

Institutional Members: CEPR, NBER and Università Bocconi

WORKING PAPER SERIES

Micro Responses to Macro Shocks

Martín Almuzara and Víctor Sancibrián

Working Paper n. 717

This Version: July 2025

IGIER – Università Bocconi, Via Guglielmo Röntgen 1, 20136 Milano – Italy http://www.igier.unibocconi.it

The opinions expressed in the working papers are those of the authors alone, and not those of the Institute, which takes non institutional policy position, nor those of CEPR, NBER or Università Bocconi.

Micro Responses to Macro Shocks*

Martín Almuzara[†] Víctor Sancibrián[‡]

July 2025
[Link to supplemental appendix]

Abstract

We study panel data regression models when the shocks of interest are aggregate and possibly small relative to idiosyncratic noise. This speaks to a large empirical literature that targets impulse responses via panel local projections. We show how to interpret the estimated coefficients when units have heterogeneous responses and how to obtain valid standard errors and confidence intervals. A simple recipe leads to robust inference: including lags as controls and then clustering at the time level. This strategy is valid under general error dynamics and uniformly over the degree of signal-to-noise of macro shocks.

Keywords: Panel data, local projections, impulse responses, aggregate shocks, inference, signal-to-noise, heterogeneity.

^{*}We have greatly benefited from discussions with Manuel Arellano, Dmitry Arkhangelsky, Richard Crump, Daniel Lewis, Mikkel Plagborg-Møller and Enrique Sentana. We also thank Stéphane Bonhomme, Òscar Jordà, Geert Mesters, seminar participants at Boston University, CEMFI, CUNY, Erasmus University Rotterdam, Philadelphia Fed, Georgetown, Princeton, Universidad Autónoma de Madrid and UPenn, and participants at the EABCN-UPF conference, Greater New York Metro Area Econometrics Colloquium, the NBER Summer Institute, the 2024 SED Annual Meeting, the 2024 System Econometrics Conference at Richmond Fed and the 2024 CEME Conference for Young Econometricians. Babur Kocaoglu provided excellent research assistance. Sancibrián acknowledges financial support from Grant PRE2022-000906 funded by MCIN/AEI/10.13039/501100011033 and by "ESF +", the María de Maeztu Unit of Excellence CEMFI MDM-2016-0684, funded by MCIN/AEI/10.13039/501100011033, Fundación Ramón Areces and CEMFI. A big part of this work was done during Sancibrián's PhD studies at CEMFI. The views expressed in this paper do not necessarily reflect the position of the Federal Reserve Bank of New York or the Federal Reserve System.

[†]FRBNY: martin.almuzara@ny.frb.org

[‡]Bocconi University and IGIER: victor.sancibrian@unibocconi.it

1 Introduction

Applied macroeconomists are increasingly interested in empirical estimates of the transmission of aggregate uncertainty to individual outcomes, often in the form of impulse responses.

A popular approach is to formulate estimating equations of the form

$$Y_{i,t+h} = \beta(h)s_i X_t + \text{controls} + v_{h,it}, \tag{1}$$

where Y_{it} is a *micro outcome* for unit i (i = 1, ..., N) at time t (t = 1, ..., T) and X_t a *macro shock* of interest. Shocks are often interacted with unit-level covariates s_i to document heterogeneity in transmission along observables. For example, Ottonello and Winberry (2020) and Crouzet and Mehrotra (2020) are interested in the heterogeneous effects of monetary policy shocks X_t on firm-level investment and sales along different margins, such as firm size or leverage. Estimates $\hat{\beta}(h)$ of the response at horizon h are then obtained via least squares; a panel local projections version of Jordà (2005).

Despite its routine application, little is known about the statistical properties of $\hat{\beta}(h)$. The way standard errors are computed in the empirical literature illustrates it well: in our own survey of almost 50 recent papers, around half compute two-way clustered standard errors, one-third cluster within units only, and many others resort to Driscoll and Kraay (1998). This reflects the vastly different ways in which researchers perceive the nature of shocks, the role of each dimension of the panel for precision, and the importance of aggregate variation in the data.

In this paper, we provide the first treatment of estimation and inference for this problem. We show how to interpret $\hat{\beta}(h)$ when impulse-response heterogeneity is unrestricted, and propose standard errors and confidence intervals that are easy to compute and robust to the signal-to-noise of macro shocks in the microdata. As a result, a very simple recipe for inference emerges: clustering standard errors at the time level and ex-ante including enough lags as controls. We refer to this as *time-clustered lag-augmented heteroskedasticityrobust* (*t-LAHR*) *inference*.¹

We establish our results in a comprehensive setup featuring observed and unobserved macro and micro shocks, cross-sectional heterogeneity in responses, general forms of

¹A full Matlab package for panel local projections — including estimation and *t*-LAHR inference — and replication files are available at https://github.com/TinchoAlmuzara/PanelLocalProjections.

serial dependence, and arbitrary signal-to-noise. Our analysis extends to a rich panel environment most of the empirical strategies from macroeconometrics (Ramey, 2016; Stock and Watson, 2018; Plagborg-Møller and Wolf, 2021), including some which are — to our knowledge — novel and empirically promising in the panel context. Specifically, we cover settings where the shock of interest is directly observed (the most prevalent assumption in applications), those where it is recoverable from a macro system (via, say, recursive or long-run identification), and those where it is contaminated with measurement error but a proxy is available (as in local projection-instrumental variables; LP-IV for short).²

In that setup, we first show that $\hat{\beta}(h)$ recovers the slope coefficient of a population linear projection of unit-specific impulse responses on the characteristics s_i , thereby formalizing what practitioners have in mind when including interactions in Equation (1). Importantly, this is also the case when the shock of interest is unobserved, such as in the instrumental variables setup where only a noisy measurement is available. If $s_i = 1$, the estimand boils down to the average response in the population. Crucially, since we place no restrictions on the underlying impulse-response heterogeneity or in s_i , our characterization of the estimand is effectively nonparametric.³

Signal-to-noise. The degree of signal-to-noise of macro shocks in the microdata is a key parameter of the problem. *Common* shocks to all units drive identification and, therefore, how sizable they are relative to micro shocks determines both the strength of identifying variation and the extent of unaccounted-for spatial dependence.⁴ Considering different signal regimes also reflects the scope of empirical work, which takes interest in atomistic and granular agents, administrative and narrow datasets, unit-specific and aggregate regression controls, etc.

Hence, one of our main contributions is to introduce a novel asymptotic framework where the signal value of aggregates may be arbitrarily low (or high) in the limit. We achieve this by indexing the relative standard deviation of macro to micro shocks by a

²Narrative approaches (as in Crouzet and Mehrotra, 2020, for monetary policy shocks) and high-frequency approaches (as in Känzig, 2021, for oil supply shocks) are examples of popular identification methods where shocks are typically treated as observable. In Section 3.4, we argue that it is sometimes more appropriate to view these as proxies and allow for measurement error. Drechsel (2023) imposes long-run restrictions on a structural VAR model to identify investment-specific technology shocks.

