# Bayesian Adaptively Updated Hamiltonian Monte Carlo with an Application to High-Dimensional BEKK GARCH Models[*]

Martin Burda,[†]      John M. Maheu [‡]

October 4, 2011

---

## Abstract

Hamiltonian Monte Carlo (HMC) is a recent statistical procedure to sample from complex distributions. Distant proposal draws are taken in a sequence of steps following the Hamiltonian dynamics of the underlying parameter space, often yielding superior mixing properties of the resulting Markov chain. However, its performance can deteriorate sharply with the degree of irregularity of the underlying likelihood due to its lack of local adaptability in the parameter space. Riemann Manifold HMC (RMHMC), a locally adaptive version of HMC, alleviates this problem, but at a substantially increased computational cost that can become prohibitive in high-dimensional scenarios. In this paper we propose the Adaptively Updated HMC (AUHMC), an alternative inferential method based on HMC that is both fast and locally adaptive, combining the advantages of both HMC and RMHMC. The benefits become more pronounced with higher dimensionality of the parameter space and with the degree of irregularity of the underlying likelihood surface. We show that AUHMC satisfies detailed balance for a valid MCMC scheme and provide a comparison with RMHMC in terms of effective sample size, highlighting substantial efficiency gains of AUHMC. Simulation examples and an application of the BEKK GARCH model show the practical usefulness of the new posterior sampler.

*JEL codes:* C01; C11; C15; C32
*Keywords: High-dimensional joint sampling, Markov chain Monte Carlo*

---

[†]Department of Economics, University of Toronto, 150 St. George St., Toronto, ON M5S 3G7, Canada; Phone: (416) 978-4479; Email: `martin.burda@utoronto.ca`
[‡]Department of Economics, University of Toronto, 150 St. George St., Toronto, ON M5S 3G7, Canada; Phone: (416) 978-1495; Email: `jmaheu@chass.utoronto.ca` and RCEA, Italy.

## 1. Introduction

Hamiltonian dynamics have been traditionally used to describe the laws of motion in molecular systems in physics. Following the recent advances in Markov chain Monte Carlo (MCMC) fuelled by increasing availability of fast computation, inferential methods based on Hamiltonian dynamic systems are becoming increasingly popular in the statistics literature (Neal, 1993, 2010; Ishwaran, 1999; Liu, 2004; Girolami and Calderhead, 2011). Hamiltonian Monte Carlo, also called Hybrid Mote Carlo, (HMC) uses Hamiltonian dynamics in constructing distant proposal draws in a sequence of steps and hence concurrently yields relatively low correlation among draws and high acceptance probabilities. Methods based on HMC have been shown to improve sampling of ill-behaved posteriors, and enabled the solution of otherwise intractable high dimensional inference problems (Neal, 2010; Girolami and Calderhead, 2011). These methods are particularly useful for the kind of problems where it is difficult to accurately approximate the surface of the (posterior) log-likelihood around the current parameter draw or the mode in real time needed for obtaining sufficiently high acceptance probabilities in importance sampling (IS) or accept-reject methods. Perpetual re-fitting of a local posterior approximating density around newly accepted draws during the MCMC run may become too costly for methods based on such mechanism to be practical. These types of problems typically arise when the log-likelihood is costly to evaluate and is near-ill-conditioned around the mode.

Even if on a small scale, with a few parameters and small sample size, such problems can be handled by standard procedures, these can become prohibitive in higher parameter dimensions and sample sizes. Examples include recursive models in finance, such as the BEKK GARCH that we treat in our application, state-space models or point process models. In such situations one would typically resort to Random walk (RW) style sampling that is fast to run and does not require the knowledge of the properties of the underlying log-likelihood. However, RW mechanisms can lead to very slow exploration of the parameter space with high autocorrelations among draws which would require a prohibitively large size of the Markov chain to be obtained in implementation to achieve satisfactory mixing and convergence. HMC combines the advantages of sampling that is relatively cheap with RW-like intensity but superior parameter space exploration.

Nonetheless, HMC uses a mechanism whose form is fixed over the parameter space, lacking adaptability to local features of the likelihood. The Riemann Manifold HMC, or RMHMC (Girolami and Calderhead, 2011), alleviates this problem and renders HMC locally adaptable which results in improved convergence and mixing properties. However, relative to HMC, RMHMC implementation requires a substantially increased computational burden with a large number of fixed point evaluations within every MC step. Crucially, a numerical estimate of the Fisher Information matrix needs to be newly evaluated in every iteration while searching for each fixed point. This can render its performance inadequate in high-dimensional problems where the likelihood is expensive to evaluate. Indeed, it is precisely this type of problems for which HMC-type methods are most useful relative to other existing methods.

In this paper we propose an alternative inferential method, the Adaptively Updated HMC (AUHMC), that is both relatively fast and locally adaptive. AUHMC is based on proposal dynamics generalizing HMC with only minimal additional functional evaluations, approximating the local adaptability properties of RMHMC. Unlike the RMHMC, AUHMC does not attempt to construct a completely locally adaptive proposal sequence, but rather a fast local approximation to the fully adaptive case. This enables AUHMC to bypass multiple fixed point evaluations in every step in the proposal sequence within every MC parameter draw that RMHMC needs to take. As a result, AUHMC features a substantial speed gain traded off for only a relatively small loss of the degree of adaptability relative to RMHMC.

From the end-user perspective AUHMC is easier to code than RMHMC, while the additional elements over HMC are simple to implement. AUHMC is *not* a special case of RMHMC as their dynamic systems are non-nested, while HMC can be obtained as a special case of AUHMC by imposing restrictions on the dynamics of the latter.

We provide a set of necessary and sufficient conditions under which AUHMC yields a valid MCMC scheme with a tractable form of its acceptance probability. The performance of AUHMC is assessed on two simulated examples: first a case with increasing dimensionality of the parameter space and fixed sample size (multivariate Normal), and second a case with increasing sample size and fixed dimensionality (GARCH(1,1)). Both examples reveal increasing relative efficiency gains of AUHMC.

We apply AUHMC to the task of model comparison in a high-dimensional BEKK GARCH environment with its highly complex likelihood. We show that AUHMC facilitates evaluation of the marginal likelihood even in the joint likelihood full BEKK GARCH model in higher dimensions than previously considered practical. Due to the inherent sampling difficulties, Bayesian estimation of multivariate GARCH models is relatively scarce (Dellaportas and Vrontos (2007), Hudson and Gerlach (2008) and Osiewalski and Pipien (2004)). Coming up with a good proposal density inside a Metropolis-Hasting procedure has been a challenge for conventional samplers. The importance of full joint likelihood BEKK inference is highlighted by a marginal likelihood comparison that clearly favors the full model version over its restricted alternatives.

AUHMC is related to but distinct from the adaptive radial-based direction sampling (ARDS) method of Bauwens, Bos, van Dijk, and van Oest (2004). While AUHMC utilizes deterministic directional derivatives (numerical or analytical) of a Hamiltonian system in order to move within hypersurfaces of approximately equal functional value, ARDS is based on a transformation into radial coordinates, stochastic sampling of directional vectors, and then applying the inverse transformation. The acceptance probability of the Metropolis-Hastings version of ARDS is a function of a numerical quadrature over the posterior in a given direction. The importance sampling version of ARDS relies on a directional approximation of the posterior. In either case, each MC draw of ARDS requires a certain type

of relatively detailed posterior approximation which AUHMC seeks to avoid in order to be applicable in problems where quadrature evaluation or importance sampling may become computationally prohibitive, as described above. Each method thus focuses on different types of applied problems.

Our work also complements other existing tailored proposal methods for posterior sampling in difficult situations such as Chib and Ramamurthy (2010), Liesenfeld and Richard (2006) and Pitt and Shephard (1997). The AUHMC is a useful addition to the applied econometrician's toolkit and can be applied to the full block of parameters as in our examples or to a sub-block of parameters in conjunction with other Gibbs and Metropolis-Hasting steps.

The paper is organized as follows: Section 2 provides an overview of useful statistical background including the detailed balance condition of the Metropolis-Hastings principle. Section 3 introduces AUHMC. Section 4 explores the properties of AUHMC on simulated examples and Section 5 details the application of AUHMC to a high-dimensional BEKK GARCH model. Section 6 concludes.

## 2. **Statistical Background**

Consider an economic model parametrized by a Euclidean vector $\theta \in \Theta$ for which all information in the sample is contained in the model posterior $\pi(\theta; \cdot)$ that we denote by $\pi(\theta)$, assumed known up to an integrating constant. Formally, a general class of such models can be characterized by a family $\mathcal{P}_\theta$ of probability measures on a measurable space $(\Theta, \mathcal{B})$ where $\mathcal{B}$ is the Borel $\sigma-$algebra.

The purpose of Markov Chain Monte Carlo (MCMC) methods is to formulate a Markov chain on the parameter space $\Theta$ for which, under certain conditions, $\pi(\theta) \in \mathcal{P}_\theta$ is the invariant (also called 'equilibrium' or 'long-run') distribution. The Markov chain of draws of $\theta$ can be used to construct simulation-based estimates of the required integrals, and functionals $h(\theta)$ of $\theta$ that are expressed as integrals. These functionals include objects of interest for inference on $\theta$ such as quantiles of $\pi(\theta)$.

The Markov chain sampling mechanism specifies a method for generating a sequence of random variables $\{\theta_r\}_{r=1}^R$, starting from an initial point $\theta_0$, in the form of conditional distributions for the draws $\theta_{r+1}|\theta_r \sim G(\theta_r)$. Under relatively weak regularity conditions (Robert and Casella, 2004), the average of the Markov chain converges to the expectation under the stationary distribution:

$$\lim_{R \to \infty} \frac{1}{R} \sum_{r=1}^R h(\theta_r) = E_\pi[h(\theta)]$$

A Markov chain with this property is called ergodic. As a means of approximation we rely on large but finite $R \in \mathbb{N}$ which the analyst has the discretion to select in applications.

The Metropolis-Hastings (M-H) principle has been the cornerstone of constructing Markov chains by sampling $\theta_{r+1}|\theta_r$ from $G(\theta_r)$; see Chib and Greenberg (1995) for a detailed overview. $G(\theta_r)$ can be obtained from a given (economic) model and its corresponding posterior $\pi(\theta)$, parametrized by $\theta$, known up to a constant of proportionality.

However, $\pi(\theta)$ typically has a complicated form which precludes direct sampling. Then the goal is to find a transition kernel $P(\theta, d\theta)$ whose $n$th iterate converges to $\pi(\theta)$ for large $n$. After this large number, the distribution of the observations generated from the Markov chain simulation is approximately the target distribution. The transition kernel $P(\theta, A)$ for $\theta \in \Theta$ and $A \subset \Theta$ is an unknown conditional distribution function that represents the probability of moving from $\theta$ to a point in the set $A$. Suppose we have a proposal-generating density $q(\theta^*_{r+1}|\theta_r)$ where $\theta^*_{r+1}$ is a proposed state given the current state $\theta_r$ of the Markov chain. The Metropolis-Hastings (M-H) principle stipulates that $\theta^*_{r+1}$ be accepted as the next state $\theta_{r+1}$ with the acceptance probability

$$(2.1) \qquad \alpha(\theta_r, \theta^*_{r+1}) = \min\left[\frac{\pi(\theta^*_{r+1})q(\theta_r|\theta^*_{r+1})}{\pi(\theta_r)q(\theta^*_{r+1}|\theta_r)}, 1\right]$$

otherwise $\theta_{r+1} = \theta_r$. Then the Markov chain satisfies the so-called detailed balance condition

$$\pi(\theta_r)q(\theta^*_{r+1}|\theta_r)\alpha(\theta_r, \theta^*_{r+1}) = \pi(\theta^*_{r+1})q(\theta_r|\theta^*_{r+1})\alpha(\theta^*_{r+1}, \theta_r)$$

which is sufficient for ergodicity. $\alpha(\theta^*_{r+1}, \theta_r)$ is the probability of the move $\theta_r|\theta^*_{r+1}$ if the dynamics of the proposal generating mechanism were to be reversed. While $\pi(\theta)$ may be difficult or expensive to sample from, the proposal-generating density $q(\theta^*_{r+1}|\theta_r)$ can be chosen to be sampled easily. The popular Gibbs sampler arises as a special case when the M-H sampler is factored into conditional densities.

