000) reports that the difference between the naïve p-value and RC test can be described as the data snooping bias in the specification search of better models than the benchmark model.

The provided p-values for the long run variance forecast comparison show that the naïve p-values are far lower than the test statistics that employ multiple testing and are thus, biased against the null. The p-values for all six test statistics i.e. RC_L , RC_C , RC_U , $Hansen_L$, $Hansen_C$, and $Hansen_U$ are fairly large and do not reject the null for the total variance forecast comparisons. In other words, the benchmark model's total variance forecast cannot be out-performed given the RV based approximation of unobservable daily volatility. In most cases the data snooping bias is at least as large as the naïve p-value itself. This is a startling inference which, when linked with the results given by the DM test, especially with flexible weighting, shows that the results from even the DM test are untenable when the data snooping bias is considered.

However, the picture is slightly different in the case of long-term volatility forecasts. The joint testing substantiates the conjecture that at least one of the competing models, which incorporate economic information, generates a more accurate long-term volatility forecasts than benchmark. The null hypothesis is rejected by both RC and Hansen SPA tests when the flexible weighting scheme is employed. For the restricted weighting scheme, the same inference follows and again the only exception is the US results.

For the US, the benchmark model's prediction for the long run volatility show superior predictive ability: the naïve p-value for the best performing model is 0.143 which is a lower bound for the more stringent RC and Hansen test values for the null. Here we also note the advantage of the Hansen test, as highlighted in Hansen (2005): its p-values are fairly lower than the corresponding RC test statistic values. For example, in the US secular variance forecast comparisons when restricted weighting is used RC_c is 0.72 while $Hansen_c$ p-value is 0.467.³¹ For the US data, our results are in accordance with findings reported in Conrad and Kleen (2020) that the GM RV model is outperformed by models involving macroeconomic information.

Overall, we infer that GM models incorporating macroeconomic information may bring benefits in forecasting the long run variances but the same cannot be concluded with respect to their total variance forecasts.

[Insert Table 9]

In order to check the sensitivity of the results reported in Table 7, we replace RV with two other variance proxies i.e. conditional variance forecasts from GARCH (1, 1) model and variance of implied volatility indices of the markets undertaken in this study.³² We recalculate the MSEs as:

$$\begin{aligned} h_{MSE}^{\sigma^2} &\equiv \frac{1}{T} (h_{i,t} - \sigma_{i,t,M}^2)^2 & \text{and} & & h_{MSE}^{\tau} &\equiv \frac{1}{T} (h_t^2 - \tau_{t,M})^2, \\ IV_{MSE}^{\sigma^2} &\equiv \frac{1}{T} (IV_{i,t}^2 - \sigma_{i,t,M}^2)^2 & \text{and} & & IV_{MSE}^{\tau} &\equiv \frac{1}{T} (IV_t^2 - \tau_{t,M})^2. \end{aligned}$$

where the $h_{MSE}^{\sigma^2}$, h_{MSE}^{τ} , and $IV_{MSE}^{\sigma^2}$, IV_{MSE}^{τ} are, respectively, the loss functions computed for total and long run variances using GARCH (1,1) and the square of implied volatilities of each of stock markets.³³

³¹ Hansen (2005) show that RC tests can be manipulated in the presence of poor and irrelevant alternatives in the sample of models. This results in less power and non-rejection of the null hypothesis. Hansen (2005) alleviated this problem in his version of the SPA test by invoking a sample-dependent distribution under the null hypothesis.

³² The benchmark model gives total and long run variance forecasts by smoothing RV in the MIDAS regression, it is potentially true that choice of variance proxy is driving the results that we report in table 9 using robust tests that account for multiple testing.

³³ We retrieve data for the implied volatility proxies for respective country equity index from DataStream and Chicago Board of Exchange website. Both, daily and monthly implied volatility indices on CAC40, DAX30, FTSE100 and S&P500 for France, Germany, the UK and the US, respectively are downloaded. The implied volatilities indices on DAX30 and S&P500 are available for the matching period, whereas for CAC40 and FTSE100 these indices begin from January 2000. The DataStream codes for the implied volatility of the CAC30, the DAX30, and the FTSE100 are the CACVOLI, the VDAXNEW, and the VFTSEIX respectively.

The results of multiple testing with these loss functions, for the total and long run variances and for the flexible and restricted weighting schemes, are presented in Table 8. In summary, the crux of evidence reported in Table 7 is replicated in Table 8. This implies that our results are insensitive to choice of variance proxy in the POS model forecast comparisons.

6.4 Out-of-sample model forecast comparisons

Until now we have carried out all tests on the pseudo out-of-sample forecasts (Engle et al. 2013, Conrad and Loch 2015). However, to test the full effect of our results reported in Tables 7 and 8, we compare the out-of-sample (OS) forecasts of all the competing models using a rolling window forecasting scheme.³⁴ To conduct the out of sample forecasts comparison, we get OS conditional total variance and long run variance forecasts for 2012-16 for one year at a time.

[Insert Tables 10 and 11]

Procedurally, the parameters are obtained using a rolling 13-year estimation window for each GM model, which are held constant during the subsequent year to compute daily (252 step-ahead) total variance and monthly (12 step-ahead) long run variance forecasts. The estimation sample is then moved one year ahead to re-estimate the parameters and the forecasts are made for the next subsequent year and so on. Thus, we perform this procedure five times to obtain total and secular variance forecasts for years 2012, 2013, 2014, 2015 and 2016. The estimation windows correspondingly cover rolling sample periods of 1999-2011, 2000-2012, 2001-2013, 2002-2014 and 2003-2015.

Here again, before getting into the implications of robust multiple SPA tests, using the naïve p-values (below 0.05) we count the number of competing models' forecasts, both for the total and long-term variances, that have SPA compared to the benchmark. The results are presented in Table 9 and these show that several competing models that might have SPA than the benchmark model. Can this be replicated using more powerful/robust tests?

We present the p-values for the RC and Hansen tests in Table 10 for the out of sample forecasts. Table 10 shows that using the flexible weighting scheme the RC tests have large p-values and null cannot be rejected at 0.05 significance p-values for the total variance forecast errors for any market. On the contrary when it comes to the $Hansen_c$ statistic: the SPA of the benchmark total variance forecast is rejected for all except for Germany. However, even for Germany the null can be rejected at 10% confidence test values: p-value for the $Hansen_c$ test statistic is 0.088. With the restricted weighting scheme, both tests reject the notion that competing models' forecast have SPA for the total variance relative to the benchmark model forecast.

The OS forecasting comparisons for the long run variance forecasts, for either type of weighting scheme, is more similar to what we have seen in Table 7 for the POS comparisons. That is, using the flexible weighting scheme, the RC_c and $Hansen_c$ statistic p-values reject null for all the markets. However, with fixed weights the SPA of the benchmark model's long run forecasts are not rejected for France or the US. The non-rejection of the null for France is borderline and can only be rejected at 10% confidence levels using Hansen central test statistic p-values. The non-rejection of the US long run variance forecast in POS and OS comparisons while adopting restricted weights signifies what we have asserted earlier: a simplified weighting scheme can contaminate the actual forecast evaluation inferences.

We redo the OS multiple tests using alternate variance proxies coming from GARCH (1, 1) and implied volatility indices. The resultant p-values are presented in Table 11. Overall, there are no differences to what we reported using RV as a proxy for the latent variances. The only exception is the non-rejection of the null for France total variance forecast comparison, whereas with RV (Table 10 using Hansen central test statistic) we were able to reject the null hypothesis.

³⁴ Multiple step ahead forecasts for long periods such as five years are not feasible in the context of the models examined in our work, knowing latent variance is time varying and the expected future value of the macroeconomic growth and volatility is not measurable given the time of forecast information filter.

To summarize, we divide our results in two parts for ease of description. In the first part, we replicate evidence from Engle et al. (2013) and Asgharian et al. (2013) using pairwise MSE ratios and DM test and provide new results using French German and British financial and economic datasets. These tests are carried out using the one step ahead GM variance forecasts. In the second part, we complement new evidence in the GM variance forecasting literature by examining data snooping issues in making the forecast comparisons among GM models.

Our results from the first part show that the MSE of the GM models that include macro information, alone or with the RV in the MIDAS regressions and with either type of weighting scheme, are smaller than the benchmark model. These findings are consistent with earlier studies. We find that total and long run variance forecasts from the GM model that use RV and PC1 give the best performance across all markets when we use a flexible weighing scheme. While results using restricted weighting scheme find even better results when using macro information in the GM models. These findings are also confirmed by the EPA DM test – only problem found is the use of restricted weight in forecasting total variance. That is, the DM tests provide overwhelming evidence that the competing models POS forecasts, for both variance components, bring efficiency gains while using the flexible weighting scheme. These results are uniform across the four countries. However, when we adopt restricted weighted scheme the consistency in results is only witnessed for the long run variance forecasts.

Drawing together the results of MSE ratios and the DM test, we note that there are differences in forecasting gains while using a particular weighting scheme in the MIDAS smoothing using either type of tests. However, the use of informative ratios in making reliable inferences is untenable. Furthermore, the DM test illustrates that the benchmark total variance forecasts underperform using unrestricted weighting scheme compared to the same tests with a restricted weighting scheme. This observation implies that projecting a weighting structure that is optimal for the GM-RV model may result in limiting the scope for the macro variables in the GM models to predict the secular variance component. Thus, we suggest that econometricians and practitioners should be wary of this simplification while searching for precise and reliable variance forecasts. Essentially, superimposing a particular evolutionary structure in MIDAS may yield results that favour null even when that is not the case. Therefore, incorporating relevant low frequency information, and the consequent evolution and approximation of variances, should not be comprised by resorting to a general weighting pattern that might not be suitable due to its data generating process.

The low reliability of these results when assessed through powerful tests that account for data mining endorse our scepticism: joint testing shows that no model outperforms the benchmark model's total variance forecast in all markets. However, evaluation of the long run variance coming from competing GM macro models there are clear gains over the benchmark model's secular trend. With respect to multiple testing we also note the value of pairwise tests: if the p-value for DM test or naïve robust test is large, there is no need to test for data snooping biases because, as White (2000) reports, p-values from multiple tests can only be larger than the naïve p-value. However, caution should be the primary concern when the p-values from the tests with EPA or naïve p-values are small. Overall, the generality of our main results with respect to scepticism to data snooping biases and adoption of a particular weighting scheme persists for using different proxies for latent variance and using a rolling OS forecasting scheme.

We summarize our results by making the following four observations. First, our results show that inferences based on informative or pairwise model comparison tests can be misleading, especially when they are about the gain from using

macro information in projecting total GM variance. This is especially dangerous when we have seen the evolving evidence studying GM models have generalized the evidence for total variance to long run variance forecasts. Second, we show that GM models are an important addition in forecasting volatility at different time frequencies but we have to be cognizant of what type of variance forecast we are looking for. Third, we should not compromise the weighting scheme for the MIDAS input variables by making ad hoc choices and, fourth, we should be wary for data mining biases that may plague a forecast that otherwise appear statistically validated.

These results are important given the demand for reliable, accurate variance forecasts for devising risk planning and management protocols at institutional and individual levels, especially since investors vary in their investment choices across time horizons; e.g. active and passive investment decisions require variance forecasts for two different investing styles. Therefore, using macro variables adds information that improve modelling and forecasting of long run or total variance, we should capitalize on it. This applies when there are several low frequency processes which contain the relevant, independent information needed to forecast latent variance for different time horizons. For example, the term structure of interest rates is widely used by analysts to predict economic expansions and contractions. Flattening of the yield curve predicts recessions while its steepening occurs prior to expansions. As equity markets become more nervous prior to an impending recession, volatility increases. These examples illustrate how the term structure is intuitively linked to the future volatility of financial markets. Furthermore, markets are always looking for employment figures in order to foresee the future performance of the economy. To this effect, our results pick up these independent information channels while taking appropriate caution to ascertain at what frequency these gains are established, and where the benefits are erroneously simplified. In the GM modelling, we recommend use of high frequency data to develop precise proxies for latent variance to compare forecasts and use of different dimensions of sentiment changes in varying market conditions are important areas for future research.

References

- Andersen, T. G., Bollerslev, T., Diebold, F. X., & Labys, P. (2001a). The distribution of realized exchange rate volatility. *Journal of the American statistical association*, 96(453), 42-55.
- Andersen TG, Bollerslev T, Diebold FX, Ebens H (2001b) The distribution of realized stock return volatility. *Journal of Financial Economics* 61, 43—76.
- Andersen, T. G. and Bollerslev, T., (1998) Answering the Skeptics: Yes, Standard Volatility Models do Provide Accurate Forecasts, *International Economic Review*, 39 (4), 885-905.
- Andersen, T.G. and Bollerslev, T (1997) Intraday Periodicity and Volatility Persistence in Financial Markets, *Journal of Empirical Finance* 4, 115-158.
- Asgharian, H., Hou, A., Javed, F. (2013) The importance of the economic variables in predicting the long term volatility; a GARCH MIDAS approach, *Journal of Forecasting*, 32, 600-612.
- Asgharian, H., Christiansen, C. Hou, A. (2016) Macro-Finance Determinants of the Long-Run Stock–Bond Correlation: The DCC-MIDAS Specification, *Journal of Financial Econometrics*, 14 (3), 617–642.
- Barndorff-Nielsen, O. E., & Shephard, N. (2001) Non-Gaussian Ornstein–Uhlenbeck-based models and some of their uses in financial economics. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2), 167-241.
- Barndorff-Nielsen, O. E., & Shephard, N. (2002) Econometric analysis of realized volatility and its use in estimating stochastic volatility models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(2), 253-280.
- Bollerslev, T., Engle, R. F. and Nelson, D. B. (1994). “ARCH Models”, in Robert F Engle and Daniel F McFadden (eds.), *Handbook of Econometrics*, 4, 2951--3038.
- Conrad, C., Loch, K., and Rittler, D. (2014) On the macroeconomic determinants of long-term volatilities and correlations in U.S. stock and crude oil markets, *Journal of Empirical Finance*, 29, 26-40.
- Conrad, C., Loch, K (2015) Anticipating Long-term Stock Market Volatility, *Journal of Applied Econometrics*. 30, 1090—1114.
- Conrad, C., Kleen, O. (2020). \Two are better than one: Volatility forecasting using multiplicative component GARCH-MIDAS models." *Journal of Applied Econometrics*, 35, 19-45.