³We discuss extensions to (exogenous) time-varying characteristics s_{it} in Section 3 (Remark 7).

⁴It is immediate that if $s_i = 1$ in Equation (1), including time fixed effects causes collinearity. If s_i varies over units, for time indicators to remove all additional aggregate variation one would need the untenable assumption that *only* impulse responses to X_t at horizon h are heterogeneous. In our exposition, we always allow for time indicators as controls when s_i displays cross-sectional variation.

parameter κ that can drift with the sample size. This device allows for a range of data generating processes (DGPs) in which estimation uncertainty is dominated by micro-level terms, a combination of micro and macro errors, or aggregate components only.⁵ On the contrary, standard asymptotic plans where κ is fixed only capture the latter and ignore idiosyncratic shocks, potentially leading to poor approximations in small samples. It is clear then that the nature of estimation error depends on κ and the question is whether inference procedures are robust to different macro signal regimes. Our main result is that t-LAHR inference is uniformly valid over κ , that is, t-LAHR confidence intervals have correct asymptotic coverage for the (nonparametric) local projection estimand uniformly over κ .

Shocks and inference. Our notion of macro shock is that of an innovation uncorrelated to its own lags and leads and other shocks, as is standard in the time series literature (Leeper, Sims, and Zha, 1996; Ramey, 2016; Stock and Watson, 2018; Montiel Olea and Plagborg-Møller, 2021). This is an identifying assumption without which the estimand of $\hat{\beta}(h)$ may not be interpretable as an impulse response. Empirical researchers widely recognize this and in practice devote great effort to constructing and motivating X_t by leveraging methods from macroeconometrics. We show that it also has important consequences for inference with microdata.

The *macro shock* nature of X_t delivers a connection which serves as a guiding principle throughout the paper: panel local projections with macro shocks are equivalent to *synthetic* time series local projections with an appropriately aggregated dependent variable. This is true even if shocks interact with covariates s_i and if unit and time effects are included. Then, aggregating the microdata by collapsing the cross-sectional dimension of the panel and treating it as a time series yields valid inference for any κ .⁶ This is precisely what t-LAHR inference does, since time-level clustering in the panel problem and heteroskedasticity-robust inference in the synthetic time series problem are equivalent.

The macro shock nature of X_t also clarifies the role of lag augmentation. In a panel

⁵Our approach also resonates with the renewed interest on the potential for unit-level shocks to explain aggregate fluctuations, as in Gabaix (2011) and subsequent literature. Our device to obtain non-negligible micro errors is closer to Jovanovic (1987) in that we rely on scaling micro variation up rather than on fat-tailed distributions. However, we conjecture that similar inference results can be obtained in the latter under appropriate regularity conditions.

⁶This synthetic time series representation is also illustrative of the fact that the concentration rate of the estimation error is at most $T^{-1/2}$, even in situations where $N \gg T$. This suggests caution regarding the conventional wisdom in many empirical applications that a larger cross-sectional dimension somehow compensates for a shorter time series.

local projection that controls for p lags of s_iX_t , the regression scores (the product of shocks and residuals) are *nearly* uncorrelated even if residuals are not. Specifically, they are a moving average of order h where the first p autocovariances are zero and the remaining ones do not depend on κ . This has two major implications. First, it confers a double layer of simplicity to inference: up to horizon $h \le p$, there is no need for unit-level clustering or heteroskedasticity and autocorrelation robust (HAR) approaches to deal with serial dependence. Second, it explains why t-LAHR inference might have only small coverage distortions even for horizons exceeding p: these distortions depend on the size of the autocorrelation coefficients of the score, which are small in low-signal environments. In fact, if the DGP is well approximated by a low-order vector autoregression (VAR), we prove t-LAHR inference is uniformly valid over both κ and $h \propto T$, a result reminiscent of those in Montiel Olea and Plagborg-Møller (2021) for time series local projections.

We complement our theoretical results with simulations for realistic designs and sample sizes, allowing for moderately long horizons and substantial persistence in micro shocks. We study the performance of a battery of approaches, including an alternative to *t*-LAHR that substitutes lag augmentation with HAR inference, and incorporating small-sample refinements (Müller, 2004; Imbens and Kolesár, 2016; Lazarus, Lewis, Stock, and Watson, 2018). We find that *t*-LAHR inference shows remarkable performance relative to all other competitors, particularly in low-signal environments, in near non-stationary scenarios, and over moderate horizons even if we do not impose a VAR on outcomes.⁷ In practice, we recommend to supplement *t*-LAHR inference (controlling for lags of both outcome and shock variables) with the refinement proposed in Imbens and Kolesár (2016).

Empirical survey and illustration. We reviewed a large body of empirical work that precedes this paper. The typical application uses administrative data on firms or households, tracks units at the quarterly or annual frequency for a limited number of periods, and estimates impulse responses to shocks via local projections. Most applications include interactions of the form $s_i X_t$ and both unit and time fixed effects, but vary widely in the number and nature of additional controls.⁸

⁷It is known that ad-hoc parameter choices and small-sample biases in sample autocovariances contribute to the subpar relative performance of HAR estimators (Herbst and Johansenn, 2023).

⁸We reproduce the full list in Supplemental Appendix ?? which includes 47 empirical papers that run panel regressions with macro shocks. A few focus on the case h=0 only, but the vast majority compute impulse responses over several horizons. The economic content of X_t is very diverse, including fiscal policy shocks, investment shocks, TFP and innovation shocks, carbon pricing shocks, temperature shocks, etc. In these applications, the cross-sectional dimension is usually orders of magnitude larger than the effective

In otherwise comparable empirical designs, we document large dispersion in the way practitioners compute standard errors: among 47 different papers, 24 compute two-way clustered standard errors (within units and time), 15 cluster within units only, 7 use Driscoll and Kraay (1998) and 1 clusters within time only.

These choices reflect diverging views on the dominant sources of statistical uncertainty, from ruling out serial dependence to ruling out spatial dependence; from a suggestion that both unit-level and aggregate dynamics matter to an explicit stance that either of the two dominates. Often, these choices are made with little discussion or citing previous work as a justification. Our framework allows us to revisit them. First, off-the-shelf autocorrelation consistent methods such as Driscoll and Kraay (1998) leave information on the table about the autocovariance function of the regression scores, which comes at a cost in small samples. Second, validity in the case where standard errors do not explicitly adjust for serial dependence (as in two-way clustering) boils down to whether a reasonable number of lags was included in estimation. Third, clustering within units is superfluous, even in low-signal regimes where the size of unit-level dynamics is comparable to that of aggregates. Fourth, clustering only within units breaks down in the face of even small amounts of spatial dependence induced by aggregate shocks, that is, in high- and moderate-signal environments.