A variation on (2.1) can be constructed by augmenting the parameter space $\Theta$ with a set of independent auxiliary stochastic parameters $\gamma \in \Gamma$ that fulfill a supplementary role in the proposal algorithm, such as facilitating the directional guidance of the proposal mechanism. The detailed balance is then satisfied using the acceptance probability

$$(2.2) \qquad \alpha(\theta_r, \gamma_r; \theta^*_{r+1}, \gamma^*_{r+1}) = \min\left[\frac{\pi(\theta^*_{r+1}, \gamma^*_{r+1})q(\theta_r, \gamma_r|\theta^*_{r+1}, \gamma^*_{r+1})}{\pi(\theta_r, \gamma_r)q(\theta^*_{r+1}, \gamma^*_{r+1}|\theta_r, \gamma_r)}, 1\right]$$

A further relevant variation on (2.1) is arises when $(\theta^*_{r+1}, \gamma^*_{r+1})$ are obtained from $(\theta_r, \gamma_r)$ using a sequence of within-proposal steps $\{\theta^k_r, \gamma^k_r\}^L_{k=1}$ with $(\theta_r, \gamma_r) = (\theta^0_r, \gamma^0_r)$ and $(\theta_r, \gamma_r) = (\theta^L_r, \gamma^L_r)$. In each case, the desired posterior can be obtained by marginalizing out $\gamma$.

## 3. Adaptively Updated Hamiltonian Monte Carlo

The original Hamiltonian (or Hybrid) Monte Carlo (HMC) algorithm has its roots in the physics literature where it was introduced as a fast method for simulating molecular dynamics (Duane, Kennedy, Pendleton, and Roweth, 1987). It has since become popular in a number of application areas including statistical physics (Akhmatskaya, Bou-Rabee, and Reich, 2009; Gupta, Kilcup, and Sharpe, 1988), computational chemistry (Tuckerman, Berne, Martyna, and Klein, 1993), or a generic tool for Bayesian statistical inference (Neal, 1993, 2010; Ishwaran, 1999; Liu, 2004; Beskos, Pillai, Roberts, Sanz-Serna, and Stuart, 2010). A separate stream of literature has developed around the Langevin diffusion mechanisms which use related proposal dynamics but utilize one-step proposals only (Roberts

and Rosenthal, 1998; Roberts and Stramer, 2003). We synthesize the HMC principles in a generally accessible form in Appendix A.

In a recent contribution to the statistics literature, Girolami and Calderhead (2011) generalize HMC to benefit from Riemannian geometry induced by the expected Fisher Information mass matrix in the HMC algorithm. The use of the Fisher Information metric tensor results in effective moves based on shortest paths (geodesics) across the induced Riemannian manifold. The geodesics across a Riemannian manifold may be described in terms of Hamilton's equations, thus providing a natural link to HMC methods. Since the Fisher Information metric is a function of the model parameters, the resulting method (RMHMC) renders the HMC algorithm adaptive to the local curvature of the posterior likelihood, in contrast to the original HMC with a constant mass matrix and hence a fixed metric over the whole parameter space.

However, relative to HMC, RMHMC implementation requires a substantially increased computational burden, resulting from the additional requirement of finding numerical solutions to two fixed points at every step $k$ of the proposal sequence $\{\theta_r^k, \gamma_r^k\}_{k=1}^L$ inside each HMC proposal step. Crucially, a numerical estimate of the Fisher Information matrix needs to be newly evaluated in every iteration while searching for each fixed point. This can render its performance inadequate in high-dimensional problems where the likelihood is expensive to evaluate.

In this paper we propose the Adaptively Updated HMC (AUHMC), an alternative HMC-based method featuring distant proposals that is locally adaptable and yet avoids the computational complexity of RMHMC. AUHMC uses an approximation to the Riemanian geometry utilized by RMHMC that is much cheaper to obtain than the exact Riemannian paths. However, AUHMC does *not* constitute a special case of RMHMC, since no simplification of the latter will yield the former. What is being simplified here is the metric tensor geometry of Riemannian manifold over which moves are proposed, which requires a distinct non-nested implementation algorithm from the previously proposed ones.

We show that AUHMC satisfies the conditions for a valid MCMC scheme in Theorem 1 below. Results of this type have been obtained for HMC and RMHMC in the literature, but the AUHMC is a non-nested distinct alternative to either of these methods and hence needs to be validated separately. Theorem 2 further provides the set of regularity conditions on the (posterior) likelihood that are sufficient for satisfying the assumptions made in Theorem 1. These conditions can be easily verified in a given application.

### 3.1. Non-separable Hamiltonian Systems

Consider a vector of parameters of interest $\theta \in \mathbb{R}^d$ distributed according to the posterior density $\pi(\theta)$. Let $\gamma \in \mathbb{R}^d$ denote a vector of auxiliary parameters with $\gamma \sim \Phi(\gamma; 0, M(\theta))$ where $\Phi$ denotes the Gaussian distribution with mean vector 0 and covariance matrix $M(\theta)$. Denote the joint density of $(\theta, \gamma)$ by $\pi(\theta, \gamma)$. Then the negative of the logarithm of the joint density of $(\theta, \gamma)$ is given by the

Hamiltonian equation

(3.1)
$$H(\theta, \gamma) = -\ln \pi(\theta) - \ln q(\gamma|\theta)$$

where

(3.2)
$$q(\gamma|\theta) = (2\pi)^{-d/2} |M(\theta)|^{-1/2} \exp\left(-\frac{1}{2}\gamma' M(\theta)^{-1}\gamma\right)$$

renders the auxiliary parameter quadratic term $\gamma' M(\theta)^{-1}\gamma/2$ as an explicit function of $\theta$ (Leimkuhler and Reich, 2004). This property leads to local adaptability of the proposal sequence but also complicates subsequent analysis. The associated Hamiltonian dynamics equations are in general given by

(3.3)
$$\frac{d\theta_i}{dt} = \frac{\partial H(\theta, \gamma)}{\partial \gamma_i} = \left[M(\theta)^{-1}\gamma\right]_i$$

$$\frac{d\gamma_i}{dt} = -\frac{\partial H(\theta, \gamma)}{\partial \theta_i} = \nabla_\theta \ln \pi(\theta) - \frac{1}{2}\mathrm{Tr}\left(M(\theta)^{-1}\frac{\partial M(\theta)}{\partial \theta_i}\right)$$

(3.4)
$$+ \frac{1}{2}\gamma' M(\theta)^{-1}\frac{\partial M(\theta)}{\partial \theta_i}M(\theta)^{-1}\gamma$$

A number of numerical methods have been devised in the physics and molecular dynamics literature to solve the differential equations (3.3)–(3.4) in order to accurately determine the position of $\theta(t+s)$ and $\gamma(t+s)$ at the next instant $t+s$ given their current position at time $t$ in the state space. These solutions include the generalized Euler and Stormer-Verlet (generalized leapfrog) methods (Hairer, Lubich, and Wanner, 2003; Leimkuhler and Reich, 2004).

## 3.2. AUHMC

The starting point for AUHMC is the non-separable Hamiltonian (3.1)-(3.2). However, instead of $M(\theta)$, for each MCMC update $r$, we use the matrix $M(\overline{\theta_r, \theta_{r+1}^*})$ that is *fixed constant* for the entire leapfrog multi-step proposal sequence $\{\theta_r^k, \gamma_r^k\}_{k=0}^L$, i.e. between $\theta_r$ and $\theta_{r+1}^*$ inclusive. Thus, for a given $r$, $M(\overline{\theta_r, \theta_{r+1}^*})$ is not a function of $\theta$. Hence the formation of the proposal sequence $\{\theta_r^k, \gamma_r^k\}_{k=1}^L$ can use the standard leapfrog integrator of HMC, with the mass matrix $M(\overline{\theta_r, \theta_{r+1}^*})$, given by,

(3.5)
$$\gamma_r^{k+1/2} = \gamma_r^k - \frac{\varepsilon}{2}\nabla_\theta \ln \pi(\theta_r^k)$$

(3.6)
$$\theta_r^{k+1} = \theta_r^k + \varepsilon\left[M(\overline{\theta_r, \theta_{r+1}^*})^{-1}\gamma_r^{k+1/2}\right]$$

(3.7)
$$\gamma_r^{k+1} = \gamma_r^{k+1/2} - \frac{\varepsilon}{2}\nabla_\theta \ln \pi(\theta_r^{k+1})$$

The value of $M(\overline{\theta_r, \theta_{r+1}^*})$ is obtained as one fixed point per proposal draw $(\theta_{r+1}^*, \gamma_{r+1}^*)$. Given $\theta_r$, and an initial guess $M(\overline{\theta_r, \theta_{r+1,j}^*})$, for $j = 1$, take $L$ steps of (3.5)-(3.7) with $k = 1, \ldots, L$, then at the resulting $\theta_{r+1,j}^*$ update $M(\overline{\theta_r, \theta_{r+1,j+1}^*})$, and keep iterating for $j = 2, \ldots$ until convergence to a fixed point $M(\overline{\theta_r, \theta_{r+1}^*})$ achieved when $(\theta_{r+1,j-1}^*, \gamma_{r+1,j-1}^*) = (\theta_{r+1,j}^*, \gamma_{r+1,j}^*)$ within some small tolerance region. The proposal sequence $\{\theta_r^k, \gamma_r^k\}_{k=1}^L$ with $(\theta_{r+1}^*, \gamma_{r+1}^*) = (\theta_r^L, \gamma_r^L)$ is then drawn by applying (3.5)-(3.7) using the the fixed-point $M(\overline{\theta_r, \theta_{r+1}^*})$. The exact AUHMC algorithm is given in Appendix B. The conditions for a contraction mapping given below ensure the existence and uniqueness of

the fixed point. In our experiments we found that only a few iterations were necessary to obtain $M(\overline{\theta_r, \theta^*_{r+1}})$, resulting in relatively rapid speed of the MCMC updates.

While $M(\overline{\theta_r, \theta^*_{r+1}})$ is fixed *within* a given MCMC proposal step from $\theta_r$ to $\theta^*_{r+1}$, $M(\overline{\theta_r, \theta^*_{r+1}})$ changes *between* any two distinct proposal steps $\{\theta_r$ to $\theta^*_{r+1}\}$ and $\{\theta_s$ to $\theta^*_{s+1}\}$ with $r \neq s$. This feature renders AUHMC adaptive to the curvature of the (posterior) likelihood for any current parameter draw $\theta_r$ over the parameter space $\Theta$.

We expect AUHMC to be more computationally efficient relative to RMHMC. First we use (3.5)-(3.7) instead of the more complex integrator associated with (3.3)-(3.4). Second, we require one fixed point in contrast to many fixed points along the proposal path in $k = 1, ..., L$. We verify this assertion numerically in the simulated examples below.

The following assumptions state sufficient conditions for AUHMC to satisfy detailed balance of a valid MCMC scheme, with a tractable acceptance probability.

**ASSUMPTION 1.** $M(\overline{\theta_r, \theta^*_{r+1}})$ *is symmetric in its arguments, satisfying*

$$M(\overline{\theta_r, \theta^*_{r+1}}) = M(\overline{\theta^*_{r+1}, \theta_r})$$

**ASSUMPTION 2.** $\nabla_\theta \ln \pi(\theta)$ *is bounded and Lipschitz continuous in $\theta$.*

**ASSUMPTION 3.** *The parameter space $\Theta$ is compact.*

Assumption 1 can be satisfied by construction when setting the functional form of $M(\cdot)$, as we do below. Assumption 2 imposes restrictions on the rate of change of the score function that are satisfied by smooth densities typically used to construct likelihood functions. Assumption 3 is standard in the literature.

The uniqueness of the AUHMC solution and its detailed balance are summarized by the following two results.

**LEMMA 1.** *Under the Assumptions 2–3, the fixed point defining $M(\overline{\theta_r, \theta^*_{r+1}})$ exists and is unique for any given $\theta_r$. In particular, for any $\delta \in (0, 1)$ there exists $\varepsilon(\delta) > 0$ dependent on $\delta$ only, such that $\forall \varepsilon^* < \varepsilon(\delta)$, $\{T_k\}^L_{k=1}$ is a contraction mapping uniquely determining $M(\overline{\theta_r, \theta^*_{r+1}})$.*

**THEOREM 1.** *Under the Assumptions 1–3, AUHMC satisfies detailed balance, with the acceptance probability given by*

$$\alpha(\theta_r, \gamma_r; \theta^*_{r+1}, \gamma^*_{r+1}) = \min\left[\frac{\pi(\theta^*_{r+1}, \gamma^*_{r+1})}{\pi(\theta^0_r, \gamma^0_r)}, 1\right]$$

(3.8)
$$= \min\left[\exp(\widetilde{\alpha}_r), 1\right]$$

*where*

$$\widetilde{\alpha}_r = \ln \pi(\theta^*_{r+1}) - \ln \pi(\theta_r) + \ln \phi(\gamma^*_{r+1}; 0, M(\overline{\theta_r, \theta^*_{r+1}})) - \ln \phi(\gamma_r; 0, M(\overline{\theta_r, \theta^*_{r+1}}))$$

The proofs are provided in Appendix C. Heuristically, we show that AUHMC implements a solution to a symmetric mapping $\widehat{\Psi}_\varepsilon$, defined in Appendix C. The symmetry of $\widehat{\Psi}_\varepsilon$ implies its time reversibility which in turn yields detailed balance. The proof of Theorem 1 closely follows the proof of symmetry of a concatenation of an explicit Euler method with an implicit Euler method (Leimkuhler and Reich, 2004, p. 84), but in defining the implicit part of each step $\widehat{\Psi}_\varepsilon$ uses the distances to the endpoints of the proposal sequence instead of the arguments of directional derivatives used in the Euler case. The resulting concatenation of explicit and implicit half-steps leading to $\widehat{\Psi}_\varepsilon$ is symmetric and hence reversible as in the Euler case, but the directional derivatives of proposal moves in $\widehat{\Psi}_\varepsilon$ are fixed at the endpoints and hence constant within the proposal sequence. This allows for HMC transitions between the proposal sequence endpoints provided by AUHMC.

There are many potential ways of specifying the functional form of $M(\overline{\theta_r, \theta_{r+1}^*})$. We take a user-friendly approach with light computational burden and set

$$(3.9) \qquad\qquad M(\overline{\theta_r, \theta_{r+1}^*}) = \frac{1}{2}\left[F(\theta_r) + F(\theta_{r+1}^*)\right]$$

where $F(\theta)$ is the Fisher information matrix evaluated at $\theta$. Using $F(\theta)$ to convey information about the curvature of $\ln \pi(\theta)$ at $\theta$ was suggested in Girolami and Calderhead (2011).

Intuitively, AUHMC amounts to running the HMC between $\theta_r$ and $\theta_{r+1}^*$, using the information about the curvature of $\ln \pi(\theta)$ at the end points $\theta_r$ and $\theta_{r+1}^*$ in a symmetric way which preserves detailed balance of the resulting Markov chain. In contrast, the local curvature information is not utilized in HMC where the mass matrix $M$ is exogenously set, often to the identity matrix. Consequently, the HMC results as a special case of AUHMC for a globally constant matrix $M$ over the entire parameter space of $(\theta, \gamma)$. As another special case when $\ln \pi(\theta)$ has a globally constant curvature with respect to $\theta$, such as when $\theta = \mu$ for data $y \sim \mathcal{N}(\mu, I)$, the AUHMC produces draws equivalent to the HMC. In general, however, when the curvature of $\ln \pi(\theta)$ changes as a function of $\theta$, such as in $\theta = (\mu, \Sigma)$ for data $y \sim \mathcal{N}(\mu, \Sigma)$, AUHMC exploits the shape of $\ln \pi(\theta)$ by locally adapting the proposal dynamics to the curvature of $\ln \pi(\theta)$.

A key feature of AUHMC, in line with other HMC-based schemes, is that it simplifies the acceptance probability (2.2) to the Metropolis form containing only the ratio of the joint densities of $(\theta, \gamma)$. This feature provides for a user-friendly implementation of the algorithm.

## 4. Simulated Examples

In this Section we assess the performance of AUHMC on two stylized illustrative examples. Girolami and Calderhead (2011) provide an excellent exposition of a series of problems that highlight the superior performance of RMHMC relative to other related samplers. Hence, to establish the performance merit of AUHMC we believe that it is sufficient to take RMHMC as the benchmark of comparison. We first examine sampling of the parameters in multivariate Normal density in Example 1, and then sampling of the parameters in a univariate GARCH(1,1) model in Example 2. In Example 1 we fix the

sample size and increase the parameter dimensionality; in Example 2 we fix the dimensionality and increase the sample size. This setup is intended to uncover any potential trends in the performance comparison.

We compare the relative efficiency of AUHMC and RMHMC by using the same approach as Girolami and Calderhead (2011) and Holmes and Held (2006) in making their comparisons. For each example and method, we calculate the effective sample size (ESS), which is is the number of effectively independent draws from the posterior distribution that the Markov chain is equivalent to. The ESS thus serves as an estimate of the number of independent samples needed to obtain a parameter estimate with the same precision as the MCMC estimate considered based on a given number of dependent samples. The nominal ESS is calculated as $ESS^* = R \left[ 1 + 2 \sum_j \gamma(j) \right]^{-1}$ where $R$ is the number of posterior samples, and $\gamma(j)$ is the monotone sample autocorrelation (Geyer, 1992). The nominal ESS is then normalized for CPU run time required to obtain the given Markov chain of posterior draws, yielding $ESS = 100 \times ESS^*/S$ where $S$ is the number of seconds of CPU run time. The MCMC chains were obtained on a 2.8 GHz unix workstation with the Intel fortran 95 compiler. For obtaining $ESS^*$ from the MCMC output chains we used the R package coda. All results reported are the averages of 10 different runs.

The results for the examples considered here are given in Tables 1 and 2 and Figures 1 and 2 below. We report the mean, standard deviation, minimum, and maximum ESS for the sampled parameter vector for each simulation setup. We also report the nominal (unnormalized) numbers along with the CPU run time as the ESS inputs. In the tables, Ratio denotes the ratio of AUHMC to RMHMC of the respective statistics. Values greater than 1 indicate better performance of AUHMC. Figures 1 and 2 plot the relative efficiency gain of AUHMC over RMHMC, calculated as the ESS means ratio for the two methods. Figure 1 shows the AUHMC relative efficiency gain for increasing dimensionality and fixed sample size in Example 1, and Figure 1 for fixed dimensionality and increasing sample size in Example 2. In each Figure, the horizontal dotted line at $y$-value 1 marks theoretical equivalence of both methods, while the region above 1 represents efficiency gains of AUHMC.

### 4.1. **Example 1: Joint Sampling of Parameters of a Multivariate Normal Density**

Let $\mathbf{y}_t \sim \mathcal{N}(\mathbf{y}|\mu, \Sigma)$ for $t = 1, \ldots, T$ with

$$\ln \pi(\theta) = -\frac{Td}{2} \ln(2\pi) - \frac{T}{2} \ln |\Sigma| - \frac{1}{2} \sum_{t=1}^{T} (\mathbf{y}_t - \mu)' \Sigma^{-1} (\mathbf{y}_t - \mu)$$

and $\theta \equiv (\mu', vech(\Sigma)')'$. Naturally, a convenient factorization of this problem is readily available, but this stylized example is meant to serve for joint sampling comparison purposes on a familiar and analytically tractable case. In general applications, a conditional factorization of the joint density $\ln \pi(\theta)$ may not be available or practical to implement (this is for instance the case of the BEKK GARCH model analyzed in the next Section). In the simulation study of Example 1, we vary $\dim(\mathbf{y})$ from 3 to 6, which corresponds to the parameter dimensionality $\dim(\theta)$ varying from 9 to 27. The

true parameter values were set to $\mu_0 = 0$, and $\Sigma$ to equal the covariance matrix of a first-order autoregressive process with correlation 0.5. Our prior restricts $\Sigma$ to be positive definite. Each chain was initialized at the true parameter values, with $L = 10$ leapfrog steps, and the stepsize $\varepsilon$ tuned to achieve acceptance rates close to 0.8. The posterior samples were obtained from 2,000 parameter draws with a 1,000 burnin section. The ESS statistics are reported in Table 1 and Figure 1.

## 4.2. Example 2: Joint Sampling of GARCH (1,1) Parameters

Let $y_t \sim \mathcal{N}(y|0, \sigma_t^2)$ with $\sigma_t^2 = \gamma + \alpha y_{t-1}^2 + \beta \sigma_{t-1}^2$ for $t = 1, \ldots, T$ and $\theta \equiv (\gamma, \alpha, \beta)$ where

$$
\ln \pi(\theta) = -\frac{T}{2} \ln(2\pi) - \frac{1}{2} \sum_{t=1}^{T} \ln(\sigma_t^2(\theta)) - \frac{1}{2} \sum_{t=1}^{T} y_t \sigma_t^{-2}(\theta)
$$

$$
\frac{\partial}{\partial \theta} \ln \pi(\theta) = -\frac{1}{2} \sum_{t=1}^{T} \frac{1}{\sigma_t^2} \frac{\partial \sigma_t^2(\theta)}{\partial \theta} + \frac{1}{2} \sum_{t=1}^{T} e_t^2 \sigma_t^{-4} \frac{\partial \sigma_t^2(\theta)}{\partial \theta}
$$

$$
\frac{\partial \sigma_t^2(\theta)}{\partial \gamma} = 1 + \beta \frac{\partial \sigma_{t-1}^2(\theta)}{\partial \gamma}
$$

$$
\frac{\partial \sigma_t^2(\theta)}{\partial \alpha} = e_{t-1}^2 + \beta \frac{\partial \sigma_{t-1}^2(\theta)}{\partial \alpha}
$$

$$
\frac{\partial \sigma_t^2(\theta)}{\partial \beta} = \sigma_{t-1}^2 + \beta \frac{\partial \sigma_{t-1}^2(\theta)}{\partial \beta}
$$

and $F(\theta)$ is consistently estimated using the average of the outer products of the scores. In this simulation study we vary the sample size $T$ from 200 to 600. The dimensionality of the parameter space of $\theta$ is kept constant at 3. The true parameter values were set to $\gamma_0 = 0.1$, $\alpha_0 = 0.05$ and $\beta_0 = 0.9$. Each chain was initialized at the true parameter values, with $L = 100$ leapfrog steps, and the stepsize $\varepsilon$ tuned to achieve acceptance rates close to 0.8. The posterior samples were obtained from 10,000 parameter draws with a 5,000 burnin section. Due to a more complex structure of the likelihood compared to Example 1, a larger number of smaller steps and longer MC run were necessary to obtain good mixing properties in this case. The ESS statistics are reported in Table 2 and Figure 2.

In summary, the improvement of AUHMC over RMHMC is substantial, with up to 50-fold efficiency gain in Example 1 and up to 11-fold efficiency gain in Example 2. In both examples, the improvement keeps increasing with increasing dimensionality and sample size, indicating sustained efficiency gain of AUHMC for more complex and sizeable problems. Both increasing the dimensionality and sample size add additional heavy computational load to the RMHMC in its fixed point iterations that AUHMC avoids. These examples highlight the benefits of AUHMC on interesting cases in order to motivate its use in applications.

Table 1: Simulation Results for Example 1

| Variable dimension | | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|
| Parameter dimension | | 9 | 14 | 20 | 27 |
| CPU Time (s) | AUHMC | 5.282 | 17.7 | 117.673 | 198.396 |
| | RMHMC | 6.886 | 22.841 | 325.452 | 952.945 |
| | Ratio | 0.767 | 0.775 | 0.362 | 0.208 |
| Nominal ESS mean | AUHMC | 157.627 | 125.713 | 94.312 | 76.721 |
| | RMHMC | 38.525 | 19.959 | 11.992 | 7.922 |
| ESS mean | AUHMC | 2991.247 | 712.301 | 80.265 | 38.731 |
| | RMHMC | 559.449 | 87.383 | 3.685 | 0.831 |
| | Ratio | 5.559 | 8.306 | 22.641 | 52.135 |
| Nominal ESS s.d. | AUHMC | 91.535 | 57.877 | 42.776 | 32.461 |
| | RMHMC | 8.271 | 5.629 | 4.011 | 2.884 |
| ESS s.d. | AUHMC | 1734.502 | 328.755 | 36.389 | 16.37 |
| | RMHMC | 120.066 | 24.652 | 1.232 | 0.303 |
| Nominal ESS min | AUHMC | 61.1 | 59.524 | 40.359 | 32.803 |
| | RMHMC | 30.342 | 15.088 | 8.22 | 5.222 |
| ESS min | AUHMC | 1160.743 | 337.019 | 34.399 | 16.622 |
| | RMHMC | 440.663 | 66.052 | 2.526 | 0.548 |
| | Ratio | 2.81 | 5.191 | 14.636 | 35.298 |
| Nominal ESS max | AUHMC | 294.058 | 233.335 | 183.867 | 148.242 |
| | RMHMC | 46.166 | 25.738 | 15.903 | 10.696 |
| ESS max | AUHMC | 5574.415 | 1325.778 | 156.48 | 74.841 |
| | RMHMC | 670.379 | 112.69 | 4.887 | 1.123 |
| | Ratio | 8.735 | 12.073 | 33.814 | 76.204 |
| Acceptance rate | AUHMC | 0.821 | 0.824 | 0.812 | 0.795 |
| | RMHMC | 0.798 | 0.785 | 0.816 | 0.836 |
| $\varepsilon$ | AUHMC | 0.11 | 0.095 | 0.085 | 0.078 |
| | RMHMC | 0.04 | 0.03 | 0.02 | 0.015 |

Figure 1: AUHMC Efficiency Gain in Example 1

Table 2: Simulation Results for Example 2

| Sample size $T$<br>Parameter dimension | | 200<br>3 | 300<br>3 | 400<br>3 | 500<br>3 | 600<br>3 |
|---|---|---|---|---|---|---|
| CPU Time (s) | AUHMC<br>RMHMC<br>Ratio | 43.722<br>131.152<br>0.333 | 45.981<br>195.712<br>0.235 | 63.116<br>260.692<br>0.243 | 77.09<br>321.97<br>0.24 | 89.092<br>381.456<br>0.234 |
| Nominal ESS mean | AUHMC<br>RMHMC | 14.724<br>20.484 | 20.939<br>20.356 | 20.199<br>19.251 | 21.656<br>16.831 | 31.02<br>17.511 |
| ESS mean | AUHMC<br>RMHMC<br>Ratio | 35.796<br>15.619<br>4.175 | 44.996<br>10.429<br>4.396 | 32.184<br>7.407<br>5.215 | 28.565<br>5.246<br>8.507 | 34.425<br>4.593<br>11.29 |
| Nominal ESS s.d. | AUHMC<br>RMHMC | 17.278<br>28.009 | 28.476<br>26.981 | 27.348<br>26.989 | 30.858<br>22.046 | 46.333<br>22.892 |
| ESS s.d. | AUHMC<br>RMHMC | 42.628<br>21.353 | 61.038<br>13.824 | 43.607<br>10.387 | 40.865<br>6.892 | 51.269<br>6.005 |
| Nominal ESS min | AUHMC<br>RMHMC | 3.596<br>3.467 | 2.993<br>3.253 | 3.029<br>2.972 | 2.913<br>3.039 | 3.429<br>2.993 |
| ESS min | AUHMC<br>RMHMC<br>Ratio | 8.446<br>2.645<br>3.442 | 6.486<br>1.665<br>3.943 | 4.781<br>1.142<br>4.345 | 3.769<br>0.937<br>4.705 | 3.884<br>0.785<br>5.682 |
| Nominal ESS max | AUHMC<br>RMHMC | 34.579<br>52.758 | 53.706<br>51.403 | 51.676<br>50.381 | 57.248<br>41.967 | 84.5<br>43.565 |
| ESS max | AUHMC<br>RMHMC<br>Ratio | 84.799<br>40.223<br>5.13 | 115.218<br>26.337<br>4.416 | 82.369<br>19.387<br>5.486 | 75.703<br>13.105<br>9.862 | 93.602<br>11.428<br>15.842 |
| Acceptance rate | AUHMC<br>RMHMC | 0.726<br>0.815 | 0.874<br>0.873 | 0.824<br>0.806 | 0.819<br>0.787 | 0.851<br>0.751 |
| $\varepsilon$ | AUHMC<br>RMHMC | 3e-05<br>2e-05 | 1e-05<br>1e-05 | 1e-05<br>1.1e-05 | 8e-06<br>1e-05 | 6e-06<br>1e-05 |

Figure 2: AUHMC Efficiency Gain in Example 2

## 5. BEKK GARCH Application

Interest in modeling the volatility dynamics of time-series data continues to grow and be important in many areas of empirical economics and finance. Generally, the literature on multivariate asset return modeling has moved to using more parsimonious models such as Engle (2002), Engle, Shephard, and Sheppard (2009) and Ding and Engle (2001). These approaches put restrictions on the volatility dynamics and feature two-step estimation and approximations to the likelihood. This makes estimation and inference feasible for a larger class of assets. However, it is desirable to consider more flexible models such as the BEKK model of Engle and Kroner (1995) and to perform full likelihood based inference. The BEKK model is one of the most flexible GARCH models that maintain positive definite conditional covariances at the expense of a large number of parameters. Although inference of the model with 2 or 3 assets have appeared in the literature we are not aware of anything beyond this asset dimension. An important question is how much do we lose in terms of statistical fit in moving from a BEKK model to a restricted model with fewer parameters to estimate. The extension to HMC discussed above provides an approach that can deal with the larger dimensions in the parameter space and jointly estimate the BEKK model in one run and compare the model to restricted versions.

Let $r_t$ be a $N \times 1$ vector of asset returns with $t = 1, \ldots, T$ and denote the information set as $\mathcal{F}_{t-1} = \{r_1, \ldots, r_{t-1}\}$. We assume returns follow

$$(5.1) \qquad\qquad r_t | \mathcal{F}_{t-1} \quad \sim \quad NID(0, H_t)$$

$$(5.2) \qquad\qquad H_t \quad = \quad CC' + F' r_{t-1} r'_{t-1} F + G' H_{t-1} G.$$

$H_t$ is a positive definite $N \times N$ conditional covariance matrix of $r_t$ given information at time $t-1$, $C$ is a lower triangular matrix and $F$ and $G$ are $N \times N$ matrices. Since our main focus is on sampling a complex posterior with many parameters we maintain a Gaussian assumption and a zero intercept for simplicity.[1] The total number of parameters in this model is $N(N+1)/2 + 2N^2$.

In the following we focus on the full BEKK model in (5.1) but also consider some restricted versions. The first imposes $F$ and $G$ to be diagonal matrices which results in $N(N+1)/2 + 2N$ parameters. The second imposes diagonal matrices on all parameter matrices $C, F$ and $G$ and has $3N$ parameters.

The data is percent log-differences of foreign exchange spot rates for AUD/USD, GBP/USD, CAD/USD, EUR/USD, and JPY/USD from 2000/01/05 - 2006/10/11, (1700 observations). A time series plot of the five ($N = 5$) series is in Figure 3 and summary statistics are in Table 3. The sample mean for all series is close to 0 and excess kurtosis is fairly small. The sample correlations indicate all series tend to move together.

With $N = 5$ there are 65 model parameters in the full BEKK model while there are 25 and 15 parameters, respectively, in the two restricted models. To start the GARCH recursion $H_1$ is set to

---

[1]Although not estimated, we expect our method could be extended to other innovation distributions such as multivariate Student-t with little modification.

the sample covariance of the first 20 observations. The priors are set to independent N(0,100). For identification, the diagonal elements of $C$ and the first element of both $F$ and $G$ are restricted to be positive (Engle and Kroner, 1995). These restrictions are enforced by dropping any parameter draw that violates this. We utilize the analytical expressions for the gradient from Hafner and Herwartz (2008), and Fisher Information matrix given in Appendix D. Starting from a point of high posterior mass we collect a total of 30,000 posterior draws for inference, with 10,000 burnin section. These computations took on the order of 2 days.

Collecting the parameters in $\theta = (vech(C)', vech(F)', vech(G)')'$, Figure 4 displays the conditional log-posterior $\log p(\theta_i|\theta_{-i}, \mathcal{F}_T)$ where $\theta_{-i}$ is set to a high probability mass point. Some of the conditional densities are approximately quadratic while others display a more complicated structure. The relatively flat regions in the log-posterior will present challenges to maximizing this function or to obtaining a hessian estimate to compute standard errors in a classical approach.

Figure 5 displays the posterior mean of the conditional correlations for the full BEKK model and the two restricted versions. The BEKK model being the most flexible displays differences with the other models most notably the version that enforces diagonal matrices on $C, G, F$. That restriction implies unconditional correlations of 0 between assets and is at odds with the sample correlations in Table 3.

These differences in the models are confirmed by the marginal likelihoods reported in Table 4. The marginal likelihoods are estimated following Gelfend and Dey (1994) using a thin tailed truncated normal following Geweke (2005). The evidence is strongly against both of the restricted diagonal models. For example, the log-Bayes factor in favor of the full BEKK model is about 35 compared to the model with diagonal $F, G$.

In conclusion, our results support the use of the most flexible BEKK model and the AUHMC sampler provides a feasible method to sample from a highly complex posterior density effectively.

## 6. Conclusion

Hamiltonian Monte Carlo (HMC) uses Hamiltonian dynamics in constructing distant proposal draws in a sequence of steps, yielding relatively low correlation among draws and high acceptance probabilities at the same time. In this paper we propose a local adaptation of HMC, the Adaptively Updated Hamiltonian Monte Carlo (AUHMC), whereby the proposal sequence follows the local evolution of the parameter space. We provide a set of sufficient conditions on the (posterior) likelihood under which AUHMC yields a valid MCMC procedure satisfying detailed balance. Simulated examples show that the performance gain of AUHMC increases with increasing dimensionality or sample size. We apply AUHMC to a high-dimensional BEKK GARCH model in 56 parameter dimensions, which substantially exceeds the dimensionality utilized in previous work. Model comparison via marginal likelihood further reveals that the full BEKK model is preferable to its restricted versions with constraints placed on various covariance components, motivating the full high-dimensional implementation of the model.

# 7. Appendix A: Hamiltonian Monte Carlo

In this Section we provide the stochastic background for HMC. This synthesis is based on previously published material, but unlike the bulk of literature presenting HMC in terms of the physical laws of motion based on preservation of total energy in the phase-space, we take a fully stochastic perspective familiar to the applied Bayesian econometrician.[2] The HMC principle is thus presented in terms of the joint density over the augmented parameter space leading to a Metropolis acceptance probability update. We hope that our synthesis of the probabilistic perspective on HMC will provide useful insights for practitioners who wish to further explore the HMC principles.

## 7.1. HMC Principle

Consider a vector of parameters of interest $\theta \in \mathbb{R}^d$ distributed according to the posterior density $\pi(\theta)$. Let $\gamma \in \mathbb{R}^d$ denote a vector of auxiliary parameters with $\gamma \sim \Phi(\gamma; 0, M)$ where $\Phi$ denotes the Gaussian distribution with mean vector 0 and covariance matrix $M$, independent of $\theta$. Denote the joint density of $(\theta, \gamma)$ by $\pi(\theta, \gamma)$. Then the negative of the logarithm of the joint density of $(\theta, \gamma)$ is given by the Hamiltonian equation[3]

$$(7.1) \qquad H(\theta, \gamma) = -\ln \pi(\theta) + \frac{1}{2} \ln\left((2\pi)^d |M|\right) + \frac{1}{2}\gamma' M^{-1}\gamma$$

Hamiltonian Monte Carlo (HMC) is formulated in the following three steps that we will describe in detail further below:

(1) Draw an initial auxiliary parameter vector $\gamma_r^0 \sim \Phi(\gamma; 0, M)$;
(2) Transition from $(\theta_r, \gamma_r)$ to $(\theta_r^L, \gamma_r^L) = (\theta_{r+1}^*, \gamma_{r+1}^*)$ according to the Hamiltonian dynamics;
(3) Accept $(\theta_{r+1}^*, \gamma_{r+1}^*)$ with probability $\alpha(\theta_r, \gamma_r; \theta_{r+1}^*, \gamma_{r+1}^*)$, otherwise keep $(\theta_r, \gamma_r)$ as the next MC draw.

*Step 1* provides a stochastic initialization of the system akin to a RW draw. This step is necessary in order to make the resulting Markov chain $\{(\theta_r, \gamma_r)\}_{r=1}^R$ irreducible and aperiodic (Ishwaran, 1999). In contrast to RW, this so-called refreshment move is performed on the auxiliary variable $\gamma$ as opposed to the original parameter of interest $\theta$, setting $\theta_r^0 = \theta_r$. In terms of the HMC sampling algorithm, the initial refreshment draw of $\gamma_r^0$ forms a Gibbs step on the parameter space of $(\theta, \gamma)$ accepted with probability 1. Since it only applies to $\gamma$, it will leave the target joint distribution of $(\theta, \gamma)$ invariant and subsequent steps can be performed conditional on $\gamma_r^0$ (Neal, 2010).

*Step 2* constructs a sequence $\{\theta_r^k, \gamma_r^k\}_{k=1}^L$ according to the Hamiltonian dynamics starting from the current state $(\theta_r^0, \gamma_r^0)$ and setting the last member of the sequence as the HMC new state proposal $(\theta_{r+1}^*, \gamma_{r+1}^*) = (\theta_r^L, \gamma_r^L)$. The role of the Hamiltonian dynamics is to ensure that the M-H acceptance probability (2.2) for $(\theta_{r+1}^*, \gamma_{r+1}^*)$ is kept close to 1. As will become clear shortly, this corresponds to maintaining the difference $-H(\theta_{r+1}^*, \gamma_{r+1}^*) + H(\theta_r^0, \gamma_r^0)$ close to zero throughout the sequence

---

[2]There are notable exceptions, such as Girolami and Calderhead (2011) who also take the statistical perspective, but their paper focuses on RMHMC while here we elaborate on the statistical background to HMC.

[3]In the physics literature, $\theta$ denotes the position (or state) variable and $-\ln \pi(\theta)$ describes its potential energy, while $\gamma$ is the momentum variable with kinetic energy $\gamma' M^{-1}\gamma/2$, yielding the total energy $H(\theta, \gamma)$ of the system, up to a constant of proportionality. $M$ is a constant, symmetric, positive-definite "mass" matrix which is often set as a scalar multiple of the identity matrix.

$\{\theta_r^k, \gamma_r^k\}_{k=1}^L$. This property of the transition from $(\theta_r, \gamma_r)$ to $(\theta_{r+1}^*, \gamma_{r+1}^*)$ can be achieved by conceptualizing $\theta$ and $\gamma$ as functions of continuous time $t$ and specifying their evolution using the Hamiltonian dynamics equations[4]

$$(7.2) \qquad \frac{d\theta_i}{dt} = \frac{\partial H(\theta, \gamma)}{\partial \gamma_i} = \left[ M^{-1}\gamma \right]_i$$

$$(7.3) \qquad \frac{d\gamma_i}{dt} = -\frac{\partial H(\theta, \gamma)}{\partial \theta_i} = \nabla_{\theta_i} \ln \pi(\theta)$$

for $i = 1, \ldots, d$. For any discrete time interval of duration $s$, (7.2)–(7.3) define a mapping $T_s$ from the state of the system at time $t$ to the state at time $t + s$. For practical applications of interest these differential equations (7.2)–(7.3) in general cannot be solved analytically and instead numerical methods are required. The Stormer-Verlet (or leapfrog) numerical integrator (Leimkuhler and Reich, 2004) is one such popular method, discretizing the Hamiltonian dynamics as

$$(7.4) \qquad \gamma(t + \varepsilon/2) = \gamma(t) + (\varepsilon/2)\nabla_\theta \ln \pi(\theta(t))$$

$$(7.5) \qquad \theta(t + \varepsilon) = \theta(t) + \varepsilon M^{-1}\gamma(t + \varepsilon/2)$$

$$(7.6) \qquad \gamma(t + \varepsilon) = \gamma(t + \varepsilon/2) + (\varepsilon/2)\nabla_\theta \ln \pi(\theta(t + \varepsilon))$$

for some small $\varepsilon \in \mathbb{R}$. From this perspective, $\gamma$ plays the role of an auxiliary variable that parametrizes (a functional of) $\pi(\theta, \cdot)$ providing it with an additional degree of flexibility to maintain the acceptance probability close to one for every $k$. Even though $\ln \pi(\theta_r^k)$ can deviate substantially from $\ln \pi(\theta_r^0)$, resulting in favorable mixing for $\theta$, the additional terms in $\gamma$ in (7.1) compensate for this deviation maintaining the overall level of $H(\theta_r^k, \gamma_r^k)$ close to constant over $k = 1, \ldots, L$ when used in accordance with (7.4)–(7.6), since $\frac{\partial H(\theta, \gamma)}{\partial \gamma_i}$ and $\frac{\partial H(\theta, \gamma)}{\partial \theta_i}$ enter with the opposite signs in (7.2)–(7.3). In contrast, without the additional parametrization with $\gamma$, if only $\ln \pi(\theta_r^k)$ were to be used in the proposal mechanism as is the case in RW style samplers, the M-H acceptance probability would often drop to zero relatively quickly.

*Step 3* applies a Metropolis correction to the proposal $(\theta_{r+1}^*, \gamma_{r+1}^*)$. In continuous time, or for $\varepsilon \to 0$, (7.2)–(7.3) would keep $-H(\theta_{r+1}^*, \gamma_{r+1}^*) + H(\theta_r, \gamma_r) = 0$ exactly resulting in $\alpha(\theta_r, \theta_{r+1}^*) = 1$ but for discrete $\varepsilon > 0$, in general, $-H(\theta^*, \gamma^*) + H(\theta, \gamma) \neq 0$ necessitating the Metropolis step. A key feature of HMC is that the generic M-H acceptance probability (2.2) can be expressed in a simple tractable form using only the posterior density $\pi(\theta)$ and the auxiliary parameter Gaussian density $\phi(\gamma; 0, M)$. The transition from $(\theta_r^0, \gamma_r^0)$ to $(\theta_r^L, \gamma_r^L)$ via the proposal sequence $\{\theta_r^k, \gamma_r^k\}_{k=1}^L$ taken according to the discretized Hamiltonian dynamics (7.4)–(7.6) is fully deterministic proposal, placing a Dirac delta probability mass $\delta(\theta_r^k, \gamma_r^k) = 1$ on each $(\theta_r^k, \gamma_r^k)$ conditional on $(\theta_r^0, \gamma_r^0)$. The system (7.4)–(7.6) is time reversible and symmetric in $(\theta, \gamma)$, which implies that the forward and reverse transition probabilities $q(\theta_r^L, \gamma_r^L | \theta_r^0, \gamma_r^0)$ and $q(\theta_r^0, \gamma_r^0 | \theta_r^L, \gamma_r^L)$ are equal: this simplifies the Metropolis-Hastings acceptance ratio in (2.2) to the Metropolis form $\pi(\theta_{r+1}^*, \gamma_{r+1}^*)/\pi(\theta_r^0, \gamma_r^0)$. From the definition of the Hamiltonian $H(\theta, \gamma)$ in (7.1) as the negative of the log-joint densities, the joint density of $(\theta, \pi)$ is given by

$$(7.7) \qquad \pi(\theta, \gamma) = \exp\left[-H(\theta, \gamma)\right] = \pi(\theta) \left( (2\pi)^d |M| \right)^{-1/2} \exp\left( -\frac{1}{2}\gamma' M^{-1}\gamma \right)$$

Hence, the Metropolis acceptance probability takes the form

---

[4]In the physics literature, the Hamiltonian dynamics describe the evolution of $(\theta, \gamma)$ that keeps the total energy $H(\theta, \gamma)$ constant.

$$
\begin{aligned}
\alpha(\theta_r, \gamma_r; \theta^*_{r+1}, \gamma^*_{r+1}) \; &= \; \min\left[\frac{\pi(\theta^*_{r+1}, \gamma^*_{r+1})}{\pi(\theta^0_r, \gamma^0_r)}, 1\right] \\
&= \; \min\left[\exp\left(-H(\theta^*_{r+1}, \gamma^*_{r+1}) + H(\theta^0_r, \gamma^0_r)\right), 1\right] \\
&= \; \min\left[\exp\left(\ln\pi(\theta^*_{r+1}) - \ln\pi(\theta^0_r) + \ln\phi(\gamma^*_{r+1}; 0, M) - \ln\phi(\gamma^0_r; 0, M)\right), 1\right]
\end{aligned}
$$

The expression for $\alpha(\theta_r, \gamma_r; \theta^*_{r+1}, \gamma^*_{r+1})$ shows, as noted above, that the HMC acceptance probability is given in terms of the difference of the Hamiltonian equations $H(\theta^0_r, \gamma^0_r) - H(\theta^*_{r+1}, \gamma^*_{r+1})$. The closer can we keep this difference to zero, the closer the acceptance probability is to one. A key feature of the Hamiltonian dynamics (7.2)–(7.3) in Step 2 is that they maintain $H(\theta, \gamma)$ constant over the parameter space in continuous time conditional on $H(\theta^0_r, \gamma^0_r)$ obtained in Step 1, while their discretization (7.4)–(7.6) closely approximates this property for discrete time steps $\varepsilon > 0$ with a global error of order $\varepsilon^2$ corrected by the Metropolis update in Step 3.

## 8. Appendix B: The AUHMC Algorithm

Initialize current $\theta$

**for** $r = 1$ **to** $R$ {

    initialize $\theta^0_r = \theta_r$, $j = 0$

    **(j loop) do while** $\left(\left(\|\theta^{L,j}_r - \theta^{L,j-1}_r\| > \delta_1\right) \textbf{ or } \left(\|\gamma^{L,j}_r - \gamma^{L,j-1}_r\| > \delta_2\right)\right)$ {

        draw $\gamma^{0,j}_r \sim q(\gamma^{0,j}_r | \theta_r) = N(0, M(\overline{\theta_r, \theta^0_r}))$ for $j = 0$ and $N(0, M(\overline{\theta_r, \theta^{L,j}_r}))$ for $j > 0$

        $j = j + 1$

        **(k loop) for** $k = 1$ **to** $L$ {

            $\gamma^{k+1/2,j}_r = \gamma^{k,j}_r + \frac{\varepsilon}{2}\nabla_\theta \ln\pi(\theta^{k,j}_r)$

            $\theta^{k+1,j}_r = \theta^{k,j}_r + \varepsilon\left[M(\overline{\theta_r, \theta^*_{r+1}})^{-1}\gamma^{k+1/2,j}_r\right]$

            $\gamma^{k+1,j}_r = \gamma^{k+1/2,j}_r + \frac{\varepsilon}{2}\nabla_\theta \ln\pi(\theta^{k+1,j}_r)$

        }

        $M(\overline{\theta_r, \theta^*_{r+1}}) = \frac{1}{2}\left[F(\theta_r) + F(\theta^{L,j}_r)\right]$

    }

    $\alpha^* = \frac{\pi(\theta^*_{r+1})q(\widetilde{\gamma}^0_r | \theta^*_{r+1})}{\pi(\theta_r)q(\gamma^0_r | \theta_r)}$

    draw $u \sim U[0,1]$

    **if** $(\alpha^* < u)$ **then** $\{\theta_{r+1} = \theta^{L,j}_r\}$ **else** $\{\theta_{r+1} = \theta_r\}$

}

## 9. Appendix C: Proof of Lemma 1 and Theorem 1

### 9.1. Proof of Lemma 1

The AUHMC mapping is a special case of an implicit Runge-Kutta method (Leimkuhler and Reich, 2004, p. 150-151). Hence, under our Assumptions 2 and 3, the proof of existence of a unique solution is given by Theorem 7.2 of Hairer, Nørsett, and Wanner (1993, p. 206). Specifically, there exists a unique solution to the mapping $T_k$ defined by (3.5)-(3.7) which can be obtained by iteration resulting in the repeated use of the triangle inequality that results from the Lipschitz condition satisfying a contraction mapping property.

### 9.2. Proof of Theorem 1

Recall that AUHMC constructs a distant proposal sequence $\{\theta^k, \gamma^k\}_{k=1}^{L}$ in a sequence of $k = 1, \ldots, L$ steps. For a given $k$ (omitting the subscripts $r$ denoting the MCMC steps), define the mapping $\Psi_\varepsilon^k$ of $(\theta^k, \gamma^k)$ into $(\theta^{k+1}, \gamma^{k+1})$ as:

$$\overline{\theta}^1 = \theta^k$$

$$\overline{\theta}^i = \theta^k + \varepsilon \widehat{a}_{i1} \nabla_\gamma H(\overline{\theta}^1, \overline{\gamma}^1) + \varepsilon \sum_{j=2}^{i-1} \widehat{a}_{ij} \nabla_\gamma H(\overline{\theta}^j, \overline{\gamma}^j), \quad \begin{matrix} \widehat{a}_{ij} = -2, \ i = 2, \ldots, k+1, \\ j = 1, \ldots, i-1 \end{matrix}$$

$$\overline{\theta}^i = \theta^k + \varepsilon \widehat{a}_{i1} \nabla_\gamma H(\overline{\theta}^1, \overline{\gamma}^1) + \varepsilon \sum_{j=k+2}^{i-1} \widehat{a}_{ij} \nabla_\gamma H(\overline{\theta}^j, \overline{\gamma}^j), \quad \begin{matrix} \widehat{a}_{ij} = 2, \ i = k+2, \ldots, L+1, \\ j = 1, \ldots, i-1 \end{matrix}$$

$$\overline{\gamma}^1 = \gamma^k + \varepsilon \widetilde{a}_{11} \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1), \ \widetilde{a}_{11} = 1$$

$$\overline{\gamma}^i = \gamma^k + \varepsilon \widetilde{a}_{i1} \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1) + \varepsilon \sum_{j=2}^{i-1} \widetilde{a}_{ij} \nabla_\theta H(\overline{\theta}^j, \overline{\gamma}^j), \quad \begin{matrix} \widetilde{a}_{i1} = 1, \ \widetilde{a}_{ij} = 2, \ i = 2, \ldots, k, \\ j = 2, \ldots, i-1 \end{matrix}$$

$$\overline{\gamma}^{k+1} = \gamma^k + \varepsilon \widetilde{a}_{k+1,1} \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1), \ \widetilde{a}_{k+1,1} = -1$$

$$\overline{\gamma}^i = \gamma^k + \varepsilon \widetilde{a}_{i1} \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1) + \varepsilon \sum_{j=k+2}^{i-1} \widetilde{a}_{ij} \nabla_\theta H(\overline{\theta}^j, \overline{\gamma}^j), \quad \begin{matrix} \widetilde{a}_{i1} = -1, \ \widetilde{a}_{ij} = -2, \ i = k+2, \ldots, L, \\ j = 2, \ldots, i-1 \end{matrix}$$

$$\overline{\gamma}^{L+1} = \gamma^k + \varepsilon \widetilde{a}_{L+1,1} \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1), \ \widetilde{a}_{L+1,1} = -1$$

$$\theta^{k+1} = \theta^k + \varepsilon \left[ \widehat{b}_{k+1} \nabla_\gamma H(\overline{\theta}^{k+1}, \overline{\gamma}^{k+1}) + \widehat{b}_{L+1} \nabla_\gamma H(\overline{\theta}^{L+1}, \overline{\gamma}^{L+1}) \right], \ \widehat{b}_{k+1} = 1/2, \ \widehat{b}_{L+1} = 1/2$$

$$\gamma^{k+1} = \gamma^k + \varepsilon \widetilde{b}_1 \nabla_\theta H(\overline{\theta}^1, \overline{\gamma}^1), \ \widetilde{b}_1 = -1$$

The coefficient notation for $\widehat{a}_{ij}, \widetilde{a}_{ij}, \widehat{b}_1, \widetilde{b}_1$ corresponds to the general setup of an implicit partitioned Runge-Kutta scheme of Leimkuhler and Reich (2004, p. 150-151). Here, all $\widehat{a}_{ij}, \widetilde{a}_{ij}, \widehat{b}_1, \widetilde{b}_1$ are equal to zero unless stated otherwise. Moreover, if in the summation sign the upper index is smaller than the lower index, then the corresponding coefficient $\widehat{a}_{ij}$ or $\widetilde{a}_{ij}$ is equal to zero. The Hamiltonian $H(\overline{\theta}^k, \overline{\gamma}^k)$ for each $k$ is given by

$$H(\overline{\theta}^k, \overline{\gamma}^k) = -\ln \pi(\overline{\theta}^k) + \frac{d}{2} \ln(2\pi) + \frac{1}{2} \ln |M_k| + \frac{1}{2} \overline{\gamma}^{k\prime} M_k^{-1} \overline{\gamma}^k$$

with

$$M_k = M(\overline{\theta^{k+1}, \overline{\theta}^{L+1}})$$

where the right-hand side is defined in (3.5), and $\overline{\theta}^{k+1}$ and $\overline{\theta}^{L+1}$ are implicitly determined in $\Psi_\varepsilon^k$.

We will next state the definitions of an adjoint mapping (Leimkuhler and Reich, 2004, p. 84).

**Definition 1.** *The mapping $\Psi_\varepsilon^{*k}$ defined by $\left[\Psi_{-\varepsilon}^k\right]^{-1}$ is called the adjoint mapping of $\Psi_\varepsilon^k$. Equivalently, given $\Psi_\varepsilon^k$, its adjoint is defined by*

$$
\begin{aligned}
(\theta^k, \gamma^k) &= \Psi_{-\varepsilon}^k(\theta^{k+1}, \gamma^{k+1}) \\
\Psi_\varepsilon^{*k}(\theta^k, \gamma^k) &= (\theta^{k+1}, \gamma^{k+1})
\end{aligned}
$$

Given $\Psi_\varepsilon^k$ as defined above, its adjoint $\Psi_\varepsilon^{*k}$ takes the form

$$
\begin{aligned}
\overline{\theta}^1 &= \theta^{k+1} \\
\overline{\theta}^i &= \theta^{k+1} + \varepsilon\widehat{a}_{i1}\nabla_\gamma H(\overline{\theta}^1,\overline{\gamma}^1) + \varepsilon\sum_{j=2}^{i-1}\widehat{a}_{ij}\nabla_\gamma H(\overline{\theta}^j,\overline{\gamma}^j), && \widehat{a}_{ij}=2,\ i=2,\ldots,k+1, \\
&&& j=1,\ldots,i-1 \\
\overline{\theta}^i &= \theta^{k+1} + \varepsilon\widehat{a}_{i1}\nabla_\gamma H(\overline{\theta}^1,\overline{\gamma}^1) + \varepsilon\sum_{j=k+2}^{i-1}\widehat{a}_{ij}\nabla_\gamma H(\overline{\theta}^j,\overline{\gamma}^j), && \widehat{a}_{ij}=-2,\ i=k+2,\ldots,L+1, \\
&&& j=1,\ldots,i-1 \\
\overline{\gamma}^1 &= \gamma^{k+1} + \varepsilon\widetilde{a}_{11}\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1),\ \widetilde{a}_{11}=1 \\
\overline{\gamma}^i &= \gamma^{k+1} + \varepsilon\widetilde{a}_{i1}\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1) + \varepsilon\sum_{j=2}^{i-1}\widetilde{a}_{ij}\nabla_\theta H(\overline{\theta}^j,\overline{\gamma}^j), && \widetilde{a}_{i1}=-1,\ \widetilde{a}_{ij}=-2,\ i=2,\ldots,k, \\
&&& j=2,\ldots,i-1 \\
\overline{\gamma}^{k+1} &= \gamma^{k+1} + \varepsilon\widetilde{a}_{k+1,1}\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1),\ \widetilde{a}_{k+1,1}=1 \\
\overline{\gamma}^i &= \gamma^{k+1} + \varepsilon\widetilde{a}_{i1}\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1) + \varepsilon\sum_{j=k+2}^{i-1}\widetilde{a}_{ij}\nabla_\theta H(\overline{\theta}^j,\overline{\gamma}^j), && \widetilde{a}_{i1}=1,\ \widetilde{a}_{ij}=2,\ i=k+2,\ldots,L, \\
&&& j=2,\ldots,i-1 \\
\overline{\gamma}^{L+1} &= \gamma^{k+1} + \varepsilon\widetilde{a}_{L+1,1}\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1),\ \widetilde{a}_{L+1,1}=1 \\
\theta^{k+1} &= \theta^k + \varepsilon\left[\widehat{b}_{k+1}\nabla_\gamma H(\overline{\theta}^{k+1},\overline{\gamma}^{k+1}) + \widehat{b}_{L+1}\nabla_\gamma H(\overline{\theta}^{L+1},\overline{\gamma}^{L+1})\right],\ \widehat{b}_{k+1}=-1/2,\ \widehat{b}_{L+1}=-1/2 \\
\gamma^{k+1} &= \gamma^k + \varepsilon\widetilde{b}_1\nabla_\theta H(\overline{\theta}^1,\overline{\gamma}^1),\ \widetilde{b}_1=1
\end{aligned}
$$

We next proceed to symmetric compositions of mappings with their adjoints.

**Definition 2.** *A mapping $\widehat{\Psi}_\varepsilon$ is called symmetric if $\widehat{\Psi}_\varepsilon = \widehat{\Psi}_\varepsilon^*$, i.e. $\widehat{\Psi}_{-\varepsilon} = \widehat{\Psi}_\varepsilon^{-1}$.*

The symmetry of $\widehat{\Psi}_\varepsilon^k$ then implies its time-reversibility (Leimkuhler and Reich, 2004, p. 87). Knowing a mapping $\Psi_\varepsilon^k$ and its adjoint $\Psi_\varepsilon^{*k}$, a symmetric mapping $\widehat{\Psi}_\varepsilon^k$ is obtained by composition (concatenation) of the two methods

(9.1) $$\widehat{\Psi}_\varepsilon^k \equiv \Psi_{\varepsilon/2}^{k,*} \circ \Psi_{\varepsilon/2}^k$$

even if neither $\Psi_{\varepsilon/2}^k$ nor $\Psi_{\varepsilon/2}^{k,*}$ are symmetric individually (Leimkuhler and Reich, 2004, p. 84). The following Lemma provides a simple extension of this result.

**LEMMA 2.** *Given a symmetric mapping $\widehat{\Psi}_\varepsilon^{m-1}$, the mapping*

$$\widehat{\Psi}_\varepsilon^m \equiv \Psi_{\varepsilon/2}^k \circ \widehat{\Psi}_\varepsilon^{m-1} \circ \Psi_{\varepsilon/2}^{k,*}$$

*is also symmetric.*

*Proof.*

$$\begin{aligned}
\widehat{\Psi}^m_{-\varepsilon} &= \Psi^k_{-\varepsilon/2} \circ \widehat{\Psi}^{m-1}_{-\varepsilon} \circ \Psi^{k,*}_{-\varepsilon/2} \\
&= \left[\Psi^{k,*}_{\varepsilon/2}\right]^{-1} \circ \left[\widehat{\Psi}^{m-1}_{\varepsilon}\right]^{-1} \circ \left[\Psi^k_{\varepsilon/2}\right]^{-1} \\
&= \left[\Psi^k_{\varepsilon/2} \circ \widehat{\Psi}^{m-1}_{\varepsilon} \circ \Psi^{k,*}_{\varepsilon/2}\right]^{-1} \\
&= \left[\widehat{\Psi}^m_{\varepsilon}\right]^{-1}
\end{aligned}$$

which satisfies the definition of a symmetric mapping. $\qquad\square$

Note that since the adjoint of the adjoint is the original mapping, i.e. $\Psi^{k,**}_{\varepsilon/2} = \Psi^k_{\varepsilon/2}$, Lemma 2 can be also equivalently stated as $\Psi^{k,*}_{\varepsilon/2} \circ \widehat{\Psi}^{m-1}_{\varepsilon} \circ \Psi^k_{\varepsilon/2}$ being symmetric.

For $L$ even, let $m = L/2$, $k = m$ and define the mapping

$$\widehat{\Psi}^{L/2+1/2}_{\varepsilon} = \Psi^{L/2}_{\varepsilon/2} \circ \Psi^{L/2,*}_{\varepsilon/2}$$

which, using (9.1), is symmetric. Then, let

$$\widehat{\Psi}^{L/2+1}_{\varepsilon} = \Psi^{L/2-1,*}_{\varepsilon/2} \circ \widehat{\Psi}^{L/2+1/2}_{\varepsilon} \circ \Psi^{L/2-1}_{\varepsilon/2}$$

and further

$$\begin{aligned}
\widehat{\Psi}^{L/2+m+1/2}_{\varepsilon} &= \Psi^{L/2+m}_{\varepsilon/2} \circ \widehat{\Psi}^{L/2+m}_{\varepsilon} \circ \Psi^{L/2+m,*}_{\varepsilon/2} \\
\widehat{\Psi}^{L/2+m+1}_{\varepsilon} &= \Psi^{L/2-m-1,*}_{\varepsilon/2} \circ \widehat{\Psi}^{L/2+1/2,L/2+1/2}_{\varepsilon} \circ \Psi^{L/2-m-1}_{\varepsilon/2}
\end{aligned}$$

for $m = 1, \ldots, L/2 - 1$. The final composite mapping $\widehat{\Psi}^L_{\varepsilon}$ then takes the form

$$\widehat{\Psi}_{\varepsilon} \equiv \widehat{\Psi}^L_{\varepsilon} = \Psi^{0,*}_{\varepsilon/2} \circ \Psi^{L-1}_{\varepsilon/2} \circ \Psi^{1,*}_{\varepsilon/2} \circ \Psi^{L-2}_{\varepsilon/2} \circ \ldots \circ \Psi^{L/2}_{\varepsilon/2} \circ \Psi^{L/2,*}_{\varepsilon/2} \circ \ldots \circ \Psi^{L-2,*}_{\varepsilon/2} \circ \Psi^1_{\varepsilon/2} \circ \Psi^{L-1,*}_{\varepsilon/2} \circ \Psi^0_{\varepsilon/2}$$

Symmetry of $\widehat{\Psi}_{\varepsilon}$ follows by repeated application of Lemma 2.

The mappings $\Psi^k_{\varepsilon}$ and $\Psi^{*,k}_{\varepsilon}$ are special cases of an implicit partitioned Runge-Kutta method (Leimkuhler and Reich, 2004, p. 150-151) and thus the existence and uniqueness of their solutions follows from Lemma 1. The uniqueness of the soluton to $\Psi^k_{\varepsilon/2}$ and $\Psi^{k,*}_{\varepsilon/2}$ for each $k$ implies that there is a unique solution to $\widehat{\Psi}_{\varepsilon}$. Such solution is equivalent to the one given by AUHMC since the AUHMC fixed-point $M(\overline{\theta_r, \theta^*_{r+1}})$ is identical to the fixed point $M_k = M(\overline{\theta}^k, \overline{\theta}^L)$ of $\widehat{\Psi}_{\varepsilon}$ that solves $\widehat{\Psi}_{\varepsilon}$. Since, by Lemma 1, the solution to AUHMC is unique, AUHMC implements $\widehat{\Psi}_{\varepsilon}$ which is a symmetric and time reversible mapping, yielding the detailed balance condition of Theorem 1.

Equivalently, from the definition of $\widehat{\Psi}_{\varepsilon}$ it follows directly that reversing the momentum at $(\theta^*_{r+1}, \gamma^*_{r+1})$ and applying AUHMC solves $\widehat{\Psi}^{-1}_{-\varepsilon}$ which, due to symmetry of $\widehat{\Psi}_{\varepsilon}$, equals $\widehat{\Psi}^*_{\varepsilon}$, following the same proposal path back to $(\theta_r, \gamma_r)$ having negated the momentum again after the final step. This satisfies the definition of reversibility for AUHMC.

We can make an analogy between the pair of Euler B and A methods (Leimkuhler and Reich, 2004, p. 84) and the pair of $\Psi^k_{\varepsilon/2}$ and $\Psi^{k,*}_{\varepsilon/2}$. In the former pair, the difference is in the point at which we evaluate directional derivatives ($\theta^k$ or $\theta^{k+1}$). In the pair of $\Psi^k_{\varepsilon/2}$ and $\Psi^{k,*}_{\varepsilon/2}$, the difference is in the number of HMC steps needed to reach $\overline{\theta}^{k+1}$ and $\overline{\theta}^{L+1}$, which at the solution equal to $\theta^0$ and

$\theta^L$ respectively, but the directional derivatives are always the same, taken with $M_k$ evaluated at the endpoints of the proposal sequence. However, neither $\Psi_{\varepsilon/2}^k$ nor $\Psi_{\varepsilon/2}^{k,*}$ is symmetric on its own, and hence we need their concatenation to attain symmetry of the composite mapping.

At the implicit solution of $\Psi_{\varepsilon/2}^k$,

$$\theta^{k+1} = \theta^k + \frac{\varepsilon}{2} M(\overline{\theta_r, \theta_{r+1}^*})^{-1} \gamma^{k+1}$$
$$\gamma^{k+1} = \gamma^k + \frac{\varepsilon}{2} \nabla_{\theta_i} \ln \pi(\theta^k)$$

in analogy to the Euler-B method. Also, due to the symmetry of $M_k$ in $\overline{\theta}^{k+1}$ and $\overline{\theta}^{L+1}$ from Assumption 1, at the solution of $\Psi_{\varepsilon/2}^{*L-k+1}$,

$$\theta^{k+1} = \theta^k + \frac{\varepsilon}{2} M(\overline{\theta_{r+1}^*, \theta_r})^{-1} \gamma^k$$
$$\gamma^{k+1} = \gamma^k + \frac{\varepsilon}{2} \nabla_{\theta_i} \ln \pi(\theta^{k+1})$$

in analogy to the Euler-A method. These half-steps are performed by AUHMC during its proposal sequence.

### 9.3. M-H Acceptance Probability

The derivation of the M-H acceptance probability form is standard in the HMC literature and we merely adapt it to the AUHMC notation below. Denote by $q(\theta_{r+1}^*, \gamma_{r+1}^*; \theta_r^0, \gamma_r^0)$ the proposal density and by $q(\theta_r^0, \gamma_r^0; \theta_{r+1}^*, \gamma_{r+1}^*)$ the reverse proposal density. Given $(\theta_r^0, \gamma_r^0)$, $q(\theta_{r+1}^*, \gamma_{r+1}^*; \theta_r^0, \gamma_r^0)$ is constructed by the method of change of variables based on the sequence of steps given by the AUHMC mapping $T_k$ for $k = 1, \ldots, L$. Since $T_k$ is deterministic, placing the Dirac delta $\delta(\cdot, \cdot) = 1$ unit probability mass at each $(\theta_r^k, \gamma_r^k)$, applying successive transformations $T_k$ yields

$$q(\theta_{r+1}^*, \gamma_{r+1}^*; \theta_r^0, \gamma_r^0) = \left| \det \nabla T(\theta_r^L, \gamma_r^L; \theta_r^{L-1}, \gamma_r^{L-1}) \right|^{-1} \times \left| \det \nabla T(\theta_r^{L-1}, \gamma_r^{L-1}; \theta_r^{L-2}, \gamma_r^{L-2}) \right|^{-1} \times \ldots$$
$$(9.2) \qquad \ldots \times \left| \det \nabla T(\theta_r^2, \gamma_r^2; \theta_r^1, \gamma_r^1) \right|^{-1} \left| \det \nabla T(\theta_r^1, \gamma_r^1; \theta_r^0, \gamma_r^0) \right|^{-1} \delta(\gamma_r^0, \theta_r^0)$$

where $\nabla T(\theta_r^k, \gamma_r^k; \theta_r^{k-1}, \gamma_r^{k-1})$ denotes the Jacobian matrix of the transformation $T_k$ with respect to $\theta_r^k$ and $\gamma_r^k$ for each $k = 1, \ldots, L$.

Denote by $\widetilde{T}_k$ the reverse mapping obtained from $T_k$ by reversing the signs in the Hamiltonian proposal dynamics. Then

$$q(\theta_r^0, \gamma_r^0; \theta_{r+1}^*, \gamma_{r+1}^*) = \left| \det \nabla \widetilde{T}(\widetilde{\theta}_r^L, \widetilde{\gamma}_r^L; \widetilde{\theta}_r^{L-1}, \widetilde{\gamma}_r^{L-1}) \right|^{-1} \times \left| \det \nabla \widetilde{T}(\widetilde{\theta}_r^{L-1}, \widetilde{\gamma}_r^{L-1}; \widetilde{\theta}_r^{L-2}, \widetilde{\gamma}_r^{L-2}) \right|^{-1} \times \ldots$$
$$(9.3) \qquad \ldots \times \left| \det \nabla \widetilde{T}(\widetilde{\theta}_r^2, \widetilde{\gamma}_r^2; \widetilde{\theta}_r^1, \widetilde{\gamma}_r^1) \right|^{-1} \left| \det \nabla \widetilde{T}(\widetilde{\theta}_r^1, \widetilde{\gamma}_r^1; \widetilde{\theta}_r^0, \widetilde{\gamma}_r^0) \right|^{-1} \delta(\widetilde{\gamma}_r^0, \theta_{r+1}^*)$$

with $(\widetilde{\theta}_r^0, \widetilde{\gamma}_r^0) = (\theta_{r+1}^*, \gamma_{r+1}^*)$. Conditional on $M(\overline{\theta_r, \theta_{r+1}^*})$ satisfying Assumption 1, the leapfrog transformation defined by (3.5)-(3.7) satisfies

$$(9.4) \qquad (\theta_r^k, \gamma_r^k) = (\widetilde{\theta}_r^{L-k+1}, \widetilde{\gamma}_r^{L-k+1}) \quad \text{for each } k = 1, \ldots, L$$

Then

$$(9.5) \quad \left| \det \nabla T(\theta_r^k, \gamma_r^k; \theta_r^{k-1}, \gamma_r^{k-1}) \right|^{-1} = \left| \det \nabla \widetilde{T}(\widetilde{\theta}_r^{L-k+1}, \widetilde{\gamma}_r^{L-k+1}; \widetilde{\theta}_r^{L-k}, \widetilde{\gamma}_r^{L-k}) \right| \quad \text{for each } k = 1, \ldots, L$$

and hence $q(\theta^*_{r+1}, \gamma^*_{r+1}; \theta^0_r, \gamma^0_r) = q(\theta^*_{r+1}, \gamma^*_{r+1}; \theta^0_r, \gamma^0_r)$.

The ratio in the acceptance probability (2.2) then satisfies detailed balance in the Metropolis form

(9.6)
$$\frac{\pi(\theta^*_{r+1}, \gamma^*_{r+1})q(\theta^0_r, \gamma^0_r; \theta^*_{r+1}, \gamma^*_{r+1})}{\pi(\theta_r, \gamma_r)q(\theta^*_{r+1}, \gamma^*_{r+1}; \theta^0_r, \gamma^0_r)} = \frac{\pi(\theta^*_{r+1}, \gamma^*_{r+1})}{\pi(\theta_r, \gamma_r)}$$

since all the Jacobian terms cancel out due to (9.5). By definition of the Hamiltonian equation in (3.1), the ratio in (9.6) is then equivalent to

$$\ln \pi(\theta^*_{r+1}) - \ln \pi(\theta^0_r) + \ln \phi(\gamma^*_{r+1}; 0, M(\overline{\theta_r, \theta^*_{r+1}})) - \ln \phi(\gamma^0_r; 0, M(\overline{\theta_r, \theta^*_{r+1}}))$$

## 10. Appendix D: Fisher Information for the Multivariate Normal Density

For the univariate case,

$$F(\theta) = N \begin{bmatrix} \sigma^{-2} & 0 \\ 0 & \frac{1}{2}\sigma^{-4} \end{bmatrix}$$

and for the multivariate case

$$F(\theta) = N \begin{bmatrix} \Sigma^{-1} & 0 \\ 0 & \frac{1}{2}D'_m\left(\Sigma^{-1} \otimes \Sigma^{-1}\right)D_m \end{bmatrix}$$

where $D_m$ is the duplication matrix (Magnus and Neudecker, 2007). In our empirical application we used the numerical approximation to the diagonal of $F(\theta)$ instead of the full matrix for faster speed of the MC runs.

# References

AKHMATSKAYA, E., N. BOU-RABEE, AND S. REICH (2009): "A comparison of generalized hybrid monte carlo methods with and without momentum flip," *Journal of Computational Physics*, 228(6), 2256–2265.

BAUWENS, L., C. S. BOS, H. K. VAN DIJK, AND R. D. VAN OEST (2004): "Adaptive radial-based direction sampling: some flexible and robust Monte Carlo integration methods," *Journal of Econometrics*, 123, 201–225.

BESKOS, A., N. S. PILLAI, G. O. ROBERTS, J. M. SANZ-SERNA, AND A. M. STUART (2010): "Optimal tuning of the Hybrid Monte-Carlo Algorithm," Working paper, arxiv:1001.4460v1 [math.pr].

CHIB, S., AND E. GREENBERG (1995): "Understanding the Metropolis-Hastings Algorithm," *American Statistician*, 49(4), 327–335.

CHIB, S., AND S. RAMAMURTHY (2010): "Tailored randomized block MCMC methods with application to DSGE models," *Journal of Econometrics*, 155(1), 19 – 38.

DELLAPORTAS, P., AND I. D. VRONTOS (2007): "Modelling volatility asymmetries: a Bayesian analysis of a class of tree structured multivariate GARCH models," *Econometrics Journal*, 10(3), 503520.

DING, Z., AND R. ENGLE (2001): "Large Scale Conditional Covariance Matrix Modeling, Estimation and Testing," *Academia Economic Papers*, 29, 157184.

DUANE, S., A. KENNEDY, B. PENDLETON, AND D. ROWETH (1987): "Hybrid Monte Carlo," *Physics Letters B*, 195(2), 216–222.

ENGLE, R. F. (2002): "Dynamic conditional correlation: A simple class of multivariate generalized autoregressive conditional heteroskedasticity models," *Journal of Business and Economic Statistics*, 20, 339–350.

ENGLE, R. F., AND K. F. KRONER (1995): "Multivariate Simultaneous Generalized ARCH," *Econometric Theory*, 11(1), 122–150.

ENGLE, R. F., N. SHEPHARD, AND K. SHEPPARD (2009): "Fitting Vast Dimensional Time-Varying Covariance Models," Available at SSRN: http://ssrn.com/abstract=1354497.

GEWEKE, J. (2005): *Contemporary Bayesian Econometrics and Statistics*. Wiley, Hoboken, New Jersey.

GEYER, C. J. (1992): "Practical Markov Chain Monte Carlo," *Statistal Science*, 7, 473483.

GIROLAMI, M., AND B. CALDERHEAD (2011): "Riemann Manifold Langevin and Hamiltonian Monte Carlo Methods (with Discussion)," *J. R. Stat. Soc. B*, 73(2), 123–214.

GUPTA, R., G. KILCUP, AND S. SHARPE (1988): "Tuning the hybrid monte carlo algorithm," *Physical Review D*, 38(4), 1278–1287.

HAFNER, C. M., AND H. HERWARTZ (2008): "Analytical quasi Maximum Likelihood Inference in Multivariate Volatility Models," *Metrika*, 67, 219–239.

HAIRER, E., C. LUBICH, AND G. WANNER (2003): "Geometric numerical integration illustrated by the Störmer–Verlet method," *Acta Numerica*, pp. 399–450.

HAIRER, E., S. NØRSETT, AND G. WANNER (1993): *Solving Ordinary Differential Equations I.* Springer-Verlag, 2 edn.

HOLMES, C. C., AND L. HELD (2006): "Bayesian Auxiliary Variable Models for Binary and Multinomial Regression," *Bayesian Analysis*, 1(1), 145–168.

HUDSON, B., AND R. GERLACH (2008): "A Bayesian approach to relaxing parameter restrictions in multivariate GARCH models," *TEST: An Official Journal of the Spanish Society of Statistics and Operations Research*, 17(3), 606–627.

ISHWARAN, H. (1999): "Applications of Hybrid Monte Carlo to Generalized Linear Models: Quasi-complete Separation and Neural Networks," *Journal of Computational and Graphical Statistics*, 8, 779–799.

LEIMKUHLER, B., AND S. REICH (2004): *Simulating Hamiltonian Dynamics.* Cambridge University Press.

LIESENFELD, R., AND J.-F. RICHARD (2006): "Classical and Bayesian Analysis of Univariate and Multivariate Stochastic Volatility Models," *Econometric Reviews*, 25(2-3), 335–360.

LIU, J. S. (2004): *Monte Carlo Strategies in Scientific Computing.* Springer Series in Statistics.

MAGNUS, J., AND H. NEUDECKER (2007): *Matrix Differential Calculus with Applications in Statistics and Econometrics.* John Wiley & Sons.

NEAL, R. M. (1993): "Probabilistic Inference Using Markov Chain Monte Carlo Methods," Technical report crg-tr-93-1, Dept. of Computer Science, University of Toronto.

———— (2010): "MCMC using Hamiltonian dynamics," in *Handbook of Markov Chain Monte Carlo*, ed. by S. Brooks, A. Gelman, G. Jones, , and X.-L. Meng. Chapman & Hall / CRC Press.

OSIEWALSKI, J., AND M. PIPIEN (2004): "Bayesian comparison of bivariate ARCH-type models for the main exchange rates in Poland," *Journal of Econometrics*, 123(2), 371 – 391.

PITT, M. K., AND N. SHEPHARD (1997): "Likelihood Analysis of Non-Gaussian Measurement Time Series," *Biometrika*, 84, 653–667.

ROBERT, C. P., AND G. CASELLA (2004): *Monte Carlo statistical methods.* Springer, New York, second edn.

ROBERTS, G., AND O. STRAMER (2003): "Langevin diffusions and Metropolis-Hastings algorithms," *Methodology and Computing in Applied Probability*, 4, 337–358.

ROBERTS, G. O., AND J. S. ROSENTHAL (1998): "Optimal Scaling of Discrete Approximations to Langevin Diffusions," *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 60(1), 255–268.

TUCKERMAN, M., B. BERNE, G. MARTYNA, AND M. KLEIN (1993): "Efficient molecular dynamics and hybrid monte carlo algorithms for path integrals," *The Journal of Chemical Physics*, 99(4), 2796–2808.

|        | Mean    | Stdev  | Skewness | Ex Kurtosis | | Sample Correlation | | | |
|--------|---------|--------|----------|-------------|---|--------|--------|--------|--------|
| AUD/USD | -0.0074 | 0.7019 | 0.4444 | 1.8868 | 1 | 0.4814 | 0.4776 | 0.5454 | 0.3614 |
| GBP/USD | -0.0074 | 0.5281 | 0.0517 | 0.6827 | | 1 | 0.3233 | 0.7123 | 0.3787 |
| CAD/USD | -0.0144 | 0.4657 | 0.0039 | 0.7231 | | | 1 | 0.3874 | 0.2698 |
| EUR/USD | -0.0115 | 0.6268 | 0.0640 | 0.6330 | | | | 1 | 0.3995 |
| JPY/USD | 0.0087 | 0.6031 | -0.2978 | 1.5496 | | | | | 1 |

Table 3: Summary Statistics

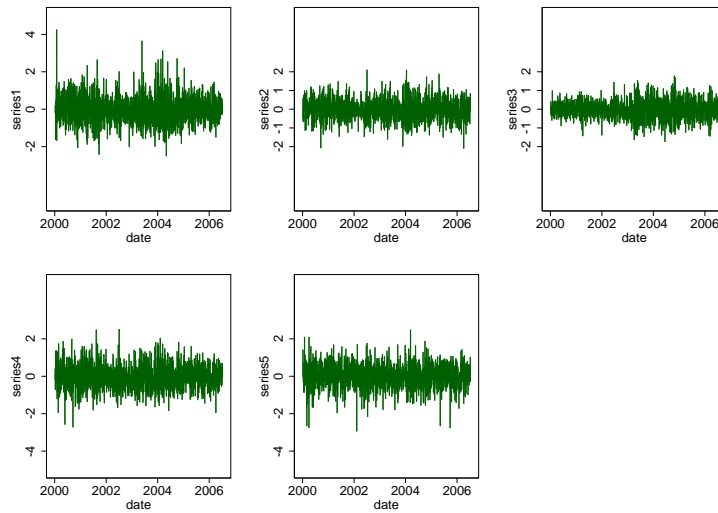| $\chi^2$ *quantile* | 0.75 | 0.90 | 0.99 |
|--------------------------------|-----------|-----------|-----------|
| *BEKK* | -99549.8 | -99549.6 | -99549.6 |
| *Diagonal BEKK (full C)* | -99584.9 | -99584.7 | -99584.6 |
| *Diagonal BEKK (diagonal C)* | -100209.9 | -100209.7 | -100209.6 |

Table 4: Log-marginal likelihoods



Figure 3: Time-series of log-differences in foreign exchange rates
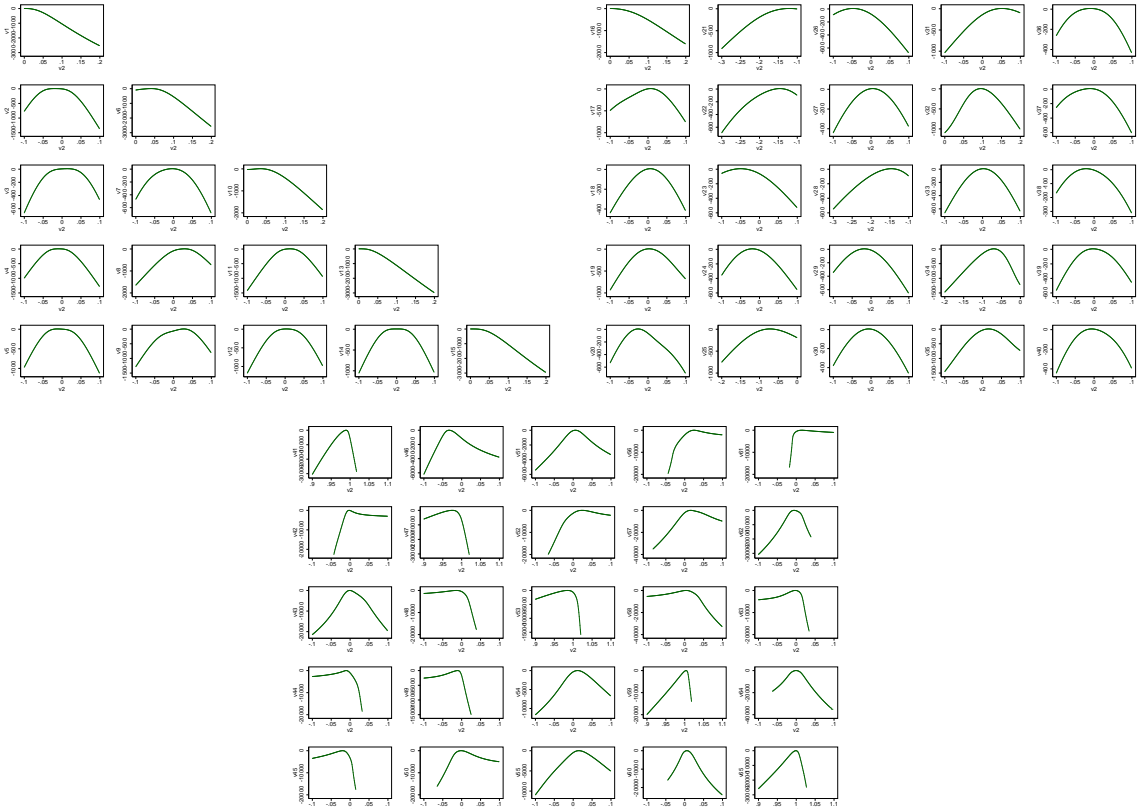
Figure 4: Conditional Log-Posterior Kernels for parameter matrices $C, F$ and $G$ from the BEKK model. Each parameter is plotted conditional on the other parameters being fixed at a point of high mass in the posterior density close to the modal values.
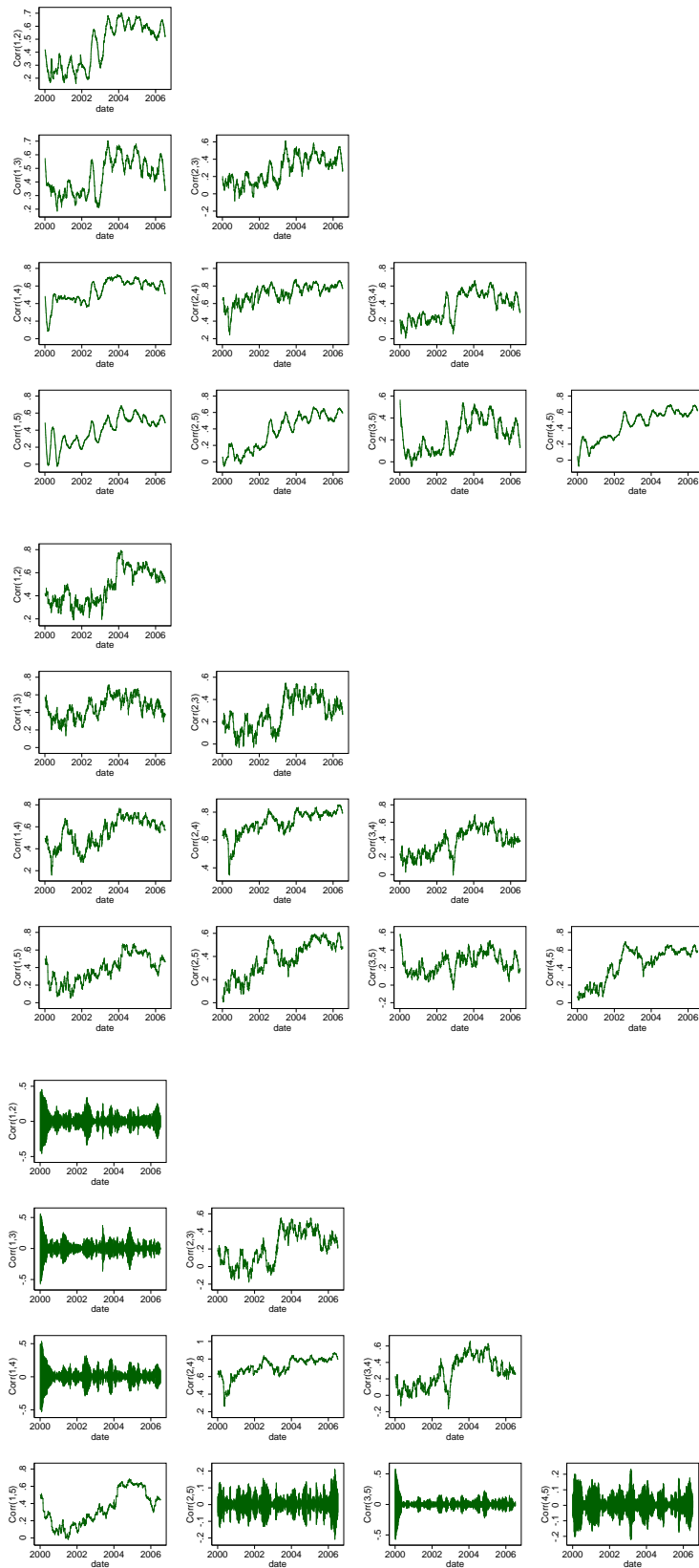
Figure 5: Conditional Correlations: BEKK; Diagonal F, G BEKK; Diagonal C, F, G BEKK