- Christie, Andrew A. (1982) The stochastic behavior of common stock variances: Value, leverage and interest rate effects, *Journal of Financial Economics* 10, 407-432.
- Diebold, F. X., and Mariano, R. S. (1995) Comparing Predictive Accuracy, *Journal of Business & Economic Statistics*, 13, 253–263.
- Diebold, F.X. and J.A. Lopez (1996) Forecast Evaluation and Combination in G.S. Maddala and C.R. Rao (eds.), *Handbook of Statistics*. Amsterdam: North-Holland, 241-268.
- Engle, R. F., Ghysels, E., & Sohn, B. (2013) Stock market volatility and macroeconomic fundamentals. *Review of Economics and Statistics*, 95(3), 776-797.
- Fang, L., Yu, H., and Huang, Y. (2018) The role of investor sentiment in the long-term correlation between U.S. stock and bond markets, *International Review of Economics & Finance*, 58, 27-139,
- Ghysels E., Santa-Clara P., Valkanov R. (2006) Predicting Volatility: Getting the Most out of Return Data Sampled at Different Frequencies, *Journal of Econometrics*, 131. 59-95.
- Ghysels, E., Santa-Clara, P., and Valkanov, R. (2004) The MIDAS touch: Mixed Data Sampling Regression. Discussion Paper UNC and UCLA.
- Ghysels, E., Sinko, A., and Valkanov R. (2007) MIDAS Regressions: Further Results and New Directions, *Econometric Reviews*, 26, 53-90.
- González-Rivera, G., Lee, T. H., & Mishra, S. (2004) Forecasting volatility: A reality check based on option pricing, utility function, value-at-risk, and predictive likelihood. *International Journal of forecasting*, 20(4), 629-645.
- Hansen, P. R. (2005) A test for superior predictive ability. *Journal of Business & Economic Statistics*, 23(4), 365-380.
- Hansen, P. R., & Lunde, A. (2005). A forecast comparison of volatility models: does anything beat a GARCH (1, 1)? *Journal of Applied Econometrics*, 20(7), 873-889.
- Lindblad, A. (2017) Sentiment indicators and macroeconomic data as drivers for low-frequency stock market volatility." MPRA Paper 80266, University Library of Munich, Germany.
- Lopez, J. A, (2001) Evaluating the Predictive Accuracy of Volatility Models *Journal of Forecasting*, 20(2), 87-109
- Meddahi, N. (2002) A theoretical comparison between integrated and realized volatility. *Journal of Applied Econometrics*, 17(5), 479-508.
- Officer, R. R. (1973) The variability of the market factor of the New York Stock Exchange. *The Journal of Business*, 46 (3), 434-453.
- Pan, Z., Wang, Y., Wu, C., and Libo, Y. (2017) Oil price volatility and macroeconomic fundamentals: A regime switching GARCH-MIDAS model, *Journal of Empirical Finance*, 43, 130-142.
- Patton, A. J. (2011) Volatility forecast comparison using imperfect volatility proxies." *Journal of Econometrics*, 160, 246-256.
- Politis, D. N., & Romano, J. P. (1994) The stationary bootstrap. *Journal of the American Statistical Association*, 89, 1303 – 1313.
- Schwert, G. W. (1989) Why does stock market volatility change over time? *Journal of Finance*, 44 (5), 1115-1153.
- Shiller, Robert J. (1981a) The use of volatility measures in assessing market efficiency, *Journal of Finance*, 36, 291-304.
- Shiller, Robert J. (1981b) Do stock prices move too much to be justified by subsequent changes in dividends, *American Economic Review* 75, 421-436.
- Virk, N. and Javed, F. (2017) European equity market integration and joint relationship of conditional volatility and correlations, *Journal of International Money and Finance*, 71, 53-77
- West, K. (1996) Asymptotic inference about predictive ability. *Econometrica*, 64, 1067 – 1084.
- White, H. (2000) A reality check for data snooping. *Econometrica*, 68(5), 1097-1126.

Table 1: The Set of Competing GARCH-MIDAS Macro Models

The first line of the table shows the table benchmark model used in this study i.e. GARCH-MIDAS (GM) model that includes the monthly rolling window (RW) realized variance (RV) in the MIDAS filter. Models 1-11 only include the level of the listed macro variable in the GM model. The last two models, i.e. 12 and 13, add the level of the first and second principal components of all the macro variables shown in the table to the benchmark GM-RV model. We estimate 13 additional models that include both the level and the volatility of the explanatory variables in each of 13 competing models shown below. Hence, a set of 26 models compete against the benchmark GM-RV.

Model	GARCH-MIDAS models	Acronym
Benchmark model	realized variance (RV)	GM-RV
1	The Term Structure of Interest rates	GM-TermStr.
2	Exchange Rate	GM-EXR
3	Narrow Money	GM-NM
4	Broad Money	GM-BM
5	Consumer Price Index	GM-CPI
6	Industrial Production	GM-IP
7	Crude Oil Prices	GM-Oil
8	Unemployment	GM-UEmp.
9	First Principal Component, PC1	GM-PC1
10	Second Principal Component, PC2	GM-PC2
11	With PC1 and PC2	GM-PC1+PC2
12	RV with PC1	GM-RV+PC1
13	RV with PC2	GM-RV+PC2

Table 2 MSE Ratios using flexible weighting scheme in the GM models

In all the estimated models, the MIDAS explanatory variables span the past five-year lagged data (from January 1994 to December 1998). The models use monthly input variables except the benchmark model where we use monthly realized variance – 22 day rolling window – that is available daily. Thus, the two-component volatility models are fitted over the period of January 1999 to December 2016 to carry out pseudo out-of-sample forecast comparisons using MSE ratios of competing model relative to GM-RV model i.e. the benchmark model. Results are presented in blocks for France, Germany, the UK and the USA. Here σ_w^2 and τ_w refer to the comparisons of total daily volatility and monthly secular volatility from competing k-models with the GM-RV model. The full sample parameters for each model are subsequently taken to compute the next period pseudo forecasts. The MSE of the forecasts given by each model is thus calculated, including the benchmark model, with respect to monthly RV. For each equity market, MSE ratios below $\sigma_{w1,w2}^2$ and $\tau_{w1,w2}$ are computed when the parameters of the beta polynomial function are unconstrained (i.e. w_l and w_v for each input variable are allowed to be free parameters). The competing model is determined by the interaction of the models below the heading “competing models” with column headings X_l and X_{l+v} . This implies that the MSE ratios are separated across k -models that use only the level of the MIDAS input variable i.e. X_l and the ones that use the level and volatility of the input variables combined in the MIDAS filter i.e. X_{l+v} . All pseudo-forecasts are generated using fixed full sample parameter estimates as in Engle et al. (2013).

Competing models	France				Germany			
	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	0.988	1.004	0.992	1.008	0.971	0.987	0.970	0.986
GM-EXR	0.995	1.004	0.999	1.008	0.994	0.981	0.992	0.980
GM-NM	0.984	0.954	0.988	0.958	0.998	0.968	0.997	0.967
GM-BM	0.986	0.989	0.990	0.993	1.003	0.994	1.001	0.992
GM-CPI	1.009	1.008	1.013	1.012	0.967	0.967	0.965	0.965
GM-IP	0.981	0.959	0.985	0.963	0.994	0.935	0.993	0.934
GM-Oil	0.979	0.983	0.983	0.987	0.982	0.995	0.981	0.993
GM-UEmp.	1.005	0.958	1.009	0.962	0.995	0.968	0.993	0.967
GM-PC1	0.995	0.977	0.999	0.981	0.995	1.059	0.994	1.058
GM-PC2	1.028	0.998	1.032	1.002	1.009	1.003	1.008	1.001
GM-PC1+PC2	1.025	0.930	1.030	0.934	1.011	0.945	1.010	0.944
GM-RV+PC1	0.916	0.901	0.920	0.905	0.972	0.934	0.970	0.933
GM-RV+PC2	0.933	0.902	0.937	0.906	0.953	0.939	0.951	0.938
Competing models	UK				USA			
	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	0.977	1.001	0.992	1.017	0.899	0.904	0.882	0.887
GM-EXR	0.998	1.031	1.014	1.047	0.964	0.966	0.945	0.947
GM-NM	1.013	0.986	1.028	1.001	0.947	0.969	0.929	0.950
GM-BM	0.999	0.999	1.014	1.015	0.992	0.988	0.973	0.970
GM-CPI	1.014	0.991	1.030	1.006	1.028	1.023	1.008	1.004
GM-IP	1.014	0.959	1.030	0.974	0.992	0.931	0.973	0.914
GM-Oil	1.011	0.979	1.027	0.994	0.975	0.996	0.956	0.977
GM-UEmp.	0.979	0.966	0.994	0.981	1.026	0.928	1.006	0.910
GM-PC1	1.007	0.992	1.022	1.007	1.006	1.014	0.987	0.994
GM-PC2	0.980	0.986	0.995	1.001	0.933	0.896	0.916	0.879
GM-PC1+PC2	1.019	0.962	1.035	0.977	0.944	0.846	0.926	0.830
GM-RV+PC1	0.932	0.901	0.946	0.915	0.939	0.898	0.921	0.881
GM-RV+PC2	0.922	0.910	0.936	0.924	0.912	0.926	0.894	0.909

Table 3 MSE Ratios using constrained weighting scheme in the GM models

In all the estimated models, the MIDAS explanatory variables span the past five-year lagged data (from January 1994 to December 1998). The models use monthly input variables except the benchmark model where we use monthly realized variance – 22 day rolling window – that is available daily. Thus, the two-component volatility models are fitted over the period of January 1999 to December 2016 to carry out pseudo out-of-sample forecast comparisons using MSE ratios of competing model relative to GM-RV model i.e. the benchmark model. Results are presented in blocks for France, Germany, the UK and the USA,. Here σ_w^2 and τ_w refer to the comparisons of total daily volatility and monthly secular volatility from competing k-models with the GM-RV model. The full sample parameters for each model are subsequently taken to compute the next period pseudo forecasts. The MSE of the forecasts given by each model is thus calculated, including the benchmark model, with respect to monthly RV – rolling window or monthly sum- as applicable to the sample of all models. For each equity market, MSE ratios below σ_{w2}^2 and τ_{w2} are computed when the parameters of the beta polynomial function are constrained: weighting scheme fixes $w_1 = 1$ and w_2 for each input variable in the MIDAS filter is estimated as a free parameter. The competing model is determined by the interaction of the models below the heading “competing models” with column headings X_l and X_{l+v} . This implies that the MSE ratios are separated across k -models that use only the level of the MIDAS input variable i.e. X_l and the ones that use the level and volatility of the input variables combined in the MIDAS filter i.e. X_{l+v} . All pseudo-forecasts are generated using fixed full sample parameter estimates as in Engle et al. (2013).

France					Germany			
Competing models	σ_{w2}^2		τ_{w2}		σ_{w2}^2		τ_{w2}	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	0.951	0.950	0.951	0.950	0.962	0.953	0.966	0.958
GM-EXR	0.943	0.975	0.943	0.975	0.963	0.981	0.967	0.985
GM-NM	0.960	0.958	0.960	0.958	0.968	0.967	0.973	0.971
GM-BM	0.952	0.955	0.952	0.955	0.963	0.963	0.967	0.967
GM-CPI	0.966	0.920	0.966	0.920	0.988	0.958	0.992	0.963
GM-IP	0.937	0.971	0.937	0.971	0.946	0.944	0.950	0.948
GM-Oil	0.944	0.972	0.944	0.972	0.955	0.971	0.960	0.976
GM-UEmp.	0.934	0.949	0.934	0.949	0.952	0.964	0.956	0.968
GM-PC1	0.974	0.963	0.974	0.963	0.987	0.980	0.992	0.985
GM-PC2	0.914	0.983	0.914	0.983	0.950	0.960	0.954	0.965
GM-PC1+PC2	0.931	0.974	0.931	0.974	0.963	0.985	0.967	0.989
GM-RV+PC1	0.946	0.942	0.946	0.942	0.956	0.955	0.961	0.959
GM-RV+PC2	1.018	1.020	1.018	1.020	0.952	0.957	0.957	0.962
UK					USA			
Competing models	σ_{w2}^2		τ_{w2}		σ_{w2}^2		τ_{w2}	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	0.907	0.914	0.908	0.916	0.808	0.873	0.813	0.878
GM-EXR	0.931	0.945	0.933	0.947	0.886	0.906	0.891	0.911
GM-NM	0.932	0.932	0.934	0.934	0.900	0.901	0.905	0.906
GM-BM	0.948	1.484	0.950	1.487	0.889	0.890	0.894	0.895
GM-CPI	0.896	0.895	0.898	0.897	0.854	0.851	0.858	0.855
GM-IP	0.916	0.919	0.918	0.921	0.882	0.864	0.887	0.869
GM-Oil	0.889	0.918	0.891	0.920	0.862	0.914	0.867	0.919
GM-UEmp.	0.910	0.906	0.912	0.908	0.847	0.865	0.852	0.869
GM-PC1	0.942	0.937	0.944	0.939	0.921	0.920	0.926	0.925
GM-PC2	0.895	0.922	0.897	0.923	0.860	0.863	0.865	0.868
GM-PC1+PC2	0.905	0.929	0.906	0.930	0.872	0.878	0.877	0.883
GM-RV+PC1	0.913	0.933	0.915	0.935	0.889	0.828	0.894	0.833
GM-RV+PC2	0.895	0.910	0.897	0.911	0.883	0.882	0.888	0.887

Table 4 The Diebold Mariano test using unconstrained weighting scheme

The table presents the t statistics of the Diebold Mariano (1995) test of equal predictive accuracy between pseudo out of sample forecasts from the GM-RV model and the competing k-models. The competing model is determined by the interaction of the models below the heading “competing models” with column headings X_l and X_{l+v} . Here the GM model below column X_l refers to the specification using the levels of the MIDAS input variables only, while X_{l+v} refers to the GM model that uses the level of the input variable together with its volatility in the MIDAS filter. The $\sigma_{w1,w2}^2$ refers to testing the predictive ability of the total daily volatility forecasts of the macro model while the $\tau_{w1,w2}$ refers to testing the accuracy of monthly volatility forecasts. This implies that MIDAS beta polynomial across all models are estimated unconstrained. The test-statistic value significant at 0.05 or below p-values are presented in bold.

Competing models	France				Germany			
	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	-2.209	-2.650	-2.557	-2.579	-4.197	-4.265	-12.815	-12.477
GM-EXR	-2.164	-3.138	-2.008	-2.873	-5.022	-6.146	-10.333	-9.833
GM-NM	-3.156	-4.460	-2.611	-3.740	-4.683	-5.327	-8.831	-6.899
GM-BM	-3.131	-3.321	-2.260	-2.237	-6.293	-4.568	-6.693	-9.461
GM-CPI	-2.378	-2.408	-2.128	-2.061	-3.904	-3.967	-4.148	-4.285
GM-IP	-2.247	-2.401	-2.394	-2.593	-4.404	-2.443	-9.221	-5.089
GM-Oil	-2.134	-3.623	-2.151	-3.287	-3.137	-4.390	-4.957	-9.711
GM-UEmp.	-2.955	-3.839	-2.655	-2.507	-5.849	-3.615	-9.658	-6.359
GM-PC1	-2.286	-3.419	-2.209	-2.796	-5.743	-2.099	-7.451	-8.704
GM-PC2	-2.582	-2.348	-2.430	-2.927	-4.745	-1.881	-11.584	-3.228
GM-PC1+PC2	-2.698	-3.837	-2.390	-4.233	-5.308	-5.790	-11.266	-7.976
GM-RV+PC1	-3.492	-3.562	-4.567	-2.426	-6.165	-3.481	-2.685	-4.778
GM-RV+PC2	-3.543	-3.369	-4.038	-5.113	-4.498	-3.918	-2.443	-1.410
Competing models	UK				USA			
	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	-2.201	-2.667	-2.597	-2.651	-2.371	-2.512	-19.330	-17.760
GM-EXR	-2.054	-3.484	-2.607	-3.220	-2.620	-2.631	-13.931	-8.475
GM-NM	-2.634	-3.472	-2.498	-3.646	-2.381	-3.830	-11.464	-4.743
GM-BM	-3.357	-3.142	-2.463	-2.467	-3.532	-3.530	-9.481	-10.069
GM-CPI	-2.740	-2.436	-2.435	-2.417	-1.914	-2.081	-11.884	-8.539
GM-IP	-2.455	-1.951	-2.586	-3.827	-3.100	-1.940	-5.995	-3.034
GM-Oil	-2.596	-3.513	-2.313	-2.999	-2.370	-4.941	-5.076	-4.107
GM-UEmp.	-2.716	-3.004	-2.931	-2.955	-2.605	-3.140	-8.616	-6.892
GM-PC1	-2.364	-3.068	-2.611	-3.097	-5.216	-5.872	-10.694	-10.088
GM-PC2	-1.515	-3.335	-2.096	-3.083	-2.655	-2.411	-12.358	-8.768
GM-PC1+PC2	-3.036	-1.639	-2.672	-3.252	-3.052	-2.548	-10.939	-5.825
GM-RV+PC1	-3.032	-2.616	-6.374	-6.154	-2.783	-2.325	-16.167	-31.927
GM-RV+PC2	-2.684	-2.778	-3.425	-2.603	-2.335	-2.647	-3.014	-1.962

Table 5 The Diebold Mariano test using constrained weighting scheme

The table presents the t statistics of the Diebold Mariano (1995) test of equal predictive accuracy between the GM-RV model and the competing k-models. The competing model is determined by the interaction of the models below the heading “competing models” with column headings X_l and X_{l+v} . Here the GM model below column X_l refers to the specification using the levels of the MIDAS input variables only, while X_{l+v} refers to the GM model that uses the level of the input variable together with its volatility in the MIDAS filter. The σ_{w2}^2 refers to testing the predictive ability of the total daily volatility forecasts of the macro model while the τ_{w2} refers to testing the accuracy of monthly volatility forecasts. This implies that MIDAS beta polynomial across all models is where weighting scheme fixes $w_l = 1$ and w_2 is estimated as a free parameter. The test-statistic value significant at 0.05 or below p-values are presented in bold.

France					Germany			
Competing models	σ_{w2}^2		τ_{w2}		σ_{w2}^2		τ_{w2}	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	-1.968	-1.893	-7.039	-7.074	-1.771	-1.873	-8.635	-9.581
GM-EXR	-1.459	-2.979	-7.774	-6.790	-2.701	-2.652	-9.019	-8.983
GM-NM	-1.672	-1.643	-6.186	-6.032	-1.589	-1.609	-8.758	-8.778
GM-BM	-2.428	-2.574	-7.119	-7.086	-2.114	-2.128	-9.313	-9.375
GM-CPI	-2.018	-1.576	-6.575	-7.406	-2.247	-1.568	-8.427	-9.459
GM-IP	-1.224	-1.719	-8.978	-6.737	-1.433	-1.219	-7.923	-9.579
GM-Oil	-1.498	-2.552	-7.771	-7.221	-1.333	-3.723	-9.933	-9.444
GM-UEmp.	-1.840	-1.904	-8.515	-9.021	-2.413	-2.667	-7.804	-7.509
GM-PC1	-2.436	-1.689	-8.369	-7.552	-2.834	-2.100	-10.873	-4.412
GM-PC2	-1.303	-2.333	-7.161	-5.158	-1.209	-2.705	-9.957	-7.548
GM-PC1+PC2	-1.328	-3.537	-8.745	-8.245	-1.828	-2.670	-11.644	-6.274
GM-RV+PC1	-1.898	-1.992	-6.924	-6.703	-1.733	-1.929	-11.440	-13.601
GM-RV+PC2	-3.418	-3.414	-14.874	-1.211	-1.727	-1.809	-7.018	-1.196
UK					USA			
Competing models	σ_{w2}^2		τ_{w2}		σ_{w2}^2		τ_{w2}	
	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}	X_l	X_{l+v}
GM-TermStr.	-1.633	-1.653	-4.506	-5.080	-1.804	-1.649	-5.430	-3.810
GM-EXR	-2.091	-2.461	-4.710	-4.716	-1.782	-1.612	-3.604	-3.613
GM-NM	-1.283	-1.272	-4.551	-4.463	-1.260	-1.448	-3.466	-3.404
GM-BM	-2.542	-2.457	-4.784	-4.811	-1.789	-1.799	-3.624	-3.637
GM-CPI	-1.383	-1.379	-4.691	-4.666	-1.578	-1.566	-3.354	-3.355
GM-IP	-1.100	-1.215	-5.210	-4.663	-1.292	-1.303	-3.810	-3.708
GM-Oil	-1.274	-1.793	-4.549	-4.197	-1.521	-1.674	-3.323	-4.102
GM-UEmp.	-1.661	-1.669	-4.446	-4.342	-1.837	-1.760	-4.784	-5.415
GM-PC1	-1.705	-1.187	-4.495	-4.216	-1.572	-1.313	-3.864	-3.983
GM-PC2	-1.126	-1.219	-4.615	-3.750	-1.898	-1.568	-4.672	-2.967
GM-PC1+PC2	-1.117	-1.740	-4.506	-3.609	-1.330	-1.366	-4.040	-5.299
GM-RV+PC1	-1.567	-1.707	-3.382	-6.947	-1.638	-1.788	-3.392	-7.692
GM-RV+PC2	-1.604	-1.694	-2.526	-6.318	-1.655	-1.720	-3.700	-7.335

Table 6 Robust pairwise testing for pseudo out-of-sample MSE losses

In this table, we report the proportion of variance forecasts – total and long run – that outperform the benchmark model’s corresponding forecasts using the naïve p-values. The naïve p-values are the bi-model – a variant of DM test – bootstrapped RC test values that are computed as if the best model is the only model in the competing model space. Thus, we count all the instances for variance forecasts and divide them by the total number of competing models (i.e., 26 in this case). To exemplify, a proportion of 0 implies that there no model that outperforms the forecast given by the benchmark model and a proportion of 1 will imply all models from the alternate model space outperform the benchmark model’s prediction. The proportions below the header $\sigma_{w1,w2}^2$ and $\tau_{w1,w2}$ are for the total variance and the secular variance while employing an unconstrained weighting scheme, whereas below headers σ_{w2}^2 and τ_{w2} the same is provided while using a constrained weighting scheme in the MIDAS smoothing.

	$\sigma_{w1,w2}^2$	$\tau_{w1,w2}$	σ_{w2}^2	τ_{w2}
	<i>RC</i>	<i>RC</i>	<i>RC</i>	<i>RC</i>
France	0.000	0.615	0.000	0.192
Germany	0.000	0.385	0.000	0.192
The UK	0.000	0.231	0.000	0.192
The US	0.000	0.038	0.000	0.000

Table 7 Multiple testing for pseudo out-of-sample MSE losses

The table presents the p-values of the White's (2000) reality check (RC) test and the Hansen's (2005) superior predictive ability (SPA) test. The forecasting errors i.e. MSE of the competing models are compared against the GM-RV benchmark model using the pseudo out-of-sample model forecasts following Engle et al. (2013). The sample period is January 1999- December 2016. Panels A, B, C and D present results for France, Germany, the UK and the US, respectively. First, we present naïve p-values that are the bootstrap RC p-values for the bi-model predictive ability test: bootstrapped p-values by comparing forecasting errors of the best model among $k = 26$ competing models relative to the benchmark model. Following Hansen, we compute three test statistics for the consistent (centre, c), upper (u) and lower (l) bounds for both RC test and Hansen test. Each test-statistic p-value is computed from 1000 bootstraps resamples and smoothing parameter $q = 0.25$. For clarity a p-value larger than 0.05 shows that no competing model outperforms the benchmark model i.e. null hypothesis cannot be rejected, whereas value <0.05 shows that at least one model has superior predictive ability (lower MSE or forecasting errors) than the benchmark GM-RV model, i.e. null is rejected. The $\sigma_{w1,w2}^2$ and the $\tau_{w1,w2}$ refers to the testing results of the forecasting errors of daily total and monthly long run volatilities when the weights for the beta polynomial for each input variable in the MIDAS filter are estimated, while σ_{w2}^2 and τ_{w2} refers to the testing results when $w_1 = 1$ and w_2 is estimated for the beta polynomial of each input variable in the MIDAS filter, respectively. These tests are carried out using daily squared returns and monthly RV as the latent variance proxies to calculate the loss functions for the full set of models across the four markets for the 26 competing models and the benchmark model GM-RV for each equity market index.

RP metric	France	$\sigma_{w1,w2}^2$ Germany	UK	US	$\tau_{w1,w2}$ France	Germany	UK	US
Panel A: flexible weighting scheme								
Naive	0.388	0.350	0.244	0.250	0.000	0.000	0.000	0.005
RC_l	0.766	0.905	0.574	0.690	0.020	0.001	0.002	0.004
RC_c	0.881	0.982	0.786	0.924	0.033	0.002	0.003	0.004
RC_u	0.883	0.982	0.786	0.924	0.036	0.004	0.004	0.004
$Hansen_l$	0.768	0.898	0.602	0.603	0.007	0.000	<0.001	0.004
$Hansen_c$	0.887	0.978	0.778	0.753	0.008	0.001	<0.001	0.004
$Hansen_u$	0.890	0.979	0.778	0.853	0.008	0.002	<0.001	0.004
Panel B: restricted weighting scheme								
		σ_{w2}^2					τ_{w2}	
Naive	0.130	0.600	0.420	0.540	0.001	0.001	0.020	0.143
RC_l	0.227	0.850	0.860	0.950	0.007	0.002	0.026	0.671
RC_c	0.470	0.800	0.760	0.880	0.008	0.003	0.039	0.720
RC_u	0.470	0.810	0.790	0.860	0.009	0.004	0.044	0.741
$Hansen_l$	0.134	0.780	0.550	0.400	0.003	0.000	0.021	0.467
$Hansen_c$	0.134	0.780	0.670	0.710	0.003	0.001	0.022	0.508
$Hansen_u$	0.134	0.750	0.630	0.780	0.005	0.002	0.029	0.526

Table 8 Multiple testing for pseudo out-of-sample MSE losses using alternate variance proxies

The table presents the p-values of the White's (2000) reality check (RC) test and the Hansen's (2005) superior predictive ability (SPA) test when we use variance proxies of GARCH (1, 1) i.e. $h_{i,t}$ and the square of implied volatility indices i.e. IV^2 for markets replacing RV. Panels A, B, C and D present results for France, Germany, the UK and the US, respectively. First, we present naïve p-values that are the bootstrap RC p-values for bi-model predictive ability test: bootstrapped p-values by comparing forecasting errors of the best model among $k = 26$ competing models relative to the benchmark model. Following Hansen (2005), we compute three test statistics for the consistent (centre, c), upper (u) and lower (l) bounds for both RC test and Hansen test. Each test-statistic p-value is computed from 1000 bootstraps resamples and smoothing parameter $q = 0.25$. For clarity a p-value larger than 0.05 shows that no competing model outperforms the benchmark model i.e. the null hypothesis cannot be rejected, whereas value <0.05 shows that at least there is one model that has superior predictive ability (lower MSE or forecasting errors) than the benchmark GM-RV model i.e. null is rejected. The $\sigma_{w1,w2}^2$ and the $\tau_{w1,w2}$ refers to the results of the forecasting errors of daily total and monthly long run volatilities when the weights for the beta polynomial for each input variable in the MIDAS filter are estimated, while σ_{w2}^2 and τ_{w2} refers to the model comparison results when $w_1 = 1$ and w_2 is estimated for the beta polynomial of each input variable in the MIDAS filter, respectively.

RP metirc	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		σ_{w2}^2		τ_{w2}	
	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2
Panel A: France								
Naive	0.367	0.454	0.001	0.001	0.143	0.208	0.000	0.000
RC_l	0.731	0.799	0.013	0.023	0.231	0.338	0.007	0.002
RC_c	0.873	0.910	0.029	0.036	0.511	0.603	0.014	0.002
RC_u	0.873	0.910	0.029	0.036	0.511	0.603	0.016	0.003
$Hansen_l$	0.722	0.793	0.002	0.002	0.147	0.212	<0.001	<0.001
$Hansen_c$	0.877	0.916	0.004	0.006	0.147	0.214	0.001	0.001
$Hansen_u$	0.877	0.916	0.007	0.007	0.147	0.214	0.001	0.002
Panel B: Germany								
Naive	0.304	0.335	0.002	0.000	0.450	0.590	0.000	0.001
RC_l	0.870	0.892	0.003	0.001	0.800	0.790	0.003	0.003
RC_c	0.976	0.976	0.007	0.003	0.810	0.780	0.005	0.004
RC_u	0.976	0.976	0.009	0.005	0.830	0.820	0.005	0.004
$Hansen_l$	0.871	0.879	0.003	0.001	0.670	0.740	<0.001	0.001
$Hansen_c$	0.972	0.968	0.003	0.002	0.710	0.650	<0.001	0.001
$Hansen_u$	0.972	0.968	0.004	0.002	0.790	0.710	0.001	0.002
Panel C: the UK								
Naive	0.217	0.204	0.000	0.000	0.350	0.390	0.023	0.019
RC_l	0.531	0.488	0.003	0.001	0.860	0.830	0.021	0.013
RC_c	0.781	0.743	0.003	0.002	0.850	0.780	0.031	0.023
RC_u	0.781	0.743	0.003	0.002	0.860	0.840	0.035	0.026
$Hansen_l$	0.586	0.533	0.002	<0.001	0.640	0.620	0.024	0.002
$Hansen_c$	0.752	0.729	0.002	<0.001	0.660	0.680	0.026	0.003
$Hansen_u$	0.752	0.729	0.002	<0.001	0.670	0.650	0.029	0.006
Panel D: the USA								
Naive	0.272	0.241	0.005	0.004	0.510	0.480	0.140	0.078
RC_l	0.737	0.693	0.004	0.004	0.870	0.880	0.666	0.556
RC_c	0.935	0.910	0.004	0.005	0.930	0.920	0.705	0.630
RC_u	0.935	0.910	0.005	0.005	0.930	0.890	0.717	0.639
$Hansen_l$	0.613	0.612	0.004	0.005	0.680	0.740	0.463	0.255
$Hansen_c$	0.799	0.691	0.004	0.004	0.700	0.530	0.506	0.309
$Hansen_u$	0.839	0.850	0.005	0.005	0.830	0.460	0.526	0.322

Table 9 Robust pairwise testing for out-of-sample MSE losses

In this table, we report the proportion of variance forecasts – total and long run – that outperform the benchmark model’s corresponding forecasts using the naïve p-values. The naïve p-values are the bi-model – a variant of DM test – bootstrapped RC test values that are computed considering that the best model is the only model in the competing model space. Thus, we count all the instances for variance forecasts and divide them by the total number of competing models (i.e., 26 in this case). To exemplify a proportion of 0 implies that there no model that outperforms the forecast given by the benchmark model and a proportion of 1 will imply all models from the alternate model space outperform the benchmark model’s prediction. The proportions below the header $\sigma_{w1,w2}^2$ and $\tau_{w1,w2}$ are for the total variance and the secular variance while employing an unconstrained weighting scheme, whereas below headers σ_{w2}^2 and τ_{w2} the same is provided while using a constrained weighting scheme in the MIDAS smoothing.

	$\sigma_{w1,w2}^2$	$\tau_{w1,w2}$	σ_{w2}^2	τ_{w2}
	<i>RC</i>	<i>RC</i>	<i>RC</i>	<i>RC</i>
France	0.615	0.462	0.615	0.577
Germany	0.577	0.385	0.654	0.538
The UK	0.577	0.500	0.654	0.654
The US	0.500	0.423	0.615	0.615

Table 10 Multiple testing for out-of-sample MSE losses

The table presents the p-values of the White's (2000) reality check (RC) test and the Hansen's (2005) superior predictive ability (SPA) test. The forecasting errors i.e. MSE of the competing models are compared against the GM-RV benchmark model using a fixed out-of-sample forecasting scheme. Parameters are estimated using for sample period January 1999- December 2011 and forecasting period is January 2012-December 2016. Panels A, B, C and D present results for France, Germany, the UK and the US, respectively. First, we present naïve p-values that are the bootstrap RC p-values for bi-model predictive ability test: bootstrapped p-values by comparing forecasting errors of the best model among $k = 26$ competing models relative to the benchmark model. Following Hansen (2005), we compute three test statistics for the consistent (center, c), upper (u) and lower (l) bounds for both RC test and Hansen test. Each test-statistic p-value is computed from 1000 bootstraps resamples and smoothing parameter $q = 0.25$. For clarity a p-value larger than 0.05 shows that no competing model outperforms the benchmark model i.e. null hypothesis cannot be rejected, whereas value <0.05 shows that at least there is one model that has superior predictive ability (lower MSE or forecasting errors) than the benchmark GM-RV model i.e. null is rejected. The $\sigma_{w1,w2}^2$ and the $\tau_{w1,w2}$ refers to the testing results of the forecasting errors of daily total and monthly long run volatilities when the weights for the beta polynomial for each input variable in the MIDAS filter are estimated, while σ_{w2}^2 and τ_{w2} refers to the testing results when $w_1 = 1$ and w_2 is estimated for the beta polynomial of each input variable in the MIDAS filter, respectively. These tests are carried out using daily squared returns and monthly RV as the latent variance proxies to calculate the loss functions for the full set of models across the four markets for the 26 competing models and the benchmark model GM-RV for each equity market index.

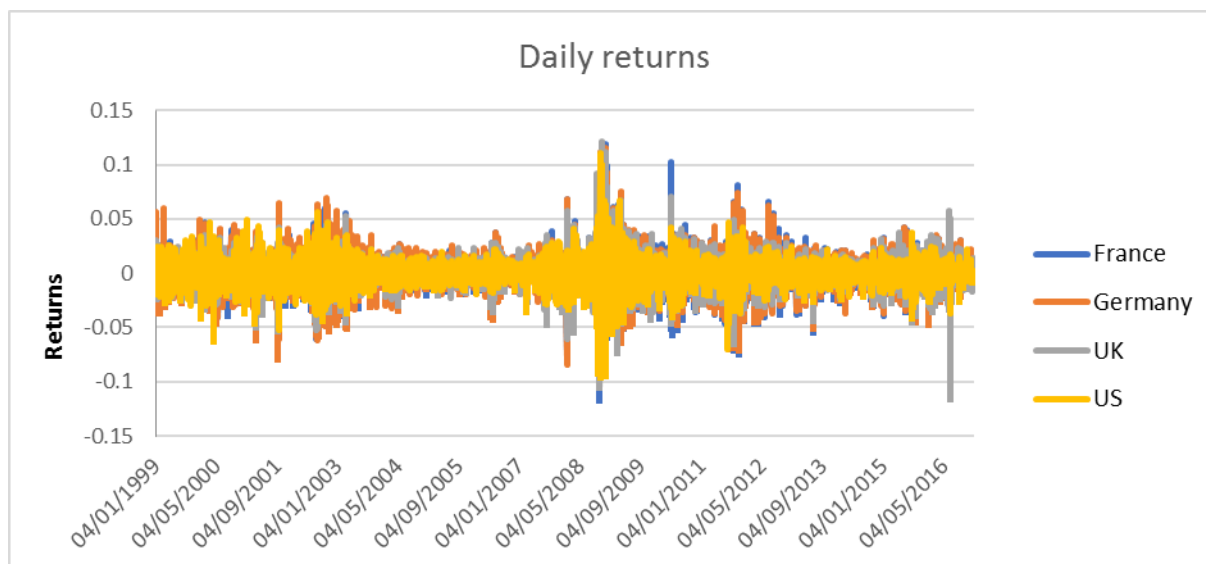
RP metric	$\sigma_{w1,w2}^2$				$\tau_{w1,w2}$			
	France	Germany	UK	US	France	Germany	UK	US
Panel A: flexible weighting scheme								
Naïve	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000
RC_l	0.096	0.380	0.061	0.094	0.334	0.002	0.061	0.063
RC_c	0.259	0.407	0.061	0.108	0.333	0.003	0.265	0.063
RC_u	0.435	0.440	0.230	0.444	0.407	0.004	0.308	0.333
$Hansen_l$	0.002	0.088	0.005	0.022	0.079	0.000	<0.001	<0.001
$Hansen_c$	0.004	0.088	0.005	0.022	0.083	0.001	<0.001	<0.001
$Hansen_u$	0.004	0.202	0.005	0.022	0.098	0.002	<0.001	<0.001
Panel B: restricted weighting scheme								
	σ_{w2}^2				τ_{w2}			
Naïve	0.000	0.600	0.451	0.000	0.000	0.000	0.000	0.000
RC_l	0.198	0.750	0.520	0.137	0.277	0.001	0.068	0.466
RC_c	0.209	0.800	0.526	0.143	0.234	0.003	0.127	0.478
RC_u	0.219	0.810	0.790	0.215	0.492	0.004	0.257	0.602
$Hansen_l$	0.108	0.780	0.594	0.096	0.092	0.000	0.004	0.111
$Hansen_c$	0.108	0.780	0.626	0.096	0.083	0.002	0.004	0.147
$Hansen_u$	0.109	0.793	0.641	0.159	0.147	0.002	0.007	0.234

Table 11 Multiple testing for out-of-sample MSE losses using alternate variance proxies

The table presents the p-values of the White's (2000) reality check (RC) test and the Hansen's (2005) superior predictive ability (SPA) test when we use variance proxies of GARCH (1, 1) i.e. $h_{i,t}$ and square of implied volatility indices i.e. IV^2 for respective markets replacing RV. The forecasting errors i.e. MSE of the competing models are compared against the GM-RV benchmark model using a fixed out-of-sample forecasting scheme. Parameters are estimated using for sample period January 1999- December 2011 and forecasting period is January 2012-December 2016. Panels A, B, C and D present results for France, Germany, the UK and the US, respectively. First we present naïve p-values that are the bootstrap RC p-values for bi-model predictive ability test: bootstrapped p-values by comparing forecasting errors of the best model among $k = 26$ competing models relative to the benchmark model. Following Hansen (2005) we compute three test statistics for the consistent (center, c), upper (u) and lower (l) bounds for both RC test and Hansen test. Each test-statistic p-value is computed from 1000 bootstraps resamples and smoothing parameter $q = 0.25$. For clarity a p-value larger than 0.05 shows that no competing model outperforms the benchmark model i.e. null hypothesis cannot be rejected, whereas value <0.05 shows that at least there is one model that has superior predictive ability (lower MSE or forecasting errors) than the benchmark GM-RV model i.e. null is rejected. The $\sigma_{w1,w2}^2$ and the $\tau_{w1,w2}$ refers to the results of the forecasting errors of daily total and monthly long run volatilities when the weights for the beta polynomial for each input variable in the MIDAS filter are estimated, while σ_{w2}^2 and τ_{w2} refers to the model comparison results when $w_1 = 1$ and w_2 is estimated for the beta polynomial of each input variable in the MIDAS filter, respectively.

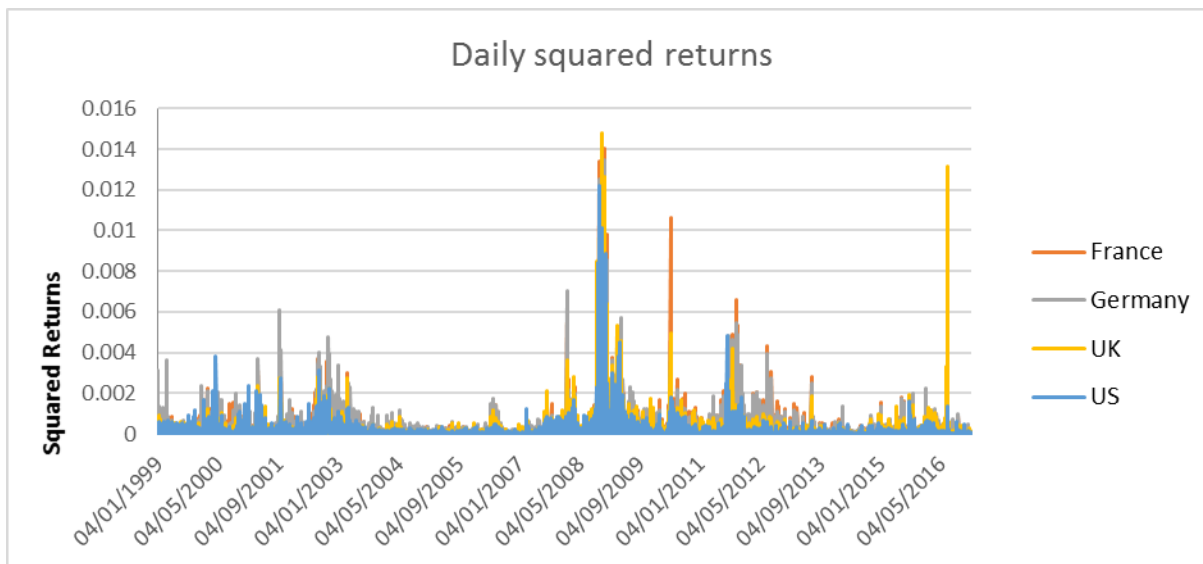
RP metirc	$\sigma_{w1,w2}^2$		$\tau_{w1,w2}$		σ_{w2}^2		τ_{w2}	
	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2	$h_{i,t}$	IV^2
Panel A: France								
Naive	0.001	0.034	0.000	0.001	0.000	0.001	0.000	0.000
RC_l	0.731	0.821	0.331	0.298	0.205	0.325	0.213	0.323
RC_c	0.873	0.773	0.331	0.345	0.205	0.232	0.267	0.317
RC_u	0.873	0.693	0.400	0.387	0.240	0.261	0.487	0.457
$Hansen_l$	0.222	0.342	0.101	0.092	0.118	0.131	0.081	0.079
$Hansen_c$	0.277	0.417	0.101	0.097	0.123	0.142	0.097	0.089
$Hansen_u$	0.277	0.439	0.101	0.108	0.126	0.166	0.147	0.134
Panel B: Germany								
Naive	0.000	0.002	0.000	0.000	0.450	0.423	0.002	0.003
RC_l	0.337	0.297	0.003	0.002	0.800	0.783	0.003	0.001
RC_c	0.413	0.421	0.005	0.002	0.810	0.880	0.007	0.004
RC_u	0.443	0.450	0.005	0.007	0.830	0.798	0.009	0.008
$Hansen_l$	0.073	0.077	<0.001	<0.002	0.670	0.621	0.003	0.001
$Hansen_c$	0.081	0.091	<0.001	<0.003	0.710	0.690	0.003	0.001
$Hansen_u$	0.116	0.120	0.001	0.007	0.790	0.743	0.004	0.008
Panel C: the UK								
Naive	0.000	0.002	0.000	0.001	0.350	0.330	0.000	0.001
RC_l	0.058	0.052	0.079	0.082	0.860	0.910	0.067	0.062
RC_c	0.059	0.057	0.137	0.101	0.850	0.920	0.153	0.193
RC_u	0.718	0.671	0.194	0.124	0.860	0.870	0.267	0.327
$Hansen_l$	0.006	0.002	<0.001	<0.001	0.640	0.662	0.004	0.003
$Hansen_c$	0.006	0.007	<0.001	0.002	0.660	0.711	0.004	0.003
$Hansen_u$	0.007	0.009	<0.001	0.002	0.670	0.742	0.004	0.001
Panel D: the USA								
Naive	0.000	0.001	0.000	0.001	0.000	0.000	0.000	0.001
RC_l	0.087	0.082	0.057	0.073	0.114	0.124	0.473	0.523
RC_c	0.137	0.141	0.057	0.092	0.133	0.154	0.504	0.484
RC_u	0.218	0.282	0.117	0.092	0.205	0.185	0.593	0.633
$Hansen_l$	0.020	0.031	<0.001	<0.001	0.073	0.077	0.108	0.098
$Hansen_c$	0.020	0.043	<0.001	0.002	0.073	0.081	0.149	0.123
$Hansen_u$	0.031	0.093	<0.001	0.004	0.437	0.399	0.411	0.401

Figure 1



The daily logarithmic returns are computed from the France, Germany, the UK and the US price indices and are presented for the sample period from January 1999 to December 2016.

Figure 2



The squared returns are computed from the France, Germany, the UK and the US price indices and are presented for the sample period from January 1999 to December 2016.

Appendix A

Table AI Summary statistics

	Index returns				Realized variance (RV)			
	France	Germany	UK	USA	France	Germany	UK	USA
average	0.010	0.014	-0.007	0.025	0.050	0.075	0.050	0.025
standard deviation	0.253	0.253	0.206	0.190	0.011	0.011	0.009	0.008
Skewness	-0.074	-0.080	-0.227	-0.201	10.043	8.529	12.495	12.144
Kurtosis	8.921	7.546	11.861	11.225	149.17	117.4	217.34	221.3
JB test	6862.7	4048.5	15400	13266	4.3e+06	2.2e+06	9.1e+06	9.4e+06
	[0.001]	[0.001]	[0.001]	[0.001]	[0.001]	[0.001]	[0.001]	[0.001]
Auto-correlations -1	-0.012	0.002	-0.010	-0.071	0.170	0.152	0.209	0.205
	[0.431]	[0.838]	[0.483]	[0.000]	[0.000]	[0.000]	[0.000]	[0.000]
Auto-correlations -3	-0.055	-0.029	-0.069	-0.0009	0.225	0.236	0.259	0.206
	[0.000]	[0.041]	[0.000]	[0.948]	[0.000]	[0.000]	[0.000]	[0.000]
Auto-correlations -6	-0.008	-0.003	-0.028	-0.004	0.175	0.149	0.199	0.281
	[0.570]	[0.841]	[0.054]	[0.782]	[0.000]	[0.000]	[0.000]	[0.000]
Auto-correlations -12	0.005	0.002	-0.0008	-0.002	0.238	0.236	0.260	0.284
	[0.712]	[0.874]	[0.953]	[0.901]	[0.000]	[0.000]	[0.000]	[0.000]