Finally, we illustrate our methods in an empirical exercise inspired by a booming literature that investigates the role of financial frictions and firm heterogeneity in the transmission of monetary policy. For instance, Crouzet and Mehrotra (2020), Ottonello and Winberry (2020), Anderson and Cesa-Bianchi (2024) and Jeenas and Lagos (2024) target impulse responses of firm investment to monetary policy shocks interacted with external covariates s_i such as firm size, default risk or stock turnover. The exercise highlights the importance of the choice of inference method, and the value of the synthetic time series representation as a way to gain intuition about the source of identifying variation.

Related literature. Our paper contributes to various strands of the literature.

First, it relates to the time series literature on local projection inference (Hansen and

time-series dimension. In our review we leave out empirical work with very small cross-sections where entities are meaningful and a unit-by-unit treatment is feasible. Nonetheless, when these units are pooled together, as in Fukui, Nakamura, and Steinsson (2023), our results still apply.

⁹The availability of a large cross-sectional dimension and the interaction of shocks with covariates s_i are also often argued as sources of large gains in statistical precision, also reflecting an implicit stance on the presence of macro shocks. We elaborate on the (im)plausibility of these notions below and in Remark 5.

Hodrick, 1980; Jordà, 2005; Stock and Watson, 2018; Montiel Olea and Plagborg-Møller, 2021; Lusompa, 2023; Xu, 2023; Montiel Olea, Plagborg-Møller, Qian, and Wolf, 2024). Relative to this literature we are the first to deal with the panel data case with aggregate shocks. We supplement the macroeconometric toolbox with results that are unique to the panel setting by accommodating cross-sectional heterogeneity in responses, spatial dependence and coexisting micro and macro shocks of any relative sizes.

In a time series finite-order VAR setup, Montiel Olea and Plagborg-Møller (2021) show the uniform validity of heteroskedasticity-robust inference on lag-augmented local projections over the persistence of the data and horizon h. They postulate mean independent innovations, the same type of assumption we impose on X_t . Our Proposition 2 can be interpreted as the panel version of their results. However, our focus is on uniformity with respect to the macro signal-noise κ , which has no obvious counterpart in the time series setup, and we derive most of our results without assuming a VAR model.

Second, we contribute to the literature on estimation and inference with aggregate shocks and microdata. In stylized models, Hahn, Kuersteiner, and Mazzocco (2020) bring attention to the drastic consequences of drawing inferences from short panels with aggregate uncertainty. Although our focus is on thought experiments where macro shocks are a key source of identification, we can connect to their results by reinterpreting confidence intervals that exploit independence across units as valid for an alternative estimand that conditions on the path of aggregate shocks. Chang, Chen, and Schorfheide (2024) propose a methodology based on functional vector autoregressions and repeated cross-sections to study the effects of aggregate shocks on cross-sectional distributions. We study the complementary problem of heterogeneity in transmission of aggregate shocks to unit-level outcomes via local projections.

An important empirical literature exploits regional-exposure designs with instruments of the form $s_i X_t$ used to identify the causal link between Y_{it} and a unit-level variable W_{it} , yielding a reduced-form equation identical to (1) for h = 0. Recent work emphasizes the role of the exogeneity of X_t for credible identification (Adão, Kolesár, and Morales, 2019;

¹⁰Our results on limited serial dependence in regression scores relate to the earlier multi-step forecast literature (Hansen and Hodrick, 1980), which relied on infinite lags to ensure forecast errors have an MA(*h*) representation. In the local projection context, Jordà (2005) arrived at a similar result under a finite-order VAR model while Lusompa (2023) provided a recent reformulation. Instead, we exploit the orthogonality properties of macro shocks to show that the *scores* have MA(*h*) dynamics. The distinction is reminiscent of the difference between design-based and model-based/conditional unconfoundedness assumptions.

¹¹Examples include regional fiscal multipliers (Nakamura and Steinsson, 2014) or the effects of foreign aid on conflict across countries (Nunn and Qian, 2014).

Arkhangelsky and Korovkin, 2023; Majerovitz and Sastry, 2023). The broad empirical literature to which our paper speaks is related, but different in fundamental ways: panel local projections interpret aggregate shocks X_t in the macroeconometric sense and aim at measuring their dynamic propagation via impulse responses, whereas regional-exposure (and shift-share) approaches target the relative effects of X_t on Y_{it} and W_{it} treating X_t as an auxiliary shifter. They are complementary tools addressing different questions, but the econometrics of panel local projections is less understood. Our paper adds to it formal inference results in a rich panel environment that allows for heterogeneity, dynamics and different macro signal regimes.

In some cases, that literature has also relied on quasi-random variation in exposures s_i for identification. In our context, this would entail very strong requirements, analogous to finding a separate independent source of exogenous variation for every horizon of interest. This point echoes a similar conclusion at which Hahn, Kuersteiner, Santos, and Willigrod (2024) arrive in shift-share setups with heterogeneous treatment effects.

Last, our paper relates to the cross-sectional dependence literature that studies models where the scores feature spatial correlation (Driscoll and Kraay, 1998; Andrews, 2005; Pesaran, 2006; Gonçalves, 2011; Pakel, 2019). Our setting lies in the polar case where the shock of interest only varies over time, precluding solutions based on partialling out the common component from the regressors, as in Pesaran (2006). Our uniformity result (that translates into robustness to the degree of spatial dependence) is new to this literature.

Outline. Section 2 provides an overview of our results in a simple static model, illustrating the role of aggregate shocks and their signal relative to micro shocks. Section 3 presents our main inference result in a general, heterogeneous dynamic model. Section 4 discusses a comprehensive simulation study and Section 5 the empirical illustration. Proofs can be found in Appendix A with additional details in the Supplemental Appendix. A MATLAB code repository is available online at https://github.com/TinchoAlmuzara/PanelLocalProjections.

 $^{^{12}}$ Consider, for instance, the role of firm size (s_i) in mediating the transmission of monetary policy shocks to firm-level outcomes, as in Crouzet and Mehrotra (2020) and our empirical application. For the cross-sectional dimension of the panel to help pin down the impulse response at, say, h=2, firm size must be orthogonal to firm-level responses at any other horizon and to any other aggregate shock. This is, of course, hardly plausible. In fact, such exclusion restrictions would rule out dynamic effects over different horizons altogether — the object of empirical interest. We elaborate on this in Remark 5.

2 Simple model

We illustrate the main points of the paper in a simple, static regression model with homogeneous responses. We keep the exposition simple and omit technical details with the goal of building insights. The general setup is studied in Section 3.

Model assumptions. We observe a micro outcome Y_{it} and a macro shock X_t for units i = 1, ..., N and over periods t = 1, ..., T. They are related by

$$Y_{it} = \beta_0 X_t + v_{it},$$

$$v_{it} = Z_t + \kappa u_{it},$$
(2)

where v_{it} is an error term including both aggregate and idiosyncratic unobservables, denoted Z_t and u_{it} , respectively. Here κ regulates their relative importance in the microdata, as explained below. The goal is to estimate and do inference on β_0 .

This simple model is a stylized representation of an empirical setting where we are interested in the transmission of aggregate uncertainty to individual outcomes; the effect of X_t on Y_{it} . Examples of the former include changes in interest rates, tax regulations or oil prices, which might leave a mark on household consumption, worker's labor income or firm sales. In fact, one could entertain any combination of macro variables and micro outcomes in these examples. When interest centers around one aggregate variable, captured by X_t , it would be hard to ex-ante rule out the presence of any others, embedded in Z_t . This basic premise is at the core of our results.

Here we are treating X_t as observable, in line with most of our empirical applications of reference (Supplemental Appendix ??), but we consider other possibilities below and in Section 3.1. We now make two sets of assumptions, later generalized in Section 3 to allow for observable and unobservable heterogeneity, flexible dynamics and serial dependence.

Assumption S1 (Stationarity and iidness in the simple model).

- (i) $\{X_t, Z_t, \{u_{it}\}_{i=1}^N\}_{t=1}^T$ is stationary.
- (ii) $\{\{u_{it}\}_{t=-\infty}^{\infty}\}_{i=1}^{N}$ are i.i.d. over i conditional on $\{X_t, Z_t\}_{t=1}^{T}$.

Assumption S1(i) implies Y_{it} is stationary too. Assumption S1(ii) simply assigns the role of inducing cross-sectional dependence in the error term v_{it} to Z_t .¹³

¹³Both assumptions can be relaxed; we briefly discuss departures from S1(i) in Section 3 and 4. Allowing

Assumption S2 (Shocks and independence in the simple model).

$$(i) \ \ E\bigg[\ X_t \ \bigg| \ \{X_\tau\}_{\tau \neq t}, \Big\{ Z_\tau, \{u_{i\tau}\}_{i=1}^N \Big\}_{\tau=1}^T \ \bigg] = 0.$$

$$(ii) \ E \left[\left. Z_t \right. \left| \left. \left\{ Z_\tau \right\}_{\tau \neq t}, \left\{ X_\tau, \left\{ u_{i\tau} \right\}_{i=1}^N \right\}_{\tau=1}^T \right. \right] = 0.$$

$$(iii) \ E \left[\left. u_{it} \right| \left\{ u_{i\tau} \right\}_{\tau \neq t}, \left\{ X_{\tau}, Z_{\tau} \right\}_{\tau = 1}^T \right] = 0.$$

Assumption S2 implies that X_t , Z_t and u_{it} are mutually unpredictable and serially uncorrelated: S2(i) is ultimately an identification condition, whereas S2(ii) and S2(iii) are made for symmetry. Indeed, mutual unpredictability of macro shocks lies at the core of macroeconometrics and is typically necessary to give structural interpretation to impulse-response calculations (see, for instance, Leeper et al., 1996; Ramey, 2016; Stock and Watson, 2016; Plagborg-Møller and Wolf, 2021). 14

Remark 1 (When X_t is not directly observable). Researchers rely on different strategies, such as narrative or high-frequency approaches, to construct measurements of the shocks of interest X_t . In practice, it is often more appropriate to treat them as imperfect measures of X_t , contaminated with measurement error or with some residual autocorrelation structure. In other cases, researchers treat X_t as recoverable from a macro system after imposing certain restrictions, as in recursive or long-run identification. In Section 3 we show that our analysis extends to all these settings.

Estimation and inference. A natural estimator of β_0 is pooled least squares,

$$\hat{\beta} = \frac{\sum_{i=1}^{N} \sum_{t=1}^{T} X_{t} Y_{it}}{\sum_{i=1}^{N} \sum_{t=1}^{T} X_{t}^{2}} = \frac{\sum_{t=1}^{T} X_{t} \left(N^{-1} \sum_{i=1}^{N} Y_{it}\right)}{\sum_{t=1}^{T} X_{t}^{2}},$$

which is also a panel local projection (LP) estimator at horizon h = 0 and the estimator in a time series regression involving the synthetic outcome $\hat{Y}_t = N^{-1} \sum_{i=1}^{N} Y_{it}$ and X_t . The double nature of $\hat{\beta}$ as panel and time series estimator arises organically in the presence of macro shocks, as we further demonstrate in Section 3.

for weak spatial dependence in u_{it} in place of S1(ii) is also possible with minor modifications.

¹⁴Mean independence assumptions with respect to past and future innovations are a slight strengthening of the more standard martingale difference assumptions, and are convenient in representations where both leads and lags of the variable might enter the model, cf. Montiel Olea and Plagborg-Møller (2021, Assumption 1) in a similar context of local projection inference. This still allows for dynamics on the second-or higher-order moments given the paths of other shocks. It permits that, say, monetary, fiscal or oil supply shocks (X_t, Z_t) increase the variance of household-level income (Y_{it}) via higher order dynamics in u_{it} .

Denote the residual by $\hat{\xi}_{it} = Y_{it} - \hat{\beta}X_t$. A key takeaway from our paper is that a reliable approach to inference uses the time-level cluster heteroskedasticity-robust standard error $\hat{\sigma}$, given by $\hat{\sigma}^2 = \hat{V}/T\hat{J}^2$ where $\hat{J} = (NT)^{-1}\sum_{i=1}^N\sum_{t=1}^TX_t^2 = T^{-1}\sum_{t=1}^TX_t^2$ and

$$\hat{V} = \frac{1}{T} \sum_{t=1}^{T} \left(\frac{1}{N} \sum_{i=1}^{N} X_{t} \hat{\xi}_{it} \right)^{2}.$$

Another sign of the duality between panel regressions with aggregate shocks and time series regression is that $\hat{\sigma}$ is also the usual Eicker–Huber–White standard error computed using the synthetic time series residuals $\hat{\xi}_t = N^{-1} \sum_{i=1}^N \hat{\xi}_{it}$.

As mentioned in the Introduction, two popular inferential choices in applications are based on one-way (unit-level) cluster and two-way (unit- and time-level) cluster standard errors, $\hat{\sigma}_{1W}$ and $\hat{\sigma}_{2W}$, given by $\hat{\sigma}_{1W}^2 = \hat{V}_{1W}/T\hat{J}^2$ and $\hat{\sigma}_{2W}^2 = \hat{V}_{2W}/T\hat{J}^2$ where

$$\hat{V}_{1W} = \frac{1}{N} \sum_{i=1}^{N} \left(\frac{1}{T} \sum_{t=1}^{T} X_{t} \hat{\xi}_{it} \right)^{2}, \quad \hat{V}_{2W} = \hat{V} + \hat{V}_{1W} - \frac{1}{NT} \sum_{t=1}^{T} \sum_{i=1}^{N} X_{t}^{2} \hat{\xi}_{it}^{2}.$$

These standard errors reflect different concerns about the nature of estimation error or, more precisely, the correlation of the regression score $X_t v_{it}$ over units and time.

Substituting (2), the estimation error decomposes as

$$\hat{\beta} - \beta_0 = \underbrace{\frac{\sum_{t=1}^T X_t Z_t}{\sum_{t=1}^T X_t^2}}_{O_p(1)} + \frac{\kappa}{\sqrt{N}} \underbrace{\left(\frac{1}{\sqrt{N}} \frac{\sum_{i=1}^N \sum_{t=1}^T X_t u_{it}}{\sum_{t=1}^T X_t^2}\right)}_{O_p(1)}, \tag{3}$$

i.e., as the sum of macro and micro components. The former induces cross-sectional correlation, the latter is uncorrelated across units and both have limited serial dependence: $E[X_tZ_t \cdot X_\tau Z_\tau] = E[X_tu_{it} \cdot X_\tau u_{i\tau}] = 0$ for $t \neq \tau$ by S2(i) and iterated expectations.¹⁵

The intuition for why $\hat{\sigma}$ gives valid inference is as follows. If the macro term is not asymptotically small, $X_t v_{it}$ displays correlation over i but not over t, the type of situation for which $\hat{\sigma}$ is designed. If, on the other hand, the micro term dominates, $X_t v_{it}$ is uncorrelated over both i and t. Yet $\hat{\sigma}$ still works: although it does not impose that the cross-sectional covariances of $X_t v_{it}$ are zero, it will correctly estimate them to be zero. One may wish

¹⁵This is a direct consequence of X_t being a shock. The lack of serial correlation would remain true even if Z_t and u_{it} were serially correlated.

to switch to a non-clustered heteroskedasticity-robust standard error in that case, but we show both analytically (Proposition 1) and in simulations (Section 4) that there is no loss in simply using $\hat{\sigma}$. Clearly, correlation over t at the unit-level is never a concern; that is why unit-level clustering either fails or is not needed. In fact, $\hat{\sigma}_{1W}$ is asymptotically equivalent to the non-clustered standard error, and the same holds for $\hat{\sigma}_{2W}$ and $\hat{\sigma}$.

Macro-micro signal-to-noise ratio. Which term dominates the decomposition (3) will depend upon κ/\sqrt{N} . We now provide another interpretation of this quantity. Consider the average outcome $\hat{Y}_t = N^{-1} \sum_{i=1}^N Y_{it}$ and, to illustrate, suppose $\text{Var}(Z_t) = \text{Var}(u_{it}) = 1$. By Assumptions S1 and S2, the proportion of the variance of \hat{Y}_t explained by the unobserved macro error can be measured as

$$\bar{R}^{2}(\kappa) = 1 - \frac{\operatorname{Var}(\hat{Y}_{t} \mid X_{t}, Z_{t})}{\operatorname{Var}(\hat{Y}_{t} \mid X_{t})} = \frac{1}{1 + \kappa^{2}/N'}$$
(4)

i.e., the signal-noise ratio is $O(N/\kappa^2)$. It increases with N since cross-sectional averaging reduces the variance from idiosyncratic errors, but decreases with $|\kappa|$.

We will study estimation and inference along sequences of data generating processes (DGPs) where κ is allowed to grow as $T,N\to\infty$. This leads, in essence, to three regimes. If $\kappa/\sqrt{N}=o(1)$, (such as if κ is fixed), $\bar{R}^2(\kappa)\to 1$ and macro shocks are the only source of aggregate variation; we call this the asymptotically high-signal case. If $\kappa \propto \sqrt{N}$, $\bar{R}^2(\kappa)$ is bounded away from 0 and 1 in the limit and both macro and micro shocks matter for aggregate fluctuations; this is the asymptotically moderate-signal case. Finally, if κ/\sqrt{N} diverges, $\bar{R}^2(\kappa)\to 0$, macro shocks are imperceptible and we are in the asymptotically low-signal case. ¹⁶

The intuitive notion of κ -regimes has a natural counterpart in our asymptotic approximations, in that there is a close relation between the contribution of macro shocks to \hat{Y}_t and the nature of estimation error for β_0 , as illustrated by (2) and (4). In particular, the macro term dominates in the high-signal case, the micro term dominates in the low-signal case, and they are roughly balanced in the moderate-signal case. Moreover, it is not always possible to consistently detect what κ -regime applies. It is important then to derive

¹⁶Of course, letting κ grow with the sample size should not be taken literally — it is simply a device to ensure our approximations suitably interpolate between high and low signal-noise environments. This type of embeddings are common in econometrics; an example which also has a low-signal interpretation is weak IV (Staiger and Stock, 1997).

inference procedures that are robust in the sense of uniform validity with respect to κ . ¹⁷

Uniformity over κ **.** From the decomposition in (3), letting $N, T \rightarrow \infty$ and under regularity conditions specified in Section 3,

$$\sigma_0(\kappa)^{-1} \sqrt{T} \left(\hat{\beta} - \beta_0 \right) \xrightarrow{d} N(0, 1),$$

where

$$\left\{E\left[X_{t}^{2}\right]\right\}^{2} \times \sigma_{0}(\kappa)^{2} = \begin{cases} E\left[X_{t}^{2}Z_{t}^{2}\right], & \text{if } \kappa/\sqrt{N} \to 0, \\ E\left[X_{t}^{2}\left(Z_{t}^{2} + \bar{\kappa}^{2}u_{it}^{2}\right)\right], & \text{if } \kappa/\sqrt{N} \to \bar{\kappa}, \\ \left(\kappa^{2}/N\right)E\left[X_{t}^{2}u_{it}^{2}\right], & \text{if } \kappa/\sqrt{N} \to \infty, \end{cases}$$

This shows two things. First, the rate of concentration of the estimation error $\hat{\beta} - \beta_0$ is either \sqrt{T} in the high- and moderate-signal cases or \sqrt{NT}/κ (i.e., slower than \sqrt{T} and possibly even zero, thus making $\hat{\beta}$ inconsistent) in the low-signal case. Second, the asymptotic distribution of $\hat{\beta}$ changes discontinuously across κ -regimes.

Despite the discontinuity, our main result is that the $(1 - \alpha)$ confidence interval given by $\hat{C}_{\alpha} = \left[\hat{\beta} \pm z_{1-\alpha/2}\hat{\sigma}\right]$, where z_q is the q-quantile of the standard normal distribution, has correct coverage for β_0 uniformly over κ ,

$$\lim_{T,N\to\infty} \sup_{\kappa} \left| P_{\kappa} \left(\beta_0 \in \hat{C}_{\alpha} \right) - (1-\alpha) \right| = 0.$$
 (5)

where P_{κ} denotes probabilities for a DGP with a given κ . This is much stronger than pointwise validity, as it implies that the quality of the asymptotic approximation to the coverage probability of \hat{C}_{α} is itself robust to the κ -regime. Statement (5) also means that if sample information about macro shocks is extremely scarce and $\hat{\beta}$ is inconsistent, the length of \hat{C}_{α} adjusts as needed to reflect the weak macro signal.

One might wonder how much the static nature of (2) limits these results. The rest of the paper will show that they extrapolate to a substantially more general and empirically realistic framework with rich forms of dynamics, serial correlation and heterogeneity.

Remark 2 (Inference conditional on aggregate shocks). Ignoring the unobservable macro

¹⁷We will consider inference procedures that are invariant to rescaling. It follows that all of our results can be equivalently obtained in an embedding that scales down the macro component of the model in (2) by κ^{-1} . Put differently, what matters is the relative size of macro and micro components.

component in (3) when doing inference is equivalent to conditioning on its realization. In that situation, $\hat{\sigma}_{1W}$ is a valid standard error for responses defined by moment restrictions that condition on the realized path of aggregate shocks during the sample period.¹⁸ In general, this induces an internal/external validity trade-off whereby practitioners may be able to pin down certain parameters very precisely but these might lack generalizability.

3 General case

In this section, we establish estimation and inference results for impulse responses to aggregate shocks in a general setup featuring observed and unobserved, macro and micro shocks, and unrestricted heterogeneity of individual responses. We introduce the setup in Section 3.1 and state the main results in Section 3.2. We treat finite-order VAR DGPs in Section 3.3 and local projections with instrumental variables (LP-IV) in Section 3.4. Proofs are developed in Appendix A with technical lemmas in Supplemental Appendix ??.

3.1 Setup

The researcher observes an outcome Y_{it} , an aggregate shock X_t and characteristics s_i for units i = 1, ..., N and over periods t = 1, ..., T. They are scalar but it is straightforward to extend our results to the multivariate case. The DGP is

$$Y_{it} = \mu_i + \sum_{\ell=0}^{\infty} \beta_{i\ell} X_{t-\ell} + v_{it},$$

$$v_{it} = \sum_{\ell=0}^{\infty} \gamma_{i\ell} Z_{t-\ell} + \kappa \sum_{\ell=0}^{\infty} \delta_{i\ell} u_{i,t-\ell},$$
(6)

where Z_t and u_{it} are unobserved serially uncorrelated aggregate and idiosyncratic errors and $\beta_{i\ell}$, $\gamma_{i\ell}$ and $\delta_{i\ell}$ are the unit-level responses to shocks X_t , Z_t and u_{it} at horizon ℓ . We denote $\beta_i = \{\beta_{i\ell}\}_{\ell=0}^{\infty}$, $\gamma_i = \{\gamma_{i\ell}\}_{\ell=0}^{\infty}$, $\delta_i = \{\delta_{i\ell}\}_{\ell=0}^{\infty}$ and $\theta_i = \{\mu_i, \beta_i, \gamma_i, \delta_i\}$. These are draws from a cross-sectional distribution and below we specify conditions so that the infinite sums in (6) are well defined with probability one.

Here, θ_i traces out cross-sectional heterogeneity in the responses to macro and micro shocks. The linear formulation (6) can accommodate flexible dynamic patterns, generating

¹⁸A proof and additional details are available upon request. As a practical example, we think of the responses of micro outcomes to monetary and fiscal policies during the COVID-19 pandemic.

rich heterogeneous versions of commonly used time series models (e.g., VARs; see Section 3.3). As a result, Y_{it} will often display non-trivial serial correlation. In fact, the infinite-order moving averages in (6) may arise from the Wold representation of more primitive, possibly serially correlated macro and micro state variables.

Access to external variables s_i allows the researcher to study the propagation of shocks along unit-level observables. Our premise is that there is usually more heterogeneity in θ_i than can be accounted for by s_i alone and our goal is to characterize estimation and inference in that context. As in Section 2, we consider a range of DGPs indexed by κ to cover different signal-to-noise environments. We also make the following assumptions:

Assumption 1 (Stationarity and iidness).

- (i) $\{X_t, Z_t, \{u_{it}\}_{i=1}^N\}_{t=-\infty}^\infty$ is stationary conditional on $\{\theta_i, s_i\}_{i=1}^N$.
- (ii) $\{\theta_i, s_i, \{u_{it}\}_{t=-\infty}^{\infty}\}_{i=1}^{N}$ is i.i.d. over i conditional on $\{X_t, Z_t\}_{t=-\infty}^{\infty}$.

Assumption 2 (Shocks and mean independence).

(i)
$$E[X_t \mid \{X_\tau\}_{\tau \neq t}, \{Z_\tau, \{u_{i\tau}\}_{i=1}^N\}_{\tau = -\infty}^{\infty}, \{\theta_i, s_i\}_{i=1}^N] = 0.$$

$$(ii) \ \ E\Big[\,Z_t \ \bigg| \ \{Z_\tau\}_{\tau \neq t}, \Big\{X_\tau, \{u_{i\tau}\}_{i=1}^N\Big\}_{\tau = -\infty}^\infty\,, \{\theta_i, s_i\}_{i=1}^N\,\Big] = 0.$$

$$(iii) \ E\Big[\,u_{it}\,\,\Big|\,\,\{u_{i\tau}\}_{\tau\neq t}, \big\{X_\tau,Z_\tau\big\}_{\tau=-\infty}^\infty\,, \theta_i,s_i\,\Big]=0.$$

Assumptions 1 and 2 extend S1 and S2 by allowing for unobserved heterogeneity and external covariates, while Equation (6) extends (2) by allowing for serial dependence in micro and aggregate components. Assumption 2 requires (θ_i , s_i) to be strictly exogenous with respect to shocks but, importantly, leaves their joint distribution (and that of { θ_i , s_i } $_{i=1}^N$ conditional on { X_t } $_{t=-\infty}^\infty$) unrestricted, as in pure fixed effects approaches.

Shock identification. Applied researchers aim at using local projections to learn about impulse responses β_{ih} . Interest in impulse responses naturally leads to the idea of shock articulated in Assumption 2(i) since then

$$\beta_{ih} = E[Y_{i,t+h} \mid X_t = 1, \theta_i] - E[Y_{i,t+h} \mid X_t = 0, \theta_i].$$

In that sense, Assumption 2 is an identifying condition without which the estimand of $\hat{\beta}(h)$ cannot be interpreted as a dynamic causal effect (Stock and Watson, 2018, Section 1.2).

A key question is how much information the researcher has about X_t . We distinguish three settings that cover all conventional macroeconometric shock identification procedures. To note, their use in the panel context entails no conceptual difference relative to the well-established time series case. The first settings is one where X_t is directly observed. Empirically, this is a very relevant case as it covers most of the applied papers surveyed in Supplemental Appendix ?? where X_t is often constructed by narrative or high-frequency methods.¹⁹ We begin our analysis by developing this setting.

The second setting is one where X_t can be recovered from a macro system via standard macroeconometric identification methods — for example, recursive, short-run, long-run, or heteroskedasticity-based methods (see Ramey (2016) for a review). The starting point is a SVAR for an n-vector \mathbf{R}_t of macro variables observable to the researcher,

$$\alpha_0 \mathbf{R}_t = \sum_{\ell=1}^{p_R} \alpha_\ell \mathbf{R}_{t-\ell} + \varepsilon_t, \tag{7}$$

where ε_t is an n-vector of shocks (in the sense of Assumption 2) whose j-th entry is the shock of interest $X_t = \varepsilon_{jt}$ and the remaining entries may overlap or be correlated with Z_t . Write $\alpha_{\ell,j\bullet}$ for the j-th row of α_{ℓ} . The identification method pins down $\{\alpha_{\ell,j\bullet}\}_{\ell=0}^{p_R}$ from the autocovariances of R_t so that knowledge of the latter delivers $X_t = \alpha_{0,j\bullet}R_t - \sum_{\ell=1}^{p_R}\alpha_{\ell,j\bullet}R_{t-\ell}$. More realistically, if the researcher does not know the population values $\{\alpha_{\ell,j\bullet}\}_{\ell=0}^{p_R}$ but has data $\{R_t\}_{t=1}^T$ (which may go beyond the span of the microdata), X_t can be estimated by

$$\hat{X}_t = \hat{\alpha}_{0,j\bullet} R_t - \sum_{\ell=1}^{p_R} \hat{\alpha}_{\ell,j\bullet} R_{t-\ell}$$

with $\hat{\alpha}_{\ell,j\bullet}$ a \sqrt{T} -consistent estimator of $\alpha_{\ell,j\bullet}$.²⁰ Our asymptotic analysis can accommodate the case where X_t in (1) is replaced by \hat{X}_t (or by $\hat{\alpha}_{0,j\bullet}R_t$) with a suitable choice of controls. For recursive identification, this strategy can be implemented by replacing X_t in (1) by the j-th entry of R_t controlling for $s_iR_{t-1},\ldots,s_iR_{t-p_R}$ and entries 1 to j-1 of s_iR_t . We return to this setting in Section 3.3 (see, in particular, Remark 8).

The third setting is one where the researcher only observes the proxy $X_t^* = X_t + v_t$ — a

¹⁹For example, Crouzet and Mehrotra (2020) and Ottonello and Winberry (2020) for monetary policy and Känzig (2021) for oil supply shocks.

²⁰For example, Drechsel (2023) follows this strategy to recover investment-specific technology shocks from a SVAR by imposing the restriction that X_t is the only driver in the long-run of the relative price of aggregate investment, which is the first entry of the vector \mathbf{R}_t .

measure of X_t contaminated with measurement error. This would apply, for instance, to the analysis of monetary policy when narrative or high-frequency approaches only yield an indirect estimate of exogenous shifts in policy. This is an empirically appealing premise which has not received much attention. In that case, it is more appropriate to think of X_t^* as an instrumental variable for an endogenous variable \tilde{X}_t (such as the policy rate) in a panel version of LP-IV with outcome $Y_{i,t+h}$ (see Ramey, 2016; Stock and Watson, 2018, for a treatment in the time series context). We study this setting in Section 3.4.

In the context of the empirical analysis of firm-level investment responses to monetary policy shocks (as in Section 5), the first setting covers panel local projections on direct measures of the shock, the second covers local projections on the interest rate controlling for other predetermined macro variables in R_t , and the third covers local projections on the interest rate instrumented by a proxy of the shock. In all of these settings, our framework allows for the inclusion of additional aggregate variables as controls which may account for the predictable part of the macro variables of interest.

3.1.1 Estimator and inference procedure

We now introduce the panel LP estimator and inference procedure. We denote by $W_{it} \in \mathbb{R}^d$ the vector of controls (d may change with the sample size). If W_{it} contains no time fixed effects, let $\hat{s}_i = s_i$ — this accommodates the case $s_i = 1$. Otherwise, let $\hat{s}_i = s_i - N^{-1} \sum_{j=1}^N s_j$ and note that if time fixed effects are included, local projections on $s_i X_t$ and $\hat{s}_i X_t$ produce numerically the same estimate $\hat{\beta}(h)$ below. In addition to unit and possibly time dummies, we consider below cases in which W_{it} contains lags of $s_i X_t$, Y_{it} and other macro and micro controls. The fitted equation for the panel LP estimator $\hat{\beta}(h)$ is

$$Y_{i,t+h} = \hat{\beta}(h)\hat{s}_i X_t + \hat{\eta}(h)' W_{it} + \hat{\xi}_{it}(h),$$

where the residual $\hat{\xi}_{it}(h)$ is orthogonal to $\hat{s}_i X_t$ and W_{it} . To characterize $\hat{\beta}(h)$ we use Frisch–Waugh–Lovell. Consider the auxiliary regression of $\hat{s}_i X_t$ on W_{it} ,

$$\hat{s}_{i}X_{t} = \hat{\pi}(h)'W_{it} + \hat{x}_{it}(h), \tag{8}$$

where the residual $\hat{x}_{it}(h)$ is orthogonal to W_{it} . Then, an explicit formula for $\hat{\beta}(h)$ is

$$\hat{\beta}(h) = \frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{x}_{it}(h) Y_{i,t+h}}{\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{x}_{it}(h)^{2}}.$$
(9)

The time-clustered heteroskedasticity-robust standard error is

$$\hat{\sigma}(h) = \sqrt{\frac{\hat{V}(h)}{(T-h)\hat{J}(h)^2}},\tag{10}$$

with

$$\hat{J}(h) = \frac{1}{N(T-h)} \sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{x}_{it}(h)^2, \quad \hat{V}(h) = \frac{1}{(T-h)} \sum_{t=1}^{T-h} \left(\frac{1}{N} \sum_{i=1}^{N} \hat{x}_{it}(h) \hat{\xi}_{it}(h) \right)^2.$$
 (11)

Finally, the $(1 - \alpha)$ confidence interval is

$$\hat{C}_{\alpha}(h) = \left[\hat{\beta}(h) \pm z_{1-\alpha/2}\hat{\sigma}(h)\right],\tag{12}$$

where z_q is the q-quantile of the standard normal distribution.

3.1.2 Additional assumptions

To establish our uniform asymptotic approximations, we need the following:

Assumption 3 (Regularity conditions).

(i) There is a positive finite constant M_8 such that, almost surely,

$$E\left[\left.X_{t}^{8} \mid \{\theta_{i}, s_{i}\}_{i=1}^{N}\right.\right] \leq M_{8}, \quad E\left[\left.Z_{t}^{8} \mid \{\theta_{i}, s_{i}\}_{i=1}^{N}\right.\right] \leq M_{8}, \quad E\left[\left.u_{it}^{8} \mid \theta_{i}, s_{i}\right.\right] \leq M_{8}.$$

(ii) There is a positive finite constant \underline{M} such that, almost surely,

$$E\left[\left.X_{t}^{2} \mid \{X_{\tau}\}_{\tau \neq t}, \{\theta_{i}, s_{i}\}_{i=1}^{N}\right.\right] \geq \underline{M}, \quad E\left[\left.Z_{t}^{2} \mid \{X_{\tau}\}, \{\theta_{i}, s_{i}\}_{i=1}^{N}\right.\right] \geq \underline{M}, \quad E\left[\left.u_{it}^{2} \mid \{X_{\tau}\}, \theta_{i}, s_{i}\right.\right] \geq \underline{M}.$$

- (iii) The conditional cumulants up to fourth-order of $\text{vec}\{(X_t, Z_t, u_{it})(X_t, Z_t, u_{it})'\}$ given $\{\theta_i, s_i\}_{i=1}^N$ are almost surely absolutely summable.
- (iv) There are positive finite constants C_{ℓ} such that $C = \sum_{\ell=0}^{\infty} C_{\ell} < \infty$ and, almost surely,

$$|\beta_{i\ell}| \le C_{\ell}$$
, $|\gamma_{i\ell}| \le C_{\ell}$, $|\delta_{i\ell}| \le C_{\ell}$, $|s_i| < C$.

(v) There is a positive finite constant \underline{C} such that, almost surely,

$$\sum_{\ell=0}^{\infty} \left(N^{-1} \sum_{i=1}^{N} \hat{s}_i \beta_{i\ell} \right)^2 \geq \underline{C}, \quad \sum_{\ell=0}^{\infty} \left(N^{-1} \sum_{i=1}^{N} \hat{s}_i \gamma_{i\ell} \right)^2 \geq \underline{C}, \quad N^{-1} \sum_{\ell=0}^{\infty} \sum_{i=1}^{N} \hat{s}_i^2 \delta_{i\ell}^2 \geq \underline{C}.$$

Our model interprets θ_i as unit-specific parameters and $\{X_t, Z_t, u_{it}\}$ as sources of uncertainty. This calls for making time series assumptions on the uncertainty given parameters (parts (i), (ii) and (iii)) while requiring that parameters ensure sufficient regularity for all units in the cross-sectional population (parts (iv) and (v)).

Parts (i), (ii) and (iii) are standard in the time series context (see, for instance, Assumption 2 in Montiel Olea and Plagborg-Møller (2021)). They put limits on the tails of the distributions of shocks, as well as the predictability and dependence of their second moments. Part (iv), on the other hand, guarantees that infinite moving averages, such as $\sum_{\ell=0}^{\infty} \beta_{i\ell} X_{t-\ell}$, are well defined for all units. Absolute summability rules out unit roots but still allows for rich persistence patterns — such as those from stationary ARMA and other short-memory processes.²¹

Lastly, part (v) requires non-zero variability given $\{\theta_i, s_i\}_{i=1}^N$ of $N^{-1} \sum_{i=1}^N \hat{s}_i \sum_{\ell=0}^\infty \beta_{i\ell} X_{t-\ell}$, $N^{-1} \sum_{i=1}^N \hat{s}_i \sum_{\ell=0}^\infty \gamma_{i\ell} Z_{t-\ell}$ and $N^{-1/2} \sum_{i=1}^N \hat{s}_i \sum_{\ell=0}^\infty \delta_{i\ell} u_{i,t-\ell}$. It is mostly a technical condition to prevent trivial cases in which the regression score has zero variance. Nevertheless, it is compatible with, say, a non-negligible fraction of units having zero exposure to macro or micro shocks. It also places no restriction on the relative importance of macro versus micro shocks which is governed by κ .

3.2 Main result

The main contribution of the paper is to characterize the large-sample properties of $\hat{\beta}(h)$, $\hat{\sigma}(h)$ and $\hat{C}_{\alpha}(h)$. In the asymptotic plan, we take $T, N \to \infty$ and we are interested in uniform approximations with respect to κ . The key result is Proposition 1 which states that $\hat{C}_{\alpha}(h)$ delivers uniformly valid inference for the coefficient in a regression of β_{ih} on \hat{s}_i if enough lags of $\hat{s}_i X_t$ are used as controls.

We describe first the estimand and then the uniform inference result. We use P_{κ} to indicate probabilities under a DGP with a given value of κ and we omit the subindex from

²¹We conjecture, however, that many of our results remain valid at moderate horizons in the presence of near unit roots and our simulation evidence supports this claim. See Section 3.3 for further discussion.

those objects that do not depend on κ (such as those in Assumptions 2 and 3).

Estimand. If s_i is not a constant and time fixed effects are included, the population object targeted by the panel LP is

$$\beta(h) = \frac{\operatorname{Cov}(s_i, \beta_{ih})}{\operatorname{Var}(s_i)}.$$
(13)

In other words, panel LPs estimate the slope in a population linear projection of β_{ih} on characteristics s_i including an intercept. Similarly, if $s_i = 1$, the estimand becomes the mean impulse response $\beta(h) = E[\beta_{ih}]$. Note that omitting either X_t or time dummies as controls in a panel LP has the effect of forcing the regression of β_{ih} on s_i through the origin, leading to the estimand $\beta(h) = (E[s_i^2])^{-1}E[s_i\beta_{ih}]$. In order to obtain a rich summary of the heterogeneity in β_{ih} , therefore, the researcher will typically need to explore different choices of s_i or allow s_i to be a vector.²²

Under the conditions of Proposition 1, $\hat{\beta}(h) = \beta(h) + o_{P_{\kappa}}(1)$ for any DGP sequence P_{κ} such that $\kappa / \sqrt{TN} = o(1)$: that is, if the panel LP estimator converges, it is to $\beta(h)$.

This clarifies the sense in which panel LPs can be interpreted when the underlying population of interest features unrestricted heterogeneity in responses to shocks, