

Estimating HANK for Central Banks

Sushant Acharya, William Chen, Marco Del Negro, Keshav Dogra, Aidan Gleich
Shlok Goyal, Ethan Matlin, Donggyu Lee, Reza Sarfati, Sikata Sengupta*

Bank of Canada, Federal Reserve Bank of New York, Harvard, MIT

May 5, 2023

please do not circulate without permission

Abstract

We provide a toolkit for efficient online estimation of HANK model, based on Sequential Monte Carlo methods. We use this toolkit to compare the out-of-sample forecasting accuracy of a prominent HANK model, Bayer et al. (2022), to that of Smets and Wouters (2007, SW)'s representative agent (RANK) model. We find that HANK's accuracy for real activity variables is notably inferior to that of SW. The results for consumption are disappointing, given that the main difference between RANK and HA is the replacement of the RA Euler equation with the aggregation of individual households' consumption policy functions, which reflects inequality.

JEL CLASSIFICATION: C11, C32, D31, E32, E37, E52

KEY WORDS: HANK, Bayesian inference, sequential Monte Carlo methods

*Correspondence: Marco Del Negro marco.delnegro@ny.frb.org, We thank participants at the XXV Annual Conference of the Central Bank of Chile "Heterogeneity in Macroeconomics: Implications for Monetary Policy," for which this paper was written, and especially our discussant, Markus Kirchner, for helpful comments. The views expressed in this paper are those of the authors and do not necessarily reflect the position of the Bank of Canada, the Federal Reserve Bank of New York, or the Federal Reserve System.

1 Introduction

Central banks are very interested in investigating questions surrounding inequality and its relationship with monetary policy. This is arguably for very good reasons. First of all, inequality has become a central issue in many if not all countries. It therefore seems important to ask how central bank policies affect inequality. Second, even if central bankers were not concerned with the answer to the above question, they ought to be concerned with the fact that inequality changes the transmission mechanism of monetary policy, as forcefully argued in Kaplan et al. (2018) and Ahn et al. (2018). Some central banks have indeed shown interest in these topics (in fact, the title of the conference on which this volume is based is “Heterogeneity in Macroeconomics: Implications for Monetary Policy”) and a few have begun to develop models that speak to the interaction of monetary policy and inequality, such as heterogeneous agents New-Keynesian (HANK) models, following the seminal work by Kaplan et al. (2018).

Models serve many purposes, and for some of these purposes a model’s ability to fit the data—that is, to adequately describe the data from a quantitative point of view—is important, especially for central banks. After all, the popularity of representative agent DSGE models such as Smets and Wouters (2007, henceforth, SW) since the beginning of the century is largely due to these models’ ability to forecast with an accuracy that is at least comparable to that of other models previously used in central banks, such as vector autoregressions. Even if forecasting is not the main purpose of a model—and arguably it is not the main purpose of DSGE models—it is a way to test its reliability in providing answers to quantitative questions: forecasting accuracy lends quantitative credibility.

These considerations prompt us to ask: What is the forecasting accuracy of HANK models? To the extent that these model have a more realistic transmission mechanism than representative agent models, one would hope that this translates into a better forecasting performance. This is particularly true for aggregate consumption, since the main difference between SW-type DSGEs and HANK models is the replacement of the representative agent Euler equation, which determines consumption in standard DSGEs, with the aggregation of individual households’ consumption policy functions. These consumption policy functions depend, among other things, on the wealth distribution in the economy, that is, on inequality. This is the first paper to our knowledge that provides an assessment of the out-of-sample forecasting accuracy of HANK models.

From a computational point of view the task of performing an out-of-sample forecasting accuracy exercise is not trivial, as it involves estimating a HANK model over and over, for each of the several vintages of data for which we want to compute forecasts. Concretely, our forecasting exercise begins in the first quarter of 2000 and ends in the last quarter of 2019, for a total of 80 periods. For each period we estimate the the model using Bayesian methods—the same approach used by SW and BBL. Each estimation is very costly in

computational terms for HANK if one calculates the likelihood using the Kalman filter, since these models have a very large state-space which includes the distribution of wealth (both liquid and illiquid, in a two asset HANK model) across households.¹

All of the growing literature estimating HANK models using Bayesian methods (see Winberry, 2018, Auclert et al., 2021, Bayer et al., 2022, and Lee, 2021, among others) use the standard Markov Chain Monte Carlo approach followed in the representative DSGE model literature to obtain draws from the posterior distribution (eg, An and Schorfheide, 2007), and featured in popular packages such as `Dynare`. This approach has two drawbacks. First, it cannot be naturally parallelized, being a Markov Chain-based algorithm. Second, one has to start from scratch every new estimation. For example, if one just estimated the model up to 2000Q1, and then adds only one more quarter of data, with Markov Chain methods one has to start anew as if they had not done any previous estimation, even though one suspects that the posterior distribution may not be all that different.

This paper deviates from this trend and uses a Monte Carlo method that can be readily parallelized—Sequential Monte Carlo. This parallelization makes it feasible to estimate models even when each likelihood computation takes a substantial amount of time. This method has another crucial advantage, namely that models can be estimated “online.” What online estimation means is that the swarm of particles describing the posterior distribution computed for the estimation up to 2000Q1 can be used to jump-start the estimation with one or more quarters of data, thereby making it considerably faster. This online feature is what makes repeated estimation, and therefore our forecasting accuracy exercise, possible.² While these methods are not new (see Cai et al., 2021), one contribution of this paper is to explain how and why they work to an audience with little or no background in Monte Carlo methods, so that this paper may serve as a blueprint for central bank researchers planning to estimate HANK models and use them in routine policy analysis and forecasting exercises. We also plan to share the code used in our forecasting exercise on `GitHub`.

As anticipated above, the other contribution of this paper is to provide a forecasting accuracy assessment of a HANK model. While several HANK models have been developed, we use that of Bayer et al. (2022, henceforth, BBL) in this paper. We do this because in their frontier contribution the authors put particular care in the empirical fit of their model, making sure that they include all the shocks and frictions that make SW-type models empirically successful. In other words, the BBL model is the closest thing to a HANK version of SW. We then ask: Does this model forecast macro time series better than the original SW?

¹The advantage of the so-called sequence-space Jacobian approach to solving and estimating HANK models championed by Auclert et al. (2021) is that it circumvents the issue of the large state-space associated with carrying around a set of distributions. As far as we understand, however, for the time being this approach cannot deal with missing data, which is an issue for our application where some data—especially those related to inequality—are not available for the entire time series.

²The methods described in this paper can be used in the context of limited information approaches, such as those used by Hagedorn et al. (2018) who estimate a HANK model using impulse-response matching as in Christiano et al. (2005).

Unfortunately, the answer seems to be no. For some series such as inflation, the forecasting performance is similar. For other series, notably for consumption growth, accuracy for the HANK model is much worse than for the representative agent model, which is particularly disappointing for the reasons discussed above.

What are the reasons for, and the implications of, the relatively worse forecasting performance of this HANK model compared to SW? We suspect that one key reason is that many parameters in HANK—namely those affecting the model’s steady state—are still calibrated. This is not necessarily for a philosophical choice on the part of the HANK modelers, but because recomputing the steady state is extremely costly. If this suspicion is correct, these findings pose a computational challenge to HANK researchers interested in estimation. For sure the findings should be interpreted as a motivation to do more research on HANK models, as opposed to sticking to representative agent models. Inequality is one of the critical issues of our times—no matter the forecasting performance of HANK models. The fact that the latter can be improved is a stimulus for further efforts, especially from central bank researchers who want to use these models for quantitative purposes.

In the remainder of the paper, section 2 presents BBL’s model and solution approach so to make the paper self contained, section 3 describes the Sequential Monte Carlo algorithm and the online estimation approach used to perform the forecasting exercise, section 4 discusses the results, and section 5 concludes.

2 Model

This paper employs a model developed in BBL, which augments standard New Keynesian DSGE models, such as those presented in SW or Christiano et al. (2005), with heterogeneous agents and incomplete markets. The model incorporates standard shocks and frictions utilized in DSGE models that match aggregate data. Moreover, it is also capable of reproducing notable characteristics of households heterogeneity that are deemed important in the literature, such as heterogeneous wealth and income composition and the presence of wealthy hand-to-mouth households. BBL show that, when the model is estimated on aggregate data, it can reproduce the business cycle dynamics of aggregate data as well as of observed U.S. inequality. As the model is entirely taken from BBL, we will provide only a brief description of the model environment below in order to make the paper self contained.³

³We use the version of model available at https://github.com/BASEforHANK/HANK_BusinessCycleAndInequality as of June 2022, when we began this project. The latest version of the model, as described in Bayer et al. (2022) has two minor differences compared to the model adopted in this paper. First, the latest version of the model has a different formulation for the liquid asset return. Specifically, BBL assume that entrepreneurs sell claims to a fraction of profits as liquid shares, and the liquid asset return is the weighted average of the interest on government bonds and the return on these shares, which consists of profit payouts and the realized capital gain. In addition, BBL assumes that time-varying income risks respond to output growth, which makes income risks either procyclical or countercyclical. These modifications allow BBL to better explain inequality

2.1 Households

There exists a unit mass of infinitely-lived households, indexed by i , that maximize their lifetime utility

$$\mathbb{E}_0 \sum_{t=0}^{\infty} \beta^t u(c_{it}, n_{it} | h_{it}), \quad (1)$$

where β denotes the discount factor, c_{it} denotes consumption, and n_{it} denotes hours worked. The instantaneous utility function $u(\cdot)$ takes the form of Greenwood et al. (1988) preference

$$u(c_{it}, n_{it} | h_{it}) = \frac{x_{it}^{1-\xi} - 1}{1-\xi}, \quad x_{it} = c_{it} - h_{it}^{1-\bar{\tau}^P} \frac{n_{it}^{1+\gamma}}{1+\gamma}, \quad (2)$$

where $\bar{\tau}^P$ is the steady-state level of tax progressivity, ξ is the coefficient of relative risk aversion, and γ is the inverse of Frisch elasticity. Regarding idiosyncratic labor productivity h_{it} , there are two types of households; workers ($h_{it} \neq 0$) and entrepreneurs ($h_{it} = 0$). Specifically, idiosyncratic productivity $h_{it} = \frac{\tilde{h}_{it}}{\int \tilde{h}_{it} di}$ evolves as follows

$$\tilde{h}_{it} = \begin{cases} \exp(\rho_h \log \tilde{h}_{it-1} + \epsilon_{it}^h) & \text{with probability } 1 - \zeta \text{ if } h_{it-1} \neq 0, \\ 1 & \text{with probability } \iota \text{ if } h_{it-1} = 0, \\ 0 & \text{else.} \end{cases} \quad (3)$$

The above equation implies that workers become entrepreneurs with the probability ζ or continue to be workers with the probability $1 - \zeta$. While being workers, labor productivity h_{it} evolves according to a log AR(1) process with the autocorrelation coefficient ρ_h . The shocks ϵ_{it}^h are normally distributed with the variance $\sigma_{h,t}^2$, which is time-varying according to the following process

$$\sigma_{h,t}^2 = \bar{\sigma}_h^2 \exp(\hat{s}_t), \quad (4)$$

$$\hat{s}_{t+1} = \rho_s \hat{s}_t + \epsilon_t^\sigma, \quad (5)$$

where the shocks ϵ_t^σ follows a normal distribution with zero mean and the standard deviation σ_s .⁴ Workers earn wage income $w_t h_{it} n_{it}$, where w_t is the real wage paid to households by labor unions. In addition, rents from unions Π_t^U are equally distributed among workers. Similar to workers, entrepreneurs become workers with the probability ι or maintain their entrepreneur status with the probability $1 - \iota$. When entrepreneurs become workers, their productivity becomes one. Entrepreneurs do not supply labor and, instead, receive profits Π_t^F generated in the firm sector, except for rents of unions.

Markets are incomplete, and thus, households self-insure against income risks by saving in two types of assets: illiquid capital and liquid bonds. Capital as an asset is illiquid in the sense that only a fraction series and their income risk estimates with their model.

⁴In the latest version of BBL, they assume that the level of income risks is affected by the output growth, i.e., $\hat{s}_{t+1} = \rho_s \hat{s}_t + \Sigma_Y \frac{Y_{t+1}}{Y_t} + \epsilon_t^\sigma$. Depending on the sign of the coefficient Σ_Y , idiosyncratic income risks are either pro or counter-cyclical in the model. This setup allows BBL to better capture the dynamics of income risks with their model.

λ of households are allowed to adjust their capital holdings each period. In contrast, households can freely adjust their liquid bond holdings. Given households' income sources and choice variables, the household's budget constraint can be written as follows.

$$c_{it} + b_{it+1} + q_t k_{it+1} = b_{it} \frac{R_t}{\pi_t} + (q_t + r_t) k_{it} + (1 - \tau_t^L)(w_t h_{it} n_{it})^{1-\bar{\tau}_t^P} + \mathbb{1}_{h_{it} \neq 0} (1 - \tau_t) \Pi_t^U + \mathbb{1}_{h_{it} = 0} (1 - \tau_t^L) (\Pi_t^E)^{1-\bar{\tau}_t^P}, \quad k_{it+1} \geq 0, \quad b_{it+1} \geq \underline{b}, \quad (6)$$

where b_{it} is real liquid bonds, k_{it} is capital stock, q_t is the price of capital, r_t is dividend on capital holdings, $\pi_t = \frac{P_t}{P_{t-1}}$ is the gross inflation rate, and $\underline{b} < 0$ is an exogenous borrowing limit. Workers' labor income and entrepreneurs' profit income are taxed progressively. The two tax rates, τ_t^L and τ_t^P , determine the degree of tax progressivity. The union profit is taxed uniformly at the average tax rate τ_t . The realized return on the liquid assets R_{it} depends on if households borrow or not

$$R_{it} = \begin{cases} A_t R_t^b & \text{if } b_{it} \geq 0 \\ A_t R_t^b + \bar{R} & \text{if } b_{it} < 0. \end{cases} \quad (7)$$

The coefficient A_t is "risk-premium shock", which reflects an intermediation efficiency, and \bar{R} is the borrowing premium. R_t^B is the nominal interest rate on government bonds, which is determined by the monetary authority. In the model described in Bayer et al. (2022), they assume that entrepreneurs sell claims to a fraction ω_Π of profits at the price of q_t^Π as liquid shares and these shares become a part of the household's liquid asset portfolio as well. Thus, the liquid asset return is the weighted average of the return on government bonds and the return on profit shares. Consequently, dynamics of profit shares also affect the liquid asset return in the model.

Since households may or may not be able to adjust their illiquid asset holdings, the household's problem is characterized by three functions; value function V_t^a when households are allowed to adjust their capital holdings, the function V_t^n when households are not allowed to adjust, and the expected value in the next period \mathbb{W}_{t+1}

$$V_t^a(b, k, h) = \max_{b', k'} u[(x(b, b', k, k', h))] + \beta \mathbb{E}_t \mathbb{W}_{t+1}(b', k', h'), \quad (8)$$

$$V_t^n(b, k, h) = \max_{b'_n} u[(x(b, b'_n, k, k, h))] + \beta \mathbb{E}_t \mathbb{W}_{t+1}(b'_n, k', h'), \quad (9)$$

$$\mathbb{W}_{t+1}(b', k', h') = \lambda V_{t+1}^a(b', k', h') + (1 - \lambda) V_{t+1}^n(b', k', h'), \quad (10)$$

where $x(b, b', k, k', h) = c(b, b', k, k', h) - h^{1-\bar{\tau}_t^P} \frac{n(w)^{1+\gamma}}{1+\gamma}$ is the household's composite demand for goods and leisure.⁵ Maximization is subject to the corresponding budget constraint described above.

⁵Because of the specific form of GHH preference used in the model, all workers supply the same amount of labor, depending on the level of real wage only.

2.2 Firms

The firm sector comprises four types of firms; 1) final goods producers, 2) intermediate goods producers, 3) capital producers, and 4) labor packers. Final goods producers transform intermediate goods into final consumption goods. Intermediate goods producers produce differentiated goods using capital and labor service as inputs. Capital producers transform final goods into new capital stock, subject to adjustment frictions, and rent out capital to intermediate goods producers. Labor packers combine differentiated labor supplied by unions and rent out homogeneous labor services to intermediate goods producers. Intermediate goods producers and unions operate under a monopolistically competitive environment and set prices subject to nominal rigidity á la Calvo (1983).

2.2.1 Final goods producers

Final goods producers combine differentiated intermediate goods and make final consumption goods according to a CES aggregation technology

$$Y_t = \left(\int y_{jt}^{\frac{\eta_t-1}{\eta_t}} dj \right)^{\frac{\eta_t}{\eta_t-1}}, \quad (11)$$

where y_{jt} is intermediate good j and η_t is the time-varying elasticity of substitution. Profit maximization yields the following individual good demand and the aggregate price index

$$y_{jt} = \left(\frac{p_{jt}}{P_t} \right)^{-\eta_t} Y_t \quad (12)$$

$$P_t = \int p_{jt}^{1-\eta_t} dj \quad (13)$$

where p_{jt} is individual good j 's price.

2.3 Intermediate goods producers

There is a continuum of intermediate good firms, indexed by j , that produce differentiated goods, using capital and labor services, according to a constant return-to-scale production functions

$$y_{j,t} = Z_t N_{jt}^\alpha (u_{jt} K_{jt})^{1-\alpha}, \quad (14)$$

where α is the labor share in production, Z_t is total factor productivity that follows a AR(1) process in logs, N_{jt} is labor input, and $u_{jt} K_{jt}$ is capital input with the utilization rate u_{jt} . Capital depreciation rate

depends on the degree of utilization, according to $\delta(u_{jt}) = \delta_0 + \delta_1(u_{jt} - 1) + \delta_2/2(u_{jt} - 1)^2$. First order conditions associated with the cost minimization are given as follows.

$$w_t^F = \alpha \text{mc}_{jt} Z_t \left(\frac{u_{jt} K_{jt}}{N_{jt}} \right)^{1-\alpha} \quad (15)$$

$$r_t + q_t \delta(u_{jt}) = u_{jt} (1 - \alpha) \text{mc}_{jt} Z_t \left(\frac{N_{jt}}{u_{jt} K_{jt}} \right)^\alpha \quad (16)$$

$$q_t [\delta_1 + \delta_2(u_{jt} - 1)] = (1 - \alpha) \text{mc}_{jt} Z_t \left(\frac{N_{jt}}{u_{jt} K_{jt}} \right)^\alpha, \quad (17)$$

where mc_{jt} is the marginal cost of production of firm j . Since the production function exhibits constant return-to-scale, the above optimality conditions imply that marginal costs are identical across producers, i.e., $\text{mc}_{jt} = \text{mc}_t$.

Firms operate under monopolistically competitive environments and set prices for their goods subject to price adjustment frictions á la Calvo (1983); only a fraction $1 - \lambda_Y$ of firms can adjust their prices, while the rest index their prices to the steady-state inflation rate $\bar{\pi}$. Thus, firms maximize the present value of real profits

$$\mathbb{E}_0 \sum_{t=0}^{\infty} \beta^t \lambda_Y^t (1 - \tau_t^L) Y_t^{1-\tau_t^P} \left\{ \left(\frac{p_{jt} \bar{\pi}^t}{P_t} - \text{mc}_t \right) \left(\frac{p_{jt} \bar{\pi}^t}{P_t} \right)^{-\eta_t} \right\}^{1-\tau_t^P}. \quad (18)$$

The corresponding optimality condition, with a first-order approximation, implies the following Phillips curve

$$\log \left(\frac{\pi_t}{\bar{\pi}} \right) = \beta \log \left(\frac{\pi_{t+1}}{\bar{\pi}} \right) + \kappa_Y \left(\text{mc}_t - \frac{1}{\mu_t^Y} \right), \quad (19)$$

where $\kappa = \frac{(1 - \lambda_Y)(1 - \lambda_Y \beta)}{\lambda_Y}$ is the slope of Phillips curve and $\mu_t^Y = \frac{\eta_t}{\eta_t - 1}$ is the target markup. The target markup follows a AR(1) process with markup shocks $\epsilon_t^{\mu^Y}$.

2.3.1 Capital producers

Capital producers transform final goods into new capital stock, subject to adjustment frictions, while taking the price of capital q_t as given. That is, they maximize

$$\mathbb{E}_0 \sum_{t=0}^{\infty} \beta^t I_t \left\{ q_t \Psi_t \left[1 - \frac{\phi}{2} \left(\log \frac{I_t}{I_{t-1}} \right)^2 \right] - 1 \right\}, \quad (20)$$

where ϕ governs the degree of investment adjustment frictions and Ψ_t represents marginal efficiency of investment á la Justiniano et al. (2011), which follows an AR(1) process in logs with shocks ϵ_t^Ψ . With a first order approximation, the optimality condition for the maximization problem is given by

$$q_t \Psi_t \left[1 - \phi \log \frac{I_t}{I_{t-1}} \right] = 1 - \beta \mathbb{E}_t \left[q_{t+1} \Psi_{t+1} \phi \log \left(\frac{I_{t+1}}{I_t} \right) \right]. \quad (21)$$

Finally, the law of motion for aggregate capital is described as follows

$$K_t - (1 - \delta(u_t))K_{t-1} = \Psi_t \left[1 - \frac{\phi}{2} \left(\log \frac{I_t}{I_{t-1}} \right)^2 \right]. \quad (22)$$

2.3.2 Unions and labor packers

There exists a unit mass of unions, indexed by l , who purchase labor services from workers and sell a different variety of labor to labor packers. Labor packers combine a different variety of labor into homogeneous labor input according to the following CES aggregation technology

$$N_t = \left(\int \hat{n}_{lt}^{\frac{\zeta_t-1}{\zeta_t}} dj \right)^{\frac{\zeta_t}{\zeta_t-1}}, \quad (23)$$

where \hat{n}_{lt} is a variety l labor service and ζ_t is the elasticity of substitution. Labor packers' cost minimization implies the following demand for each variety l of labor services

$$\hat{n}_{lt} = \left(\frac{W_{lt}}{W_t^F} \right)^{-\zeta_t} N_t, \quad (24)$$

where W_{lt} is the nominal wage set by union l and W_t^F is the nominal wage at which labor packers sell labor input to intermediate goods producers.

Unions have monopolistic power and maximize their stream of profits by setting prices w_{lt} for their labor variety, subject to nominal rigidity á la Calvo (1983). Specifically, only $1 - \lambda_w$ fraction of unions can adjust wages, while the rest of unions index wages to the steady-state wage inflation rate. Thus, they maximize

$$\mathbb{E}_0 \sum_{t=0}^{\infty} \lambda_w^t \frac{W_t^F}{P_t} N_t \left\{ \left(\frac{W_{lt}^{\bar{\pi}_t^W}}{W_t^F} - \frac{W_t}{W_t^F} \right) \left(\frac{W_{lt}^{\bar{\pi}_t^W}}{W_t^F} \right)^{-\zeta_t} \right\}. \quad (25)$$

From the optimality condition for the maximization problem, we obtain, with a first-order approximation, the wage Phillips curve

$$\log \left(\frac{\pi_t^W}{\bar{\pi}_W} \right) = \beta \mathbb{E}_t \log \left(\frac{\pi_{t+1}^W}{\bar{\pi}_W} \right) + \kappa_\omega \left(\text{mc}_t^\omega - \frac{1}{\mu_t^\omega} \right), \quad (26)$$

where $\kappa_\omega = \frac{(1 - \lambda_w)(1 - \lambda_w \beta)}{\lambda_w}$ is the slope of Phillips curve and $\pi_t^W \equiv \frac{W_t^F}{W_{t-1}^F} = \frac{w_t^F}{w_{t-1}^F} \pi_t^Y$ is the gross wage inflation rate with w_t and w_t^F being the real wages for households and firms, respectively. $\text{mc}_t^\omega = \frac{w_t}{w_t^F}$ is the actual and $\frac{1}{\mu_t^\omega} = \frac{\zeta_t - 1}{\zeta_t}$ is the target mark-down of wages that unions pay to households relative to wages they charge to intermediate goods producers. The target mark-down follows a log AR(1) process with wage markup shocks ϵ_t^w .

2.4 Governments

The government sector consists of a fiscal and a monetary authority. The fiscal authority issues government bonds, levies taxes, and make government purchases. The issuance of government bonds is governed by the following rule

$$\frac{B_{t+1}}{B_t} = \left(\frac{B_t}{\bar{B}}\right)^{-\gamma_B} \left(\frac{\pi_t}{\bar{\pi}}\right)^{\gamma_\pi} \left(\frac{Y_t}{Y_{t-1}}\right)^{\gamma_Y} D_t, \quad D_t = D_{t-1}^{\rho_D} \epsilon_t^D, \quad (27)$$

where D_t is a shock to the government structural deficit. The parameters γ_B , γ_π , and γ_Y represents how sensitively the deficit responds to the existing debt, the evolution of the inflation rate, and the output growth, respectively. Similalry, the government also sets the average tax rate according to a rule

$$\frac{\tau_t}{\bar{\tau}} = \left(\frac{\tau_{t-1}}{\bar{\tau}}\right)^{\rho_\tau} \left(\frac{B_t}{B_{t-1}}\right)^{(1-\rho_\tau)\gamma_B} \left(\frac{Y_t}{Y_{t-1}}\right)^{(1-\rho_\tau)\gamma_Y}. \quad (28)$$

The fiscal authority ensures that the average tax rate equals the target tax rate τ_t by adjusting the level parameter τ_L

$$\tau_t = \frac{\mathbb{E}_t(w_t n_{it} h_{it} + \mathbb{1}_{h_{it}=0} \Pi_t^E) - \tau_L \mathbb{E}_t(w_t n_{it} h_{it} + \mathbb{1}_{h_{it}=0} \Pi_t^E)^{\bar{\tau}^P}}{\mathbb{E}_t(w_t n_{it} h_{it} + \mathbb{1}_{h_{it}=0} \Pi_t^E)}. \quad (29)$$

Then, the total tax revenue is $T_t = \tau_t(w_t n_{it} h_{it} + \mathbb{1}_{h_{it}=1} \Pi_t^U + \mathbb{1}_{h_{it}=0} \Pi_t^E)$ and government purchase is determined by the balanced budget constraint, i.e., $G_t = B_{t+1} + T_t - R_t^b / \pi_t B_t$.

The monetary authority determines the nominal interest rate on government bonds according to the following Taylor rule with interest rate smoothing.

$$\frac{R_{t+1}^b}{R_t^b} = \left(\frac{R_t^b}{\bar{R}^b}\right)^{\rho_R} \left(\frac{\pi_t}{\bar{\pi}}\right)^{(1-\rho_R)\theta_\pi} \left(\frac{Y_t}{Y_{t-1}}\right)^{(1-\rho_R)\theta_Y} \epsilon_t^R, \quad (30)$$

where \bar{R}^b is the steady-state nominal interest rate. The coefficients ϕ_π and ϕ_Y represents the sensitivity of the policy rate to the evolution of price and output gap, respectively. The parameres ρ_R represents the degree of interest rate smoothing.

2.5 Market clearing conditions

The model has four markets; goods, labor, liquid and illiquid asset markets. Regarding liquid assets, there are two kinds, government bonds and profit shares. Thus, the liquid asset market clearing condition is given by

$$B_{t+1} + q_t^\Pi = B_{t+1}^d \equiv \int \{\lambda b_a^*(b, k, h) + (1-\lambda)b_n^*(b, k, h)\} d\phi_t(b, k, h), \quad (31)$$

where b_a^* and b_n^* are the optimal liquid asset choice of adjusting and non-adjusting households with liquid asset holding b , illiquid asset holding k , and productivity level h , respectively, and ϕ_t is the distribution of households over the idiosyncratic state space. The left-hand and right-hand side of the above equation

represent the aggregate liquid asset supply and demand, respectively. Similarly, the illiquid asset, i.e., capital, market clearing condition is given by

$$K_{t+1} = K_{t+1}^d \equiv \int \{\lambda k_a^*(b, k, h) + (1 - \lambda)k_n^*(b, k, h)\} d\phi_t(b, k, h), \quad (32)$$

with k_a^* and k_n^* are being optimal capital holding of adjusting and non-adjusting households with liquid asset holding b , illiquid asset holding k , and productivity level h .

The labor market clears at the wage w_t^F and w_t once the following condition equation holds.

$$\int \hat{n}_{lt} dl = D_t^w N_t = \int h_{it} n_{it} d\phi_t(b, k, h), \quad (33)$$

where $D_t^w = \int \left(\frac{W_{lt}}{W_t^F}\right)^{-\zeta_t} dl$ is the dispersion of wages set by unions and N_t is the aggregate labor input. The first two items of the above equation represent the demand of intermediate goods producers for a variety of labor, while the last item is the aggregate labor variety supplied by households. Once assets and labor markets clear, goods market also clears because of Walras' law.

2.6 Numerical method

Following Reiter (2009), we solve the model using a linearized solution technique. First step is to write the equilibrium as a system of non-linear difference equations as follows

$$\mathbb{E}_t F(X_t^*, X_{t+1}^*) = 0, \quad (34)$$

where X_t^* is a vector of state and control variables in period t . Then, we linearize the above system around the non-stochastic steady-state and apply a standard perturbation method, such as ones proposed by Klein (2000). However, without any further treatments, applying a standard perturbation method is infeasible since the size of the above system is very large due to many idiosyncratic state variables such as asset holdings, productivity levels, and working statuses. Thus, we follow Bayer and Luetticke (2020) and reduce the size of the two biggest components of the system, value functions and the households distribution.

For the value functions, BBL use a discrete cosine transformation (DCT), as proposed by Bayer and Luetticke (2020). All the value functions are written as linear interpolants based on a set of nodal values and these nodal values are represented by DCT coefficients of Chebyshev polynomials as follows

$$\hat{\mathbb{W}}_{b/k,t}(b_i, k_j, h_l) = \sum_{p,q,r} \theta_{\mathbb{W}_{b/k,t}}^{p,q,r} T_p(i) T_q(j) T_r(l), \quad (35)$$

where $\hat{\mathbb{W}}_{b/k}$ is the partial derivative of the continuation value \mathbb{W} with respect to bond b , (capital k) holdings, $T_{p/q/r}(\cdot)$ are Chebyshev polynomials, and $\theta_{\mathbb{W}}^{p,q,r}$ are the corresponding DCT coefficients. In the above expression, BBL force very small coefficients to be zero for the size reduction. Then, they only keep a small

number of these coefficients, instead of the original value functions, required to recover the original functions with a certain threshold level of precision. In perturbing the system, they perturb these coefficients instead of the function values themselves.

BBL reduce the size of the distribution in a similar way. For the distribution, they only keep marginal distributions F_t^b , F_t^k , and F_t^h in the system and use a copula $C_t(\cdot)$, a functional relationship between marginals and the joint distribution, to recover the joint distribution from marginals. Then, the copula $C_t(\cdot)$ at time t is approximatged using Chebyshev polynomials

$$\hat{C}_t(F_i^b, F_j^k, F_l^h) = \sum_{p,q,r} \theta_{C,t}^{p,q,r} T_p(i) T_q(j) T_r(l), \quad (36)$$

where $\hat{C}(\cdot)$ is the deviation of the copula at time t from its steady-state counterpart. Again, BBL reduce the size of the system by keeping only a small number of DCT coefficients $\theta_C^{p,q,r}$.

After the state space reduction, the dimension of the system decreases substantially and one can find a linerized solution rather quickly. However, for the purpose of estimation, further acceleration of the solution method is required since one needs to efficiently evaluate the model's likelihood. To this end, we also follow Bayer and Luetticke (2020) and only estimate a subset of parameters that do not affect the households' problem. BBL first partion X^* into the part related to households choices f and the aggregate part X

$$\mathbb{E}_t F(X_t^*, X_{t+1}^*) = \mathbb{E}_t F(f_t, X_t, f_{t+1}, X_{t+1}). \quad (37)$$

Then, they obtain the following linerized system

$$\begin{bmatrix} A_{ff} & A_{fX} \\ A_{Xf} & A_{XX} \end{bmatrix} \begin{bmatrix} f_t \\ X_t \end{bmatrix} = -\mathbb{E}_t \begin{bmatrix} B_{ff} & B_{fX} \\ B_{Xf} & B_{XX} \end{bmatrix} \begin{bmatrix} f_{t+1} \\ X_{t+1} \end{bmatrix}. \quad (38)$$

In the above system, one only needs to update A_{XX} and B_{XX} during the estimation if only the parameters that do not affect household problem are estimated. Since the size of aggregate blocks A_{XX} and B_{XX} is relatively small, one can update the Jacobian rather quickly.

Finally, BBL perform a further model reduction, which relies on a factor representation of the idiosyncratic model part, i.e., the part related to household choices. Once they define objects in a way such that $B_{fX} = B_{Xf} = 0$, they reduce the size of the system by applying a singular value decomposition (SVD) on the idiosyncratic model part. Specifically, they rewrite the linearized system as

$$\begin{bmatrix} B_{ff}^{-1} A_{ff} & B_{ff}^{-1} A_{fX} \\ A_{Xf} & A_{XX} \end{bmatrix} \begin{bmatrix} f_t \\ X_t \end{bmatrix} = \begin{bmatrix} \tilde{A}_{ff} & \tilde{A}_{fX} \\ \tilde{A}_{Xf} & \tilde{A}_{XX} \end{bmatrix} \begin{bmatrix} f_t \\ X_t \end{bmatrix} = -\mathbb{E}_t \begin{bmatrix} f_{t+1} \\ B_{XX} X_{t+1} \end{bmatrix}. \quad (39)$$

Then, by applying a SVD on \tilde{A}_{ff} , i.e., $\tilde{A}_{ff} = U\Sigma V'$, and the Eckart-Young-Mirsky theorem, they obtain

$$\begin{bmatrix} V_1' U \Sigma_1 & V_1' \tilde{A}_{fX} \\ \tilde{A}_{Xf} V_1 & \tilde{A}_{XX} \end{bmatrix} \begin{bmatrix} Y_t \\ X_t \end{bmatrix} \approx -\mathbb{E}_t \begin{bmatrix} Y_{t+1} \\ B_{XX} X_{t+1} \end{bmatrix}, \quad (40)$$

where V_1 refers to the rows in V that correspond to the largest singular values and $Y_t = V'f_t$. Since \tilde{A}_{ff} is independent of the estimated parameters, the SVD needs to be performed only infrequently. With this second-stage model reduction the size of the model decreases drastically once again and QZ-decomposition can take place within a relatively short amount of time, which makes the estimation feasible.

3 Online estimation of HANK models

The goal of this section is to describe a Monte Carlo approach that makes it possible what we call “online” estimation of HANK models. By online estimation, we mean estimation that can be conducted without starting from scratch as the dataset changes because, say, a new quarter of data is available. If estimating a model from scratch is nowadays a relatively trivial computational task for (linear) medium scale DSGE models of the size of Smets and Wouters (2007), it becomes much more time consuming and computer intensive when the size of the state space becomes very large as is the case for HANK models.

Online estimation can be useful to central bank researchers who would like to use HANK models for forecasting. It can also be useful for academics who intend to run pseudo out-of-sample forecasting comparisons to assess the forecasting ability of HANK models, as we do in this paper, as these comparisons involve re-estimating the model(s) for each vintage of data.⁶ Finally, online estimation can also be used to quickly re-estimate a model after small changes such as, for instance, modification of the prior, or any other relatively minor (or perhaps even major) alterations of the model.

The first part of this section describes the estimation problem and why the way it is currently handled by popular DSGE estimation packages such as `Dynare` may not be ideal for HANK models. The following subsection provides an intuitive description of an alternative estimation method—Sequential Monte Carlo (henceforth, SMC)—and explains why this approach is suitable for online estimation. While this section borrows much of the material from Cai et al. (2021), it strives to be accessible to an audience with little or no prior knowledge of Monte Carlo methods.⁷

⁶Edge and Gürkaynak (2010), and Del Negro and Schorfheide (2013), are examples of forecasting comparisons using medium scale DSGEs.

⁷A terrific introduction to such methods is provided in textbooks such as Gelman et al. (1995), Geweke (2005), and Herbst and Schorfheide (2015), where the latter focuses specifically on DSGE model estimation. We refer the reader to these textbooks for a more formal treatment of the ideas described below.

3.1 Bayesian estimation of HANK models

The solution of the log-linearized version of the model described in section 2 produces the following transition equation

$$s_t = T(\theta)s_{t-1} + R(\theta)\varepsilon_t, \quad t = 1, \dots, T, \quad (41)$$

where s_t is the vector of states, θ is the parameter vector, and the shocks ε_t are independently and identically distributed according to $\varepsilon_t \sim N(0, Q(\theta))$. The measurement equation

$$y_t = Z(\theta)s_t + D(\theta) + u_t, \quad t = 1, \dots, T, \quad (42)$$

connects the latent states s_t to the vector of observables y_t , where the measurement error shocks are independently and identically distributed according to $u_t \sim N(0, H(\theta))$. The likelihood of this linear, Gaussian state-space model $p(y_{1:T}|\theta)$ can be readily computed via the Kalman filter, where we use the notation $y_{1:T}$ to denote the sequence of observations $\{y_1, \dots, y_T\}$. Using Bayes' law, the posterior distribution of the parameters $p(\theta|y_{1:T})$ is obtained from

$$p(\theta|y_{1:T}) = \frac{p(y_{1:T}|\theta)p(\theta)}{\int p(y_{1:T}|\theta)p(\theta)d\theta}, \quad (43)$$

where $p(\theta)$ represents our prior for the parameters (Del Negro and Schorfheide, 2009, discusses the choice of priors for DSGE models and Müller, 2012, provides an easy way to assess their influence on the results). The discussion so far applies to any log-linearized DSGE model, and follows closely An and Schorfheide (2007) and Del Negro and Schorfheide (2010). The peculiarity of HANK models is that the state-space vector s_t is extremely large, making the Kalman filter and hence the evaluation of the likelihood $p(y_{1:T}|\theta)$ very costly.⁸

Since the posterior $p(\theta|y_{1:T})$ does not follow any known distribution, we need Monte Carlo methods in order to obtain draws from it and describe the results of our inference on θ (that is, tabulate the posterior mean, the 90 percent posterior coverage intervals, *et cetera*). The most standard Monte Carlo algorithm used for this purpose when estimating DSGE model, and used in *Dynare*, is the *Random-Walk Metropolis Hastings* (RWMH). This algorithm is a so-called Markov Chain algorithm, in that it produces a chain of draws from the posterior distribution $\{\theta^{(1)}, \dots, \theta^{(j)}, \dots, \theta^{(J)}\}$, and, loosely speaking, works as follows: in order to obtain the draw $\theta^{(j)}$, you take the previous draw $\theta^{(j-1)}$, add some randomness to generate a proposal θ^* , and then either accept (that is set $\theta^{(j)} = \theta^*$) or reject (that is set $\theta^{(j)} = \theta^{(j-1)}$) this proposal according to a formula that guarantees convergence of the chain to the desired ergodic distribution, that is, $p(\theta|y_{1:T})$ (again, see An and Schorfheide, 2007, or Del Negro and Schorfheide, 2010).

⁸Herbst (2015) shows how the so-called ‘‘Chandrasekhar Recursions’’ formulas can substantially reduce the computational burden of evaluating the likelihood. One issue with these formulas is that they are far less generous than standard formulas in terms of accommodating missing data, which is why we do not use them in this paper.

This is the algorithm used by almost all papers doing Bayesian estimation of DSGE models, including BBL, and for medium sized models this algorithm has been shown to work reasonably well. It has a few downsides however: 1) it is well known that RWMH may get stuck and fail to explore the entirety of the parameter space, especially in presence of multi-modality (see for instance Herbst and Schorfheide, 2014); 2) it cannot be parallelized, since it is a Markov chain; and 3) one has to start from scratch for any new estimation, even if the changes in the estimation settings are relatively minor so that one would not expect a major change in the posterior distribution (eg, adding one more quarter of data). These issues are particularly serious for HANK models because their posterior distribution is harder to evaluate. For instance, one approach to dealing with problem (1) amounts to running very long chains, which increases the chances of visiting the entirety of the parameter space. Of course this approach is less appealing when computing $p(\theta^{(j)}|y_{1:T})$ is very costly. Similarly, the fact that the algorithm cannot be parallelized limits the extent to which one can take advantage of computer power to speed up the algorithm. While recent developments in Monte Carlo methods, such as Hamiltonian Monte Carlo (eg, Duane et al., 1987; Neal et al., 2011; Stan Development Team, 2015, henceforth, HMC), have made Markov Chain methods more efficient and to some extent amenable to parallelization, problem (3)—the fact that one has to start each estimation from scratch—makes SMC methods appealing. We turn to describing these methods in the next section.

3.2 The Sequential Monte Carlo algorithm

In order to appreciate how and why Sequential Monte Carlo works, it may be useful to take a brief detour into the early history of Monte Carlo methods, and discuss an approach called Importance Sampling (see Hammersley and Morton, 1954 for an early example and the textbooks mentioned in footnote 7 for a more modern treatment). Let's say you do not know how to draw from the posterior $p(\theta|y_{1:T})$, but you can draw very efficiently from a *proposal* distribution $q(\theta)$. For example, $q(\theta)$ could be a Gaussian with mean $\hat{\theta} = \operatorname{argmax}_{\theta} p(\theta|y_{1:T})$, the peak of the posterior, and with variance proportional to the negative of the inverse of the numerical second derivative of the posterior evaluated at $\hat{\theta}$. Then you can obtain $\{\theta^{(1)}, \dots, \theta^{(j)}, \dots, \theta^{(J)}\}$ independent draws from $q(\theta)$ and assign to each of these draws a weight $W^{(j)} = w^{(j)} / (\frac{1}{J} \sum_{j=1}^J w^{(j)})$, where

$$w^{(j)} = p(y_{1:T}|\theta^{(j)})p(\theta^{(j)})/q(\theta^{(j)}) \propto p(\theta^{(j)}|y_{1:T})/q(\theta^{(j)}).$$

In other words, the idea behind Importance Sampling is to draw from $q(\theta)$ and then do a change of measure from $q(\cdot)$ to the so-called *target* distribution (the actual posterior) by reweighing these draws. Note that the denominator in (43) is irrelevant in the computation of $w^{(j)}$ since it does not depend on θ , and that the $W^{(j)}$ are in any case re-normalized to sum up to J (the choice of J , as opposed to the more conventional 1, as the normalization constant is driven by numerical reasons). Given the *swarm of particles* $\{\theta^{(j)}, W^{(j)}\}_{j=1}^J$

produced by this approach, one can then approximate any object of interest $h(\theta)$ using the Monte Carlo average

$$\bar{h}_J = \frac{1}{J} \sum_{j=1}^J W_n^{(j)} h(\theta^{(j)}),$$

where for instance $h(\theta) = \theta$ if one wants to compute the mean.

This may sound like a very reasonable approach except that the accuracy of this approximation does not just depend on J , which can be easily increased, but by the effective particle sample size

$$\widehat{ESS} = J / \left(\frac{1}{J} \sum_{j=1}^J (W^{(j)})^2 \right).$$

In other words, if $q(\theta)$ is a good proposal (in the example above, if the posterior is approximately Gaussian), then for J reasonably large Importance Sampling delivers a good approximation of the object of interest. If it is not, and \widehat{ESS} is low, then it fails. When the posterior is irregular, as is the case for many DSGEs, coming up with a good (global) approximation is nearly impossible, and this may partly explain why in DSGE estimation these methods have been abandoned in favor of Markov Chain approaches such as RWMH.⁹

SMC brings Importance Sampling and the use of swarms of particles back into play for DSGE estimation thanks to two ideas. The first idea is that if you can pick the posterior you want to approximate then the problem of choosing a suitable proposal becomes much easier. For instance, if the posterior is

$$p_n(\theta|y_{1:T}) = \frac{p(y_{1:T}|\theta)^{\phi_n} p(\theta)}{\int p(y_{1:T}|\theta)^{\phi_n} p(\theta) d\theta} \quad (44)$$

with ϕ_n being a very small number, then the prior $p(\theta)$ is likely to work pretty well as a proposal: by construction the target is almost the same as the proposal. Of course, $p_n(\theta|y_{1:T})$ constructed with ϕ_n close to zero is not what we want to approximate in the end. So we can increase ϕ_n progressively toward 1, and use the $n - 1$ swarm as a proposal for the stage n target, making sure that at each stage n the target and the proposal remain reasonably close. If the swarm of particles is still that generated by the prior, all this slicing into intermediate steps amounts to nothing: the prior is a poor proposal for the eventual posterior, and the effective sample size will likely be very low. But the second idea, which borrows from Markov Chain methods, comes to the rescue: from one stage to the other particles can travel. Just like a single particle in RWMH travels around the posterior, and naturally tends to visit regions of the parameter space where the posterior places non negligible mass, so can each of the particles in the swarm $\{\theta_n^{(j)}, W_n^{(j)}\}_{j=1}^J$. In other words, the particles can adapt as ϕ_n increases toward 1, so that in the end we have a good approximation of the posterior distribution.¹⁰

Formally, the SMC algorithm goes as follows:

⁹Importance Sampling inspired approaches have remained very popular for filtering problems however, such as the particle filter, see Fernández-Villaverde and Rubio-Ramírez, 2007.

¹⁰A little bit of history: In the statistics literature, Chopin (2002) showed how to adapt particle filtering techniques to

Algorithm 1 (SMC Algorithm).

1. **Initialization.** ($\phi_0 = 0$). Draw the initial particles from the prior: $\theta_1^{(j)} \stackrel{iid}{\sim} p(\theta)$ and $W_1^{(j)} = 1$, $j = 1, \dots, J$.

2. **Recursion.** For $n = 1, \dots, N_\phi$,

(a) **Correction.** Re-weight the particles from stage $n - 1$ by defining the incremental weights

$$\tilde{w}_n^{(j)} = p(y_{1:T} | \theta_{n-1}^{(j)})^{\phi_n - \phi_{n-1}} \quad (45)$$

and the normalized weights

$$\tilde{W}_n^{(j)} = \frac{\tilde{w}_n^{(j)} W_{n-1}^{(j)}}{\frac{1}{J} \sum_{j=1}^J \tilde{w}_n^{(j)} W_{n-1}^{(j)}}, \quad j = 1, \dots, J. \quad (46)$$

(b) **Selection (Optional).** Resample the swarm of particles $\{\theta_{n-1}^{(j)}, \tilde{W}_n^{(j)}\}_{j=1}^J$ and denote resampled particles by $\{\hat{\theta}_n^{(j)}, W_n^{(j)}\}_{j=1}^J$, where $W_n^{(j)} = 1$ for all j .

(c) **Mutation.** Propagate the particles $\{\hat{\theta}_n^{(j)}, W_n^{(j)}\}$ via N_{MH} steps of an MH algorithm with transition density $\theta_n^{(j)} \sim K_n(\theta_n | \hat{\theta}_n^{(j)}; \zeta_n)$ and stationary distribution $p_n(\theta | y_{1:T})$.¹¹

3. For $n = N_\phi$ ($\phi_{N_\phi} = 1$) the final importance sampling approximation of $\mathbb{E}_\pi[h(\theta)]$ is given by:

$$\bar{h}_{N_\phi, N} = \sum_{j=1}^J h(\theta_{N_\phi}^{(j)}) W_{N_\phi}^{(j)}. \quad (48)$$

Step 2a is the same as in Importance Sampling, where the proposal is the previous stage's posterior $p_{n-1}(\theta | y_{1:T})$ and the target is $p_n(\theta | y_{1:T})$. Step 2c is the Metropolis Hastings step where each particle is given conduct posterior inference for a static parameter vector. John Geweke played an important role popularizing these techniques in economics (eg, Durham and Geweke, 2014), and Creal (2007) was the first paper that applied SMC techniques to posterior inference in a DSGE model. Herbst and Schorfheide (2014) was quite impactful, as it showed that a properly tailored SMC algorithm delivers more reliable posterior inference for the Smets and Wouters (2007) DSGE model with loose priors and a multimodal posterior distribution than the widely used RWMH algorithm. They also provide some convergence results for an adaptive version of the algorithm building on Chopin (2004).

¹¹The transition kernel $K_n(\theta_n | \hat{\theta}_n; \zeta_n)$ needs to have the following invariance property:

$$p_n(\theta_n | y_{1:T}) = \int K_n(\theta_n | \hat{\theta}_n; \zeta_n) p_n(\hat{\theta}_n | y_{1:T}) d\hat{\theta}_n. \quad (47)$$

Thus, if $\hat{\theta}_n^{(j)}$ is a draw from the stage n posterior $p_n(\theta_n | y_{1:T})$, then so is $\theta_n^{(j)}$. The MH accept-reject probabilities insure that such property is satisfied. In our application we follow Herbst and Schorfheide (2014) and Cai et al. (2021) in our choice of $K_n(\theta_n | \hat{\theta}_n; \zeta_n)$, but developments in MH algorithms, such as HMC, can be used to make this step, and hence the whole SMC algorithm, more efficient. Farkas and Tatar (2020) is an example of a paper that combines SMC with HMC.

a chance to adapt to the new posterior. Step 2b needs some discussion. Its purpose is to make sure that if the weights of the particles in the swarm become very uneven, and effective particle sample size

$$\widehat{ESS}_n = J / \left(\frac{1}{J} \sum_{j=1}^J (\tilde{W}_n^{(j)})^2 \right)$$

falls below threshold \underline{N} , a new swarm of particles are generated from the old swarm so that all the particles have the same weight.¹² Loosely speaking, particles with relatively large weight $\tilde{W}_n^{(j)}$ —that is, that are in high posterior regions of the parameter space—are given an opportunity to “procreate” (generate a number of children that is in expected values proportional $\tilde{W}_n^{(j)}$), while particles with relatively small weight ($\tilde{W}_n^{(j)} < 1$)—that is, that are in regions of the parameter space with very little mass—are “killed” with probability $1 - \tilde{W}_n^{(j)}$.

One aspect of the algorithm we have not yet discussed is the number of stages N_ϕ as well as the schedule $\{\phi_1, \dots, \phi_{N_\phi}\}$. In estimating the Smets and Wouters (2007) model, Herbst and Schorfheide (2014) find that $N_\phi = 500$ and a schedule given by the function $\phi_n = (n/N_\phi)^{2.1}$ works well. The convexity of the schedule implies that ϕ_n increases very slowly at the beginning, and faster at the end. Of course, it is far from obvious that whatever schedule works well for the Smets and Wouters (2007) model also works well for a HANK or any other DSGE model. In this respect, Cai et al. (2021) improve upon Herbst and Schorfheide (2014) by making the schedule $\{\phi_1, \dots, \phi_{N_\phi}\}$ adaptive—that is, endogenous to the difficulty of the problem. Recall that the ESS measures, loosely speaking, the deterioration of the quality of the swarm $\{\theta_n^{(j)}, W_n^{(j)}\}_{j=1}^J$: if ESS is low, the swarm has essentially “lost” most of its particles as the weights have become very uneven. Adaptation is then achieved by setting at each stage $\phi_n = \phi$ where ϕ solves

$$\widehat{ESS}(\phi) - (1 - \alpha)\widehat{ESS}_{n-1} = 0,$$

where

$$\tilde{w}^{(j)}(\phi) = [p(y_{1:T}|\theta_{n-1}^{(j)})]^{\phi - \phi_{n-1}}, \quad \tilde{W}^{(j)}(\phi) = \frac{\tilde{w}^{(j)}(\phi)W_{n-1}^{(j)}}{\frac{1}{J} \sum_{j=1}^J \tilde{w}^{(j)}(\phi)W_{n-1}^{(j)}}, \quad \widehat{ESS}(\phi) = N / \left(\frac{1}{J} \sum_{j=1}^J (\tilde{W}_n^{(j)}(\phi))^2 \right).$$

The above formulas can be understood as follows: Pick a desired deterioration α of the effective sample size between stages $n - 1$ and n , and set ϕ_n so as to achieve exactly such deterioration (see Cai et al., 2021, for a more detailed description). The parameter α expresses the degree of “carefulness” of the researchers, bearing in mind that lower α ’s imply a longer estimation time (in light of the results in Cai et al., 2021, we choose $\alpha = 5$ percent in this application). Once α is chosen, the schedule becomes endogenous to the difficulty of the problem, as measured by the deterioration of the ESS.

¹²There are many resampling schemes (see Liu, 2001, or Cappé et al., 2005) We use systematic resampling in the applications below.

This section concludes with a description of the some of virtues of SMC. First, for a suitably large choice of the size of the swarm J , it is robust to irregular shapes of the posterior such as multi-modality, as shown in Herbst and Schorfheide (2014) and Cai et al. (2021) among others. This is because the initial swarm $\{\theta_0^{(j)}, W_0^{(j)}\}_{j=1}^J$ is drawn from the prior, and hence covers for large enough J any region of the parameter space where the prior places non negligible mass. Hence if the posterior has many modes, there ought to be some initial particles in the neighborhood of such modes. Second, most of the SMC steps, such as the computation of the incremental weights in Step 2a and most important the mutation step in Step 2c, can be parallelized. Third, the algorithm produces an approximation of the marginal likelihood as a by-product. In fact, using the definitions of $\tilde{w}_n^{(j)}$ and $\tilde{W}_{n-1}^{(j)}$ one can see that:

$$\frac{1}{J} \sum_{i=1}^N \tilde{w}_n^{(j)} \tilde{W}_{n-1}^{(j)} \approx \int \frac{p(y_{1:T}|\theta)^{\phi_n}}{p(y_{1:T}|\theta)^{\phi_{n-1}}} \left[\frac{p(y_{1:T}|\theta)^{\phi_{n-1}}}{\int p(y_{1:T}|\theta)^{\phi_{n-1}} p(\theta) d\theta} \right] d\theta = \frac{\int p(y_{1:T}|\theta)^{\phi_n} p(\theta) d\theta}{\int p(y_{1:T}|\theta)^{\phi_{n-1}} p(\theta) d\theta}. \quad (49)$$

This implies that the product $\prod_{n=1}^{N_\phi} \left(\frac{1}{J} \sum_{j=1}^J \tilde{w}_n^{(j)} W_{n-1}^{(j)} \right)$ approximates the marginal likelihood as long as the prior is proper ($\int p(\theta) d\theta = 1$), since all the terms cancel out except for $\int p(y_{1:T}|\theta) p(\theta) d\theta / \int p(\theta) d\theta$. Fourth, and perhaps most important for this application, the final swarm of particles $\{\theta_{N_\phi}^{(j)}, W_{N_\phi}^{(j)}\}_{j=1}^J$ can be reused, making recursive estimation of the model very convenient. We are going to turn to this feature next.

3.3 Online estimation

Imagine you have run the SMC algorithm 1 and have a swarm of particles $\{\theta^{(j)}, W^{(j)}\}_{j=1}^J$ that approximates well the posterior $p(\theta|y_{1:T})$. Expression (45) in Step 2a of the algorithm can be generalized as

$$\tilde{w}_n^{(j)} = \frac{p_n(Y|\theta_{n-1}^{(j)})}{p_{n-1}(Y|\theta_{n-1}^{(j)})}, \quad (50)$$

where we now use the more generic notation Y for $y_{1:T}$ for reasons that will soon become apparent. Note that in (45) we considered the special case where the stage- n likelihood $p_n(Y|\theta) = p(Y|\theta)^{\phi_n}$.

Imagine that you now want to obtain the posterior for a different model $\tilde{p}(\cdot|\theta)$ (but with the same parameter vector θ) estimated on a different dataset \tilde{Y} :

$$\tilde{p}(\theta|\tilde{Y}) = \frac{\tilde{p}(\tilde{Y}|\theta)p(\theta)}{\int \tilde{p}(\tilde{Y}|\theta)p(\theta)d\theta}. \quad (51)$$

The simplest possible case is the one where the model is the same ($\tilde{p}(\cdot|\theta) = p(\cdot|\theta)$), and the dataset has one more time series observation ($\tilde{Y} = y_{1:T+1}$), but the algorithm can accommodate situations where the data has been revised, or the model changed. Draws for the posterior $\tilde{p}(\theta|\tilde{Y})$ can be readily obtained from algorithm 1 after replacing expression (45) with expression (50), and using the stage- n likelihood function

$$\tilde{p}_n(\tilde{Y}|\theta) = \tilde{p}(\tilde{Y}|\theta)^{\phi_n} p(Y|\theta)^{1-\phi_n}. \quad (52)$$

In other words, we use the posterior distribution $p(\theta|Y)$ as a “bridge” to obtain the new posterior $\tilde{p}(\theta|\tilde{Y})$, as opposed to start from the prior distribution. To the extent that the differences between $\tilde{p}(\tilde{Y}|\theta)$ and $p(Y|\theta)$ are not large, the swarm from $p(\theta|Y)$ should offer a fairly good starting point for the SMC algorithm.¹³ This is the approach we use to estimate the BBL HANK model recursively. In particular, we start from the end-of-sample estimation $p(\theta|y_{1:T})$, and then proceed backward using formula (52) with $\tilde{Y} = y_{1:T-\tau}$ and $Y = y_{1:T-\tau+1}$, for $\tau = 1, \dots, \bar{\tau}$. We should stress that doing the online recursive estimation backward or forward—that is, starting from $p(\theta|y_{1:T-\bar{\tau}})$ and using this as bridge to obtain $p(\theta|y_{1:T-\bar{\tau}+1})$, and so on—should make no difference, as both procedures recover $p(\theta|y_{1:T-\tau})$.¹⁴

We conclude this section by highlighting some of the potentials of this approach beside the online estimation of HANK models. Mlikota and Schorfheide (2022) introduce the notion of “model tempering.” If a model is very costly to estimate from scratch, one can save a lot of time by first estimating a coarser version of the model that is much cheaper to estimate (eg, the linearized version of a non linear model), and then use that as a bridge to estimate the full model. Mlikota and Schorfheide (2022) use this approach to estimate a non linear model.

4 Results

This section describes the forecasting results and is divided into three parts. The first part describes the setup of the exercise, including the data. The second part discusses the estimates of the parameters, focusing on the differences between the original BBL results and those obtained using the SMC algorithm. The last part covers the forecasting horse race between BBL and SW.

4.1 Setup

For our exercise we use the dataset made available by BBL online at https://github.com/BASEforHANK/HANK_BusinessCycleAndInequality as of June 2022, when we begun this project. This dataset comprises the seven variables used by SW in the estimation of their model, namely the growth rate of per capita real output, consumption, investment, and wages, the logarithm of hours worked per capita, GDP deflator inflation, and the federal funds rate (during the zero lower bound period the authors use the shadow rate measure created by Wu and Xia, 2016). In their database, these variables are available at the quarterly frequency from 1954Q3 to 2019Q4. In addition, BBL estimate their model adding four variables that reflect

¹³The initialization step in algorithm 1 needs to be modified so that the swarm $\{\theta^{(j)}, W^{(j)}\}_{j=1}^J$, possibly after a selection step 2b so that all the $W^{(j)}$'s equal 1, replaces the swarm drawn from the prior.

¹⁴In particular there is no sense in which the backward procedure introduces any hindsight bias: by the time that ϕ_n in (52) reaches 1, the posterior draws no longer condition on $Y = y_{1:T-\tau+1}$.

various aspects inequality and are not used in standard representative agent DSGE estimation. These are the wealth and income shares of the top 10 percent, estimates of tax progressivity constructed following Ferriere and Navarro (2018), and estimates of idiosyncratic income risk from Bayer et al. (2019). The top 10 percent shares are available annually from 1954 to 2019, the tax progressivity measure is available annually from 1954 to 2017, while the idiosyncratic income risk measure is available from 1983Q1 to 2013Q1.¹⁵ The likelihood computation of the state space model easily accommodates missing data.

BBL demean all the time series prior to estimation. While this is not standard practice in the DSGE estimation and forecasting literature and in central banks’ practice (eg, see Del Negro and Schorfheide, 2013; Cai et al., 2019), we follow BBL because adding a constant would imply introducing steady state growth and inflation, and therefore altering their model. We chose not to do this in order to remain as close as possible to BBL’s specification. This choice has two implications. First, we have to use their dataset also for the forecasting exercise—that is, the demeaned data is what the model’s forecasts are evaluated upon. Second, we estimate the competitor in the horse race—the SW model—also on demeaned data, which implies that we drop the constant from SW’s measurement equations.¹⁶

The out-of-sample forecasting exercise begins in the first quarter of 2000 (in the notation of section 3, $T - \bar{\tau} = 2000Q1$) and ends in the last quarter of 2019 ($T = 2019Q4$), for a total of 80 periods. In order to avoid hindsight bias, for each period we re-estimate the model using only data available up to that period.¹⁷ For each model \mathcal{M}_m under consideration (BBL, SW), we then generate horizon- h mean forecasts $\mathbb{E}[y_{T-\tau+h}|y_{1:T-\tau}, \mathcal{M}_m]$ for the variables of interest using the state space model consisting of equations (41) and (42), and compare these forecasts with actual outcomes $y_{T-\tau+h}$.¹⁸ As discussed above in section 3, we estimate the model in period $T - \tau$ using the posterior distribution for $T - \tau + 1$ as a bridge in the SMC algorithm. For the sake of robustness we start this process from two different posterior distributions $p(\theta|y_{1:T})$. One distribution consists of the draws made available by BBL on GitHub (which is the reason we do the online estimation backward, since these draws are only available for $T = 2019Q4$), and the other is based on an SMC estimation starting from the prior.¹⁹ The next section discusses these posterior estimates

¹⁵See Bayer et al. (2022) for a more detailed description of the dataset.

¹⁶We should note that BBL have a representative agent version of their HANK model, which they show fits the seven macro variable worse than the heterogeneous agents version in terms of marginal likelihood. We see no point in using this model as the alternative in the horse race given that the actual SW model is available. From the perspective of a central bank choosing whether to use a representative agent or a HANK model for predictions, arguably the choice is between SW and BBL.

¹⁷In the forecasting literature it is customary to perform so-called pseudo real time forecasting, where the data vintage available the time $T - \tau$ is used for estimation, as opposed to the revised data (here, the T vintage). The demeaning of the data, and the fact that there are no vintages for the inequality series, makes this pseudo real time exercise not possible.

¹⁸In order to compute the expectation $\mathbb{E}[y_{T-\tau+h}|y_{1:T-\tau}, \mathcal{M}_m]$ using (41) and (42), only the filtered states $s_{T-\tau|T-\tau} = \mathbb{E}[s_{T-\tau}|y_{1:T-\tau}, \mathcal{M}_m]$ are needed, which are obtained from the Kalman filter.

¹⁹For all SMC estimations we use a swarm of $J = 1000$ particles. While $J = 1000$ is not a large number (in Cai et al., 2021, we use $J = 12000$ for the SW model), this is as much as the computer cluster we have at our disposal could handle. However

in some detail.

4.2 Estimation

This section presents the results from our estimation. Specifically, we discuss the prior distribution, which coincides with BBL, and posteriors from the full sample SMC estimation using the eleven variables described in the previous section. Also, we show the posteriors from SMC estimation when using only seven variables, excluding data on inequality, tax progressivity estimates, and income risk estimates. For comparison, we present BBL’s estimation results obtained from BBL’s [GitHub](#) page. As mentioned, BBL made a few changes to their model and calibrated parameters since June 2022. Hence these MH estimates do not replicate the results presented in the most recent version of the paper. Table 1 summarizes the calibration.

Table 1: Calibration

Par.	Value	Description
Households: Income process		
ρ_h	0.980	Persistence labor productivity
σ_h	0.120	Std. dev. labor productivity
ι	0.063	Trans. prod. from entrepreneurs to workers
ζ	1/3750	Trans. prod. from workers to entrepreneurs
Households: Financial frictions		
λ	0.095	Portfolio adj. prob.
\bar{R}	0.017	Borrowing premium
Households: Preferences		
β	0.984	Discount factor
ξ	4.000	Relative risk aversion
γ	2.000	Inverse of Frisch elasticity
Firms		
α	0.682	Share of labor
δ_0	0.022	Depreciation rate
$\bar{\eta}$	11.000	Elasticity of substitution
$\bar{\zeta}$	11.000	Elasticity of substitution
Government		
$\bar{\tau}^L$	0.175	Tax rate level
$\bar{\tau}^P$	0.12	Tax progressivity
\bar{R}^b	1	Gross nominal rate
$\bar{\pi}$	1	Steady-state inflation rate

we have performed for T an estimation for $J = 10000$ and found that it produces very similar forecasts.

Table 2: Prior and posterior distributions: policies and frictions (2000Q1 estimation)

Par	Dist	Prior		Posterior		
		Mean	Std. Dev	BBL (MH)	Backward from MH	Backward from SMC
Monetary policy						
ρ_R	Beta	0.50	0.20	0.785 (0.754, 0.814)	0.733 (0.700, 0.766)	0.796 (0.761, 0.829)
σ_R	Inv-Gamma	0.10	2.00	0.243 (0.224, 0.269)	0.303 (0.271, 0.329)	0.242 (0.220, 0.263)
θ_π	Normal	1.70	0.30	2.237 (2.044, 2.424)	2.1911 (2.025, 2.395)	1.535 (1.355, 1.718)
θ_Y	Normal	0.13	0.05	0.287 (0.223, 0.361)	0.2810 (0.2120, 0.3422)	0.170 (0.110, 0.238)
Fiscal policy: deficit						
ρ_D	Beta	0.50	0.20	0.965 (0.950, 0.980)	0.986 (0.977, 0.998)	0.813 (0.772, 0.846)
σ_D	Inv-Gamma	0.10	2.00	0.310 (0.277, 0.342)	0.318 (0.278, 0.354)	0.717 (0.601, 0.853)
γ_B	Gamma	0.10	0.08	0.031 (0.008, 0.047)	0.048 (0.016, 0.073)	0.039 (0.016, 0.060)
γ_π	Normal	0.00	1.00	-1.601 (-1.778, -1.452)	-1.530 (-1.674, -1.375)	-3.010 (-3.400, -2.531)
γ_Y	Normal	0.00	1.00	-0.350 (-0.418, -0.309)	-0.232 (-0.292, -0.165)	-0.788 (-0.939, -0.672)
Fiscal policy: taxes						
ρ_τ	Beta	0.50	0.20	0.653 (0.440, 0.961)	0.507 (0.348, 0.680)	0.722 (0.598, 0.835)
γ_B^τ	Normal	0.00	1.00	0.166 (0.110, 0.217)	0.176 (0.107, 0.234)	-1.589 (-1.958, -1.181)
γ_Y^τ	Normal	0.00	1.00	-0.148 (-0.410, 0.038)	-2.297 (-3.144, -1.229)	0.262 (-0.894, 1.800)
Income risk						
ρ_s	Beta	0.50	0.20	0.663 (0.606, 0.727)	0.647 (0.539, 0.760)	0.9818 (0.970, 0.991)
σ_s	Gamma	65.00	30.00	64.08 (55.91, 71.06)	53.25 (44.14, 61.34)	50.02 (43.96, 56.89)
Frictions						
δ_s	Gamma	5.00	2.00	0.456 (0.278, 0.631)	1.298 (1.031, 1.494)	1.9130 (1.407, 2.452)
ϕ	Gamma	4.00	2.00	0.787 (0.373, 1.244)	0.355 (0.336, 0.363)	3.759 (3.092, 4.582)
κ_p	Gamma	0.10	0.03	0.111 (0.094, 0.125)	0.130 (0.116, 0.145)	0.129 (0.112, 0.143)
κ_w	Gamma	0.10	0.03	0.112 (0.095, 0.128)	0.100 (0.084, 0.113)	0.099 (0.084, 0.114)

Notes: The table shows prior and posterior distributions of the estimated parameters. 10 and 90 percentile of the distributions are in parenthesis.

Table 3: Prior and posterior distributions: structural shocks (2000Q1 estimation)

Par	Dist	Prior		Posterior		
		Mean	Std. Dev	BBL (MH)	Backward from MH	Backward from SMC
Structural shocks						
ρ_A	Beta	0.50	0.20	0.954 (0.925, 0.976)	0.951 (0.911, 0.989)	0.976 (0.970, 0.985)
σ_A	Inv-Gamma	0.10	2.00	0.162 (0.133, 0.194)	0.170 (0.129, 0.203)	0.074 (0.054, 0.091)
ρ_Z	Beta	0.50	0.20	0.998 (0.996, 0.999)	0.997 (0.994, 0.999)	0.966 (0.949, 0.980)
σ_Z	Inv-Gamma	0.10	2.00	0.569 (0.526, 0.624)	0.610 (0.542, 0.665)	0.640 (0.589, 0.692)
ρ_Ψ	Beta	0.50	0.20	0.848 (0.790, 0.904)	0.933 (0.918, 0.951)	0.390 (0.308, 0.470)
σ_Ψ	Inv-Gamma	0.10	2.00	3.814 (2.820, 4.982)	3.154 (2.817, 3.436)	15.79 (12.88, 18.69)
ρ_{μ_p}	Beta	0.50	0.20	0.862 (0.824, 0.907)	0.822 (0.770, 0.876)	0.879 (0.814, 0.936)
σ_{μ_p}	Inv-Gamma	0.10	2.00	1.563 (1.404, 1.714)	1.357 (1.183, 1.520)	1.254 (1.091, 1.395)
ρ_{μ_w}	Beta	0.50	0.20	0.862 (0.826, 0.907)	0.939 (0.920, 0.961)	0.924 (0.898, 0.951)
σ_{μ_w}	Inv-Gamma	0.10	2.00	6.142 (5.385, 6.916)	4.443 (3.922, 5.001)	4.168 (3.576, 4.630)
ρ_P	Beta	0.50	0.20	0.961 (0.943, 0.981)	0.957 (0.936, 0.978)	0.963 (0.944, 0.981)
σ_P	Inv-Gamma	0.10	2.00	3.534 (2.938, 4.192)	3.798 (3.133, 4.603)	3.647 (2.971, 4.173)

Notes: The table shows prior and posterior distributions of the estimated parameters. 10 and 90 percentile of the distributions are in parenthesis. Standard deviations are multiplied by 100 for readability.

Priors

For priors, they assume the same distributions used in BBL. Regarding variable capital utilization, they assume a gamma distribution with a mean of 5.0 and standard deviation of 2.0 for $\delta_s = \delta_2/\delta_1$. Similarly, they impose a gamma distribution with a mean of 4.0 and standard deviation for ϕ , the parameter that governs investment adjustment costs. For the slopes of price and wage Phillips curves, κ_Y and κ_w , they adopt Gamma priors with a mean of 0.10 and standard deviation of 0.03. The prior mean for these parameters implies that the average duration of price and wage is four quarters.

For parameters related to monetary policy, they impose normal distribution with a mean of 1.7 and standard deviation of 0.3 for θ_π , while imposing normal distribution with a mean of 0.13 and standard deviation of 0.05 for θ_Y . For the interest rate smoothing parameter ρ_R , they assume a beta distribution with parameters (0.5, 0.2).

Regarding fiscal policy, the debt-feedback parameter γ_B in the bond issuance rule is assumed to follow a gamma distribution with a mean of 0.10 and standard deviation of 0.08, which implies that the prior for the autocorrelation of government debt is 0.9. For the responsiveness of government debt to inflation and output growth, γ_π and γ_Y , they impose standard normal distributions. Similarly, they assume beta distributions with a mean of 0.5 and standard deviation of 0.2 for the autoregressive parameters in the tax rules, ρ_P , and ρ_τ . The feedback parameters for average tax rates, γ_Y^τ , and γ_B^τ , are assumed to follow standard normal distributions.

For the structural shocks, they assume beta distributions with a mean of 0.5 and a standard deviation of 0.2 for the autocorrelation parameters and inverse-gamma distributions with a mean of 0.001 and a standard deviation of 0.02 for standard deviations of shocks. Finally, for idiosyncratic income risks, we follow BBL and impose a beta distribution with a mean of 0.7 and a standard deviation of 0.2 for autocorrelation parameters and a gamma distribution with a mean of 0.65 and a standard deviation of 0.03.

Posteriors

Posterior distributions from different estimations are displayed in Table 4 and 5. Column 5 shows BBL's original posteriors that they obtained by the RWMH algorithm (MH estimation hereinafter). Columns 6 and 7 show the posterior distributions we obtained from SMC estimation using eleven and seven variables (11 and 7 var SMC estimations hereafter), respectively. In 7 var SMC estimation, we follow BBL and shut down income risk and tax progressivity shocks as we do not use the related data in the estimation.

First, the posteriors from 11 and 7 var SMC estimations exhibit only small differences, which is consistent with the findings of BBL and Chang et al. (2021). Adding data on inequality to the estimation does not significantly affect the estimation results of parameters that govern the aggregate dynamics of the model.

Table 4: Prior and posterior distributions: policies and frictions

Par	Dist	Prior		Posterior		
		Mean	Std. Dev	BBL (MH)	BBL (SMC)	BBL (7 Var)
Monetary policy						
ρ_R	Beta	0.50	0.20	0.785 (0.754, 0.814)	0.825 (0.803, 0.848)	0.841 (0.823, 0.862)
σ_R	Inv-Gamma	0.10	2.00	0.243 (0.224, 0.269)	0.215 (0.201, 0.232)	0.209 (0.193, 0.222)
θ_π	Normal	1.70	0.30	2.237 (2.044, 2.424)	1.635 (1.475, 1.760)	1.593 (1.399, 1.795)
θ_Y	Normal	0.13	0.05	0.287 (0.223, 0.361)	0.176 (0.115, 0.247)	0.185 (0.116, 0.256)
Fiscal policy: deficit						
ρ_D	Beta	0.50	0.20	0.965 (0.950, 0.980)	0.802 (0.775, 0.829)	0.7745 (0.743, 0.807)
σ_D	Inv-Gamma	0.10	2.00	0.310 (0.277, 0.342)	0.664 (0.576, 0.759)	0.984 (0.814, 1.143)
γ_B	Gamma	0.10	0.08	0.031 (0.008, 0.047)	0.030 (0.014, 0.045)	0.024 (0.003, 0.042)
γ_π	Normal	0.00	1.00	-1.601 (-1.778, -1.452)	-2.936 (-3.235, -2.605)	-3.939 (-4.593, -3.385)
γ_Y	Normal	0.00	1.00	-0.350 (-0.418, -0.309)	-0.770 (-0.868, -0.631)	-1.223 (-1.462, -1.071)
Fiscal policy: taxes						
ρ_τ	Beta	0.50	0.20	0.653 (0.440, 0.961)	0.788 (0.695, 0.874)	0.694 (0.515, 0.911)
γ_B^τ	Normal	0.00	1.00	0.166 (0.110, 0.217)	-1.543 (-1.842, -1.223)	-1.222 (-1.530, -0.894)
γ_Y^τ	Normal	0.00	1.00	-0.148 (-0.410, 0.038)	-0.008 (-1.196, 1.071)	1.137 (-0.172, 2.500)
Income risk						
ρ_s	Beta	0.50	0.20	0.663 (0.606, 0.727)	0.693 (0.601, 0.990)	- -
σ_s	Gamma	65.00	30.00	64.08 (55.91, 71.06)	62.49 (55.81, 70.02)	- -
Frictions						
δ_s	Gamma	5.00	2.00	0.456 (0.278, 0.631)	2.106 (1.683, 2.575)	2.343 (1.894, 2.778)
ϕ	Gamma	4.00	2.00	0.787 (0.373, 1.244)	4.024 (3.214, 4.712)	6.954 (5.727, 8.022)
κ_p	Gamma	0.10	0.03	0.111 (0.094, 0.125)	0.116 (0.106, 0.138)	0.113 (0.101, 0.128)
κ_w	Gamma	0.10	0.03	0.112 (0.095, 0.128)	0.124 (0.106, 0.138)	0.112 (0.096, 0.129)

Notes: The table shows prior and posterior distributions of the estimated parameters. 10 and 90 percentile of the distributions are in parenthesis.

Table 5: Prior and posterior distributions: structural shocks

Par	Dist	Prior		Posterior		
		Mean	Std. Dev	BBL (MH)	BBL (SMC)	BBL (7 Var)
Structural shocks						
ρ_A	Beta	0.50	0.20	0.954 (0.925, 0.976)	0.974 (0.965, 0.983)	0.985 (0.976, 0.992)
σ_A	Inv-Gamma	0.10	2.00	0.162 (0.133, 0.194)	0.082 (0.055, 0.103)	0.062 (0.039, 0.080)
ρ_Z	Beta	0.50	0.20	0.998 (0.996, 0.999)	0.970 (0.961, 0.981)	0.973 (0.965, 0.984)
σ_Z	Inv-Gamma	0.10	2.00	0.569 (0.526, 0.624)	0.620 (0.577, 0.664)	0.611 (0.566, 0.654)
ρ_Ψ	Beta	0.50	0.20	0.848 (0.790, 0.904)	0.497 (0.433, 0.560)	0.418 (0.343, 0.504)
σ_Ψ	Inv-Gamma	0.10	2.00	3.814 (2.820, 4.982)	14.04 (11.16, 16.44)	24.32 (20.55, 27.86)
ρ_{μ_p}	Beta	0.50	0.20	0.862 (0.824, 0.907)	0.965 (0.953, 0.980)	0.968 (0.955, 0.981)
σ_{μ_p}	Inv-Gamma	0.10	2.00	1.563 (1.404, 1.714)	1.455 (1.327, 1.583)	1.459 (1.319, 1.591)
ρ_{μ_w}	Beta	0.50	0.20	0.862 (0.826, 0.907)	0.898 (0.862, 0.925)	0.877 (0.846, 0.910)
σ_{μ_w}	Inv-Gamma	0.10	2.00	6.142 (5.385, 6.916)	4.707 (3.921, 5.195)	5.114 (4.505, 5.807)
ρ_P	Beta	0.50	0.20	0.961 (0.943, 0.981)	0.966 (0.947, 0.981)	-
σ_P	Inv-Gamma	0.10	2.00	3.534 (2.938, 4.192)	3.570 (2.959, 4.094)	-

Notes: The table shows prior and posterior distributions of the estimated parameters. 10 and 90 percentile of the distributions are in parenthesis. Standard deviations are multiplied by 100 for readability.

The investment adjustment cost is estimated to be higher in 7 var SMC estimation, but overall, the estimated posteriors between the two versions are close to each other.

Posteriors from MH and SMC estimations are also broadly similar, but they exhibit notable differences for some parameters. Regarding monetary policy, posteriors from SMC estimations imply a slightly higher interest rate inertia and lower sensitivities of the interest rate with respect to the inflation rate and output growth relative to those from MH estimation. The interest rate smoothing parameters are 0.83 and 0.84 at the mean in SMC estimations, while the corresponding mean is 0.79 in the posterior from MH estimation. The estimated Taylor rule coefficients on inflation are relatively low in SMC estimations, with the 10 to 90 percentile range being from around 1.4 to 1.8. In contrast, the corresponding range is from 2.0 to 2.4 in the posterior from MH estimation. Similarly, the coefficients on output growth are around 0.18 at the mean in the posteriors from SMC estimations, but the corresponding mean is a bit higher at 0.29 in MH estimation.

In the case of fiscal policy parameters, differences between MH and SMC estimation are more pronounced. Regarding the bond issuance rule, SMC estimations imply much less persistence of structural deficit with the autoregressive coefficient of around 0.8 at the posterior mode, while the corresponding mean is 0.97 in MH estimation. Also, posteriors from SMC estimations imply a much stronger counter-cyclical response of government debt to the inflation rate and output growth. The elasticities of the bond issuance with respect to inflation and output growth are -2.94 and -0.77 at the posterior mean in 11 var SMC estimation and -3.94 and -1.22 in 7 var SMC estimation. Posteriors for parameters governing tax rules show bigger differences. While posteriors from SMC estimations imply counter-cyclical tax rate responses with respect to the growth rate of government debt, the posterior from MH estimation implies pro-cyclical responses. Also, tax rates are estimated to be more persistent in SMC estimations than in MH estimation.

Among parameters governing model frictions, the posterior distributions of variable capital depreciation and investment adjustment cost parameters show significant differences. In MH estimation, the capital depreciation parameter is 0.46 at the posterior mode. In contrast, the modes for the same parameter are 2.11 and 2.34 in 11 and 7 var SMC estimations, respectively. Similarly, the posterior modes for the capital adjustment cost parameter are 4.02 and 6.95 in SMC estimations, while the mean for the same parameter is 0.79 in MH estimation.

The posterior distributions for the rest of the parameters, including income risk parameters, slopes of price and wage Phillips curve, and structural shocks, are broadly similar, except for the autocorrelation of the MEI shock. In SMC estimations, MEI shocks are estimated to be much less persistent, with autocorrelation of around 0.5 at the modes, than in MH estimation, where the posterior mean for the same parameter is 0.85.

The differences between the MH and the SMC estimation results obviously lead to the question as to which method is the most accurate. This is a non trivial question to address, since doing so would involve

repeated (independent) estimations of the HANK model as done for instance in Cai et al. (2021). This is computationally very costly. We therefore sidestep this issue entirely and use both in our forecasting comparison exercise. It turns out that the accuracy of the BBL model estimated using SMC is better than that using the MH draws. While this evidence is no proof that the SMC estimation is more reliable, it seems to point in that direction.

Finally, in Table 2 and 3, we present the posterior distributions from the estimations using the data up to 2000Q1, which we obtain using the approach described in section 3. Columns 6 and 7 show the posterior distributions from the backward estimation starting from the 2019Q4 MH and SMC draws, respectively. For comparison, we also show the posteriors from the original 2019Q4 BBL’s estimation in column 4. The 2000Q1 posterior is close to the MH BBL posterior when the original MH draws were used as a starting point. Similarly, when using the full sample SMC estimation result as a starting point, the 2000Q1 results are close to the estimates obtained for the 2019Q4 SMC estimation.

4.3 Assessing HANK’s out-of-sample forecasting accuracy

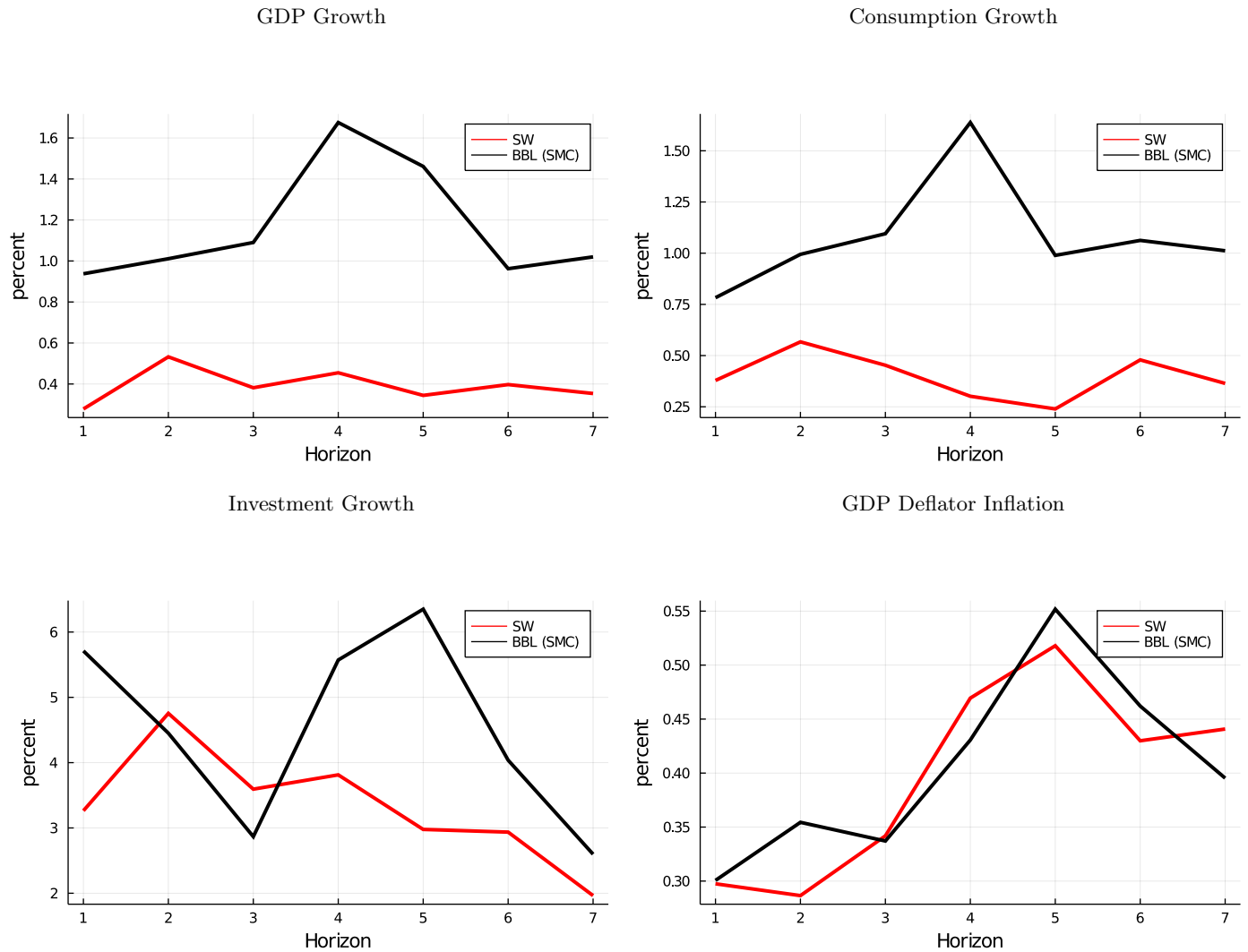
Figure 1 shows the results of the horse race between BBL and SW focusing on four variables of interest: output, consumption, and investment growth, and the GDP deflator inflation. For each of these variables the figure displays the root mean square errors (RMSEs), expressed in percent, computed as

$$RMSE_{i,h,\mathcal{M}_m} = \frac{1}{\bar{\tau} - h + 1} \sum_{\tau=h}^{\bar{\tau}} (y_{i,T-\tau+h} - \mathbb{E}[y_{i,T-\tau+h}|y_{1:T-\tau}, \mathcal{M}_m])^2 \quad (53)$$

where i indicates the variable being forecast, h is the forecast horizon, which ranges from 1 to 7 quarters ahead, and \mathcal{M}_m is the model. The model set is $\mathcal{M} = \{\text{BBL}, \text{SW}\}$, with the BBL RMSEs shown by the solid black line and the SW RMSEs by the solid red line. The BBL model is referred to as BBL(SMC) as it uses the posterior computed from the online estimation starting from the SMC draws, as opposed to the original MH draws from BBL, as the SMC estimation performs better than the MH one, as shown later.

For the variables measuring real activity, in particular output and consumption growth, the results of the horse race are not very kind to the BBL model. This is especially true for consumption growth, where the RMSEs are roughly between about two (for both short and longer horizons) and six ($h = 4$) times larger for BBL. The differences in forecasting performance for consumption growth largely translate into similar differences in RMSEs for GDP growth, given that consumption represents its largest component. For investment growth, the forecasting accuracy of the two models is similar for shorter horizons, but is again worse for BBL for medium horizons. One piece of good news for the BBL model is that for the GDP deflator inflation its RMSEs are comparable to those of SW for all forecast horizons.

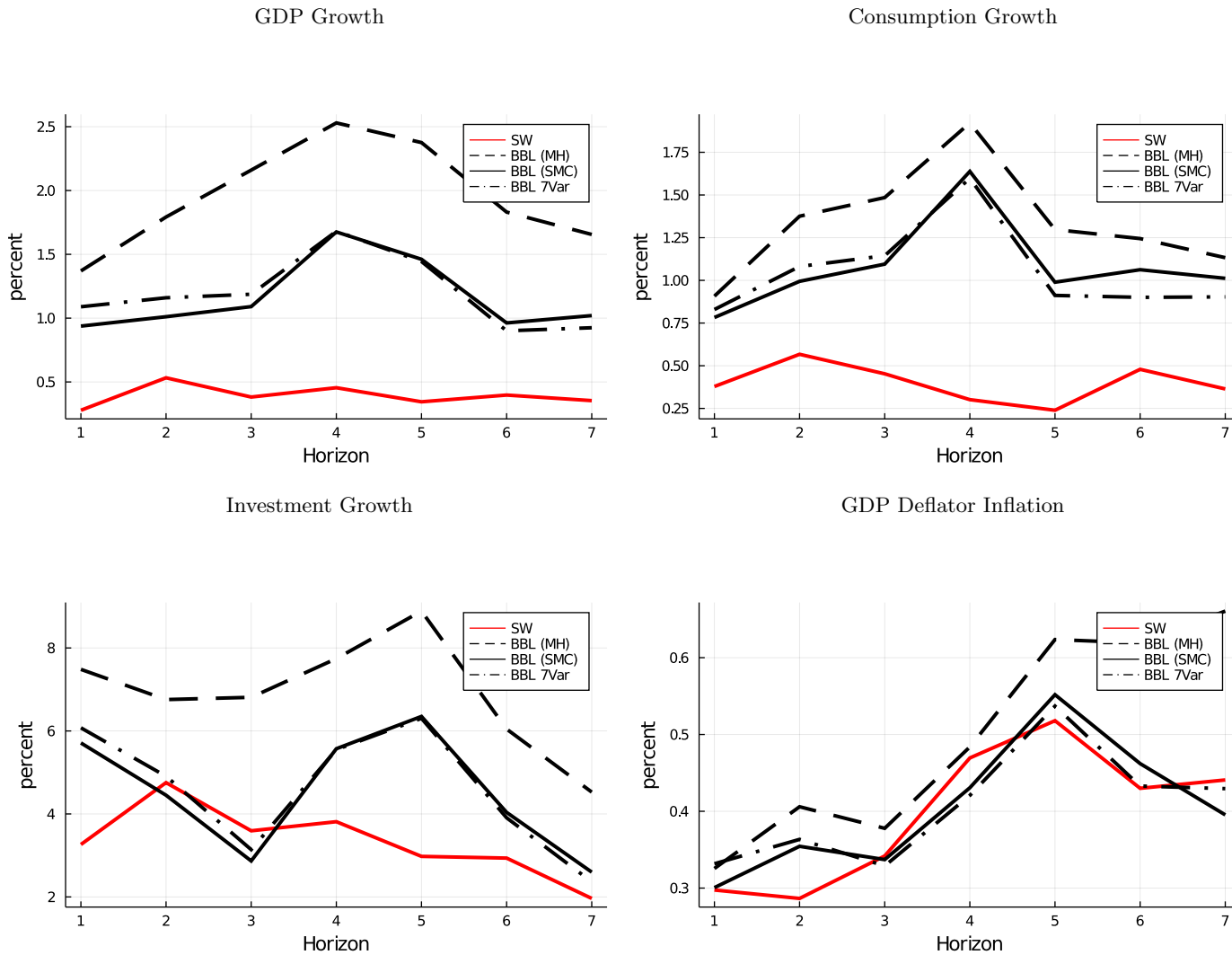
Figure 1: RMSEs: BBL vs SW



Note: The figure plots $RMSE_{i,h,\mathcal{M}_m}$ computed using expression (53) for the BBL (solid black lines) and the SW (solid red lines) models. The BBL model is referred to as BBL(SMC) as it uses the posterior computed from the online estimation starting from the SMC draws.

The much worse forecasting performance for BBL compared to SW for consumption growth is particularly disappointing. The key difference between HANK and SW-type models is the following: in HANK models the representative agent Euler equation, which determines consumption in standard DSGEs, is replaced with the aggregation of individual households' consumption policy functions. These consumption policy functions reflect inequality in both income and wealth: poor agents are hand-to-mouth, or close to, and have a high marginal propensity to consume out of income while richer agents can substitute intertem-

Figure 2: RMSEs: Robustness



Note: The figure plots $RMSE_{i,h,\mathcal{M}_m}$ computed using expression (53) for the seven-variable BBL model (BBL 7Var, dash-and-dotted black lines), the BBL model using the posterior computed from the online estimation starting from the MH draws, and referred to as BBL(MH) (dotted black lines), the BBL model using the posterior computed from the online estimation starting from the SMC draws, and referred to as BBL(SMC) (solid black lines), and the SW model (solid red lines).

porally and have lower marginal propensities to consume. The BBL version we use to compute the RMSEs in figure 1 incorporates as observables the variables reflecting inequality, such as the top 10 percent income and wealth shares. One would have hoped that this much more realistic view of the world translated into a better quantitative understanding of the behavior of aggregate consumption, and hence a better forecasting performance. This does not seem to be the case, at least for the BBL model.

Before discussing possible reasons for these findings, we show in figure 2 that the results are robust to using the results from i) the online estimation starting from the Metropolis Hastings (MH) draws, which we refer to as BBL(MH) (dotted black lines), and ii) the model using only the seven aggregate macro variables, and no measure of inequality, as observables (BBL 7Var, dash-and-dotted black lines). We find that the RMSEs obtained using the MH draws are uniformly worse than those obtained from the MH draws. We also find that the .RMSEs for the eleven and the seven-variable BBL are almost indistinguishable from one another. This is somewhat disappointing from the perspective of the HANK literature, as it suggests that measures of inequality matter little for the dynamics of macroeconomic aggregates, at least for this model. The result is reminiscent of the findings in Chang et al. (2021), who use functional vector autoregressions to argue that there is limited feedback between inequality and aggregate macro time series.

What are the possible reasons for these somewhat negative results? First, while BBL is *a priori* an ideal candidate for this forecasting comparison given that it incorporates SW's shocks and frictions, perhaps other HANK models may perform better than BBL from a forecasting point of view. Seen from this perspective, the results in this paper are an invitation to HANK modelers to use the methodology (and the code) described in this paper to see how well their model fares in terms of forecasting accuracy.

Second, the good (at least relative to VARs) forecasting performance of representative agent DSGEs à la SW was not achieved overnight, but resulted from a decade of advancement in modeling, crystallized in Christiano et al. (2005).²⁰ It may be that HANK models need to go through a similar process. There is also evidence (eg, Del Negro et al., 2007) that some of the reasonable forecasting performance of representative agent DSGEs is due to features like habit persistence that i) according to some may not have particularly strong micro-foundations, and ii) may be difficult to replicate in HANK models.

Finally, as mentioned in the introduction and discussed in section 4.2, the parameters in HANK affecting the model's steady state are calibrated, not estimated. This is for a computational reason: recomputing the steady state is extremely costly. But the estimated DSGE literature has shown that not estimating parameters, perhaps not too surprisingly, hurts the fit of DSGE models and their forecasting performance. If this is the reason why BBL forecasts worse than SW, these findings pose a computational challenge to HANK researchers interested in estimation: finding ways of computing the steady state more efficiently and/or using estimation algorithms that do not require recomputing the steady state too many times.

²⁰There is a perception among macroeconomists that the reasonable forecasting performance of DSGEs is the result of hindsight: model features are chosen *ex post* so that these models produce reasonably good RMSEs. More than ten years of actual (*ex ante*) forecasting with DSGE models at the NY Fed arguably shows that this perception is unfounded (Cai et al., 2019).

5 Conclusion

This paper had two objectives. One was to provide a toolkit for efficient repeated estimation of HANK model that can be used by researchers at central banks and in academia. We argued that online estimation using Sequential Monte Carlo provides such a toolkit, and we explained how it works. The second objective was to “kick the tires” of HANK models by comparing the out-of-sample forecasting accuracy of a prominent example of such models, Bayer et al. (2022), to that of the Smets and Wouters (2007) model. HANK models did not fare too well: their forecasting performance for real activity variables, especially GDP and consumption growth, is notably inferior to that of SW. The results for consumption are particularly disappointing, given that the main difference between SW-type DSGEs and HANK models is the replacement of the representative agent Euler equation with the aggregation of individual households’ consumption policy functions, which reflects inequality.

These findings should be interpreted as a motivation to do more research on HANK models. First, no matter the forecasting performance of HANK models, inequality is one of the critical issues of our times and features prominently in the transmission of policies. There are questions, such as investigating the effect on growth and inflation of the government transfers during the COVID pandemics, that representative agent models simply cannot adequately answer. Kaplan et al. (2020) and Auclert et al. (2023) are recent examples of quantitative research based on HANK models that focuses on some of these salient policy issues. Second, since all models are misspecified, model diversity should play an important role for policymakers that use models to inform their decisions. Finally, the fact that the forecasting performance of HANK models can be improved is a just stimulus for further efforts, in terms of both modeling and making computations more efficient.

References

- Ahn, SeHyouun, Greg Kaplan, Benjamin Moll, Thomas Winberry, and Christian Wolf**, “When inequality matters for macro and macro matters for inequality,” *NBER macroeconomics annual*, 2018, *32* (1), 1–75.
- An, Sungbae and Frank Schorfheide**, “Bayesian Analysis of DSGE Models,” *Econometric Reviews*, 2007, *26* (2-4), 113–172.
- Auclert, Adrien, Bence Bardóczy, Matthew Rognlie, and Ludwig Straub**, “Using the sequence-space Jacobian to solve and estimate heterogeneous-agent models,” *Econometrica*, 2021, *89* (5), 2375–2408.
- , **Matthew Rognlie, and Ludwig Straub**, “The Trickle Up of Excess Savings,” Technical Report, National Bureau of Economic Research 2023.
- Bayer, Christian and Ralph Luetticke**, “Solving heterogeneous agent models in discrete time with many idiosyncratic states by perturbation methods,” *Quantitative Economics*, 2020, *11*, 1253–1288.

- , **Benjamin Born**, and **Ralph Luetticke**, “Shocks, frictions, and inequality in US business cycles,” 2022.
- , **Ralph Lütticke**, **Lien Pham-Dao**, and **Volker Tjaden**, “Precautionary savings, illiquid assets, and the aggregate consequences of shocks to household income risk,” *Econometrica*, 2019, *87* (1), 255–290.
- Cai, Michael**, **Marco Del Negro**, **Edward Herbst**, **Ethan Matlin**, **Reca Sarfati**, and **Frank Schorfheide**, “Online estimation of DSGE models,” *The Econometrics Journal*, 2021, *24* (1), C33–C58.
- , – , **Marc P Giannoni**, **Abhi Gupta**, **Pearl Li**, and **Erica Moszkowski**, “DSGE forecasts of the lost recovery,” *International Journal of Forecasting*, 2019, *35* (4), 1770–1789.
- Calvo, Guillermo**, “Staggered Prices in a Utility Maximizing Framework,” *Journal of Monetary Economics*, 1983, *12* (3), 383–398.
- Cappé, Olivier**, **Eric Moulines**, and **Tobias Ryden**, *Inference in Hidden Markov Models*, Springer Verlag, New York, 2005.
- Chang, Minsu**, **Xiaohong Chen**, and **Frank Schorfheide**, “Heterogeneity and aggregate fluctuations,” Technical Report, National Bureau of Economic Research 2021.
- Chopin, Nicolas**, “A Sequential Particle Filter for Static Models,” *Biometrika*, 2002, *89* (3), 539–551.
- , “Central Limit Theorem for Sequential Monte Carlo Methods and its Application to Bayesian Inference,” *Annals of Statistics*, 2004, *32* (6), 2385–2411.
- Christiano, Lawrence J.**, **Martin Eichenbaum**, and **Charles L. Evans**, “Nominal Rigidities and the Dynamic Effects of a Shock to Monetary Policy,” *Journal of Political Economy*, 2005, *113*, 1–45.
- Creal, Drew**, “Sequential Monte Carlo Samplers for Bayesian DSGE Models,” *Manuscript, University Chicago Booth*, 2007.
- Del Negro, Marco** and **Frank Schorfheide**, “Monetary Policy with Potentially Misspecified Models,” *American Economic Review*, 2009, *99* (4), 1415–1450.
- and – , “Bayesian Macroeconometrics,” in Herman K. van Dijk, Gary Koop, and John Geweke, eds., *Handbook of Bayesian Econometrics*, Oxford University Press, 2010.
- and – , “DSGE Model-Based Forecasting,” in Graham Elliott and Allan Timmermann, eds., *Handbook of Economic Forecasting, Volume 2*, Elsevier, 2013.
- , – , **Frank Smets**, and **Raphael Wouters**, “On the Fit of New Keynesian Models,” *Journal of Business and Economic Statistics*, 2007, *25* (2), 123 – 162.
- Duane, Simon**, **Anthony D Kennedy**, **Brian J Pendleton**, and **Duncan Roweth**, “Hybrid monte carlo,” *Physics letters B*, 1987, *195* (2), 216–222.
- Durham, Garland** and **John Geweke**, “Adaptive Sequential Posterior Simulators for Massively Parallel Computing Environments,” in Ivan Jeliazkov and Dale Poirier, eds., *Advances in Econometrics*, Vol. 34, Emerald Group Publishing Limited, West Yorkshire, 2014, chapter 6, pp. 1–44.
- Edge, Rochelle** and **Refet Gürkaynak**, “How Useful Are Estimated DSGE Model Forecasts for Central Bankers,” *Brookings Papers of Economic Activity*, 2010, p. forthcoming.
- Farkas, Mátyás** and **Balint Tatar**, “Bayesian estimation of DSGE models with Hamiltonian Monte Carlo,” Technical Report, IMFS Working Paper Series 2020.
- Fernández-Villaverde, Jesús** and **Juan F Rubio-Ramírez**, “Estimating macroeconomic models: A likelihood approach,” *The Review of Economic Studies*, 2007, *74* (4), 1059–1087.

- Ferriere, Axelle and Gaston Navarro**, “The Heterogeneous Effects of Government Spending: It’s All About Taxes,” *FRB International Finance Discussion Paper*, 2018, (1237).
- Gelman, Andrew, John B Carlin, Hal S Stern, and Donald B Rubin**, *Bayesian data analysis*, Chapman and Hall/CRC, 1995.
- Geweke, John**, *Contemporary Bayesian econometrics and statistics*, John Wiley & Sons, 2005.
- Greenwood, Jeremy, Zvi Hercowitz, and Gregory W. Huffman**, “Investment, Capacity Utilization, and the Real Business Cycle,” *The American Economic Review*, 1988, 78 (3), 402–417.
- Hagedorn, Marcus, Iourii Manovskii, and Kurt Mitman**, “Monetary Policy in Incomplete Market Models: Theory and Evidence,” Technical Report, University of Pennsylvania Working Paper 2018.
- Hammersley, John M and K William Morton**, “Poor man’s monte carlo,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 1954, 16 (1), 23–38.
- Herbst, Edward**, “Using the “Chandrasekhar Recursions” for Likelihood Evaluation of DSGE Models,” *Computational Economics*, 2015, 45 (4), 693–705.
- and **Frank Schorfheide**, “Sequential Monte Carlo Sampling for DSGE Models,” *Journal of Applied Econometrics*, 2014, 29 (7), 1073–1098.
- and – , *Bayesian Estimation of DSGE Models*, Princeton University Press, 2015.
- Justiniano, Alejandro, Giorgio E. Primiceri, and Andrea Tambalotti**, “Investment shocks and the relative price of investment,” *Review of Economic Dynamics*, 2011, 14 (1), 102–121.
- Kaplan, Greg, Benjamin Moll, and Giovanni L Violante**, “Monetary policy according to HANK,” *American Economic Review*, 2018, 108 (3), 697–743.
- , – , and – , “The great lockdown and the big stimulus: Tracing the pandemic possibility frontier for the US,” Technical Report, National Bureau of Economic Research 2020.
- Klein, Paul**, “Using the generalized Schur form to solve a multivariate linear rational expectations model,” *Journal of Economic Dynamics and Control*, 2000, 24 (10), 1405–1423.
- Lee, Donggyu**, “Quantitative easing and inequality,” Technical Report, mimeo 2021.
- Liu, Jun S**, *Monte Carlo Strategies in Scientific Computing*, Springer Verlag, New York, 2001.
- Mlikota, Marko and Frank Schorfheide**, “Sequential Monte Carlo With Model Tempering,” *arXiv preprint arXiv:2202.07070*, 2022.
- Müller, Ulrich K**, “Measuring prior sensitivity and prior informativeness in large Bayesian models,” *Journal of Monetary Economics*, 2012, 59 (6), 581–597.
- Neal, Radford M et al.**, “MCMC using Hamiltonian dynamics,” *Handbook of markov chain monte carlo*, 2011, 2 (11), 2.
- Reiter, Michael**, “Solving heterogeneous-agent models by projection and perturbation,” *Journal of Economic Dynamics and Control*, 2009, 33 (3), 649–665.
- Smets, Frank and Raf Wouters**, “Shocks and Frictions in US Business Cycles: A Bayesian DSGE Approach,” *American Economic Review*, 2007, 97 (3), 586 – 606.
- Stan Development Team**, “Stan: A C++ library for probability and sampling, version 2.8. 0,” 2015.
- Winberry, Thomas**, “A method for solving and estimating heterogeneous agent macro models,” *Quantitative Economics*, 2018, 9 (3), 1123–1151.
- Wu, Jing Cynthia and Fan Dora Xia**, “Measuring the macroeconomic impact of monetary policy at the zero lower bound,” *Journal of Money, Credit and Banking*, 2016, 48 (2-3), 253–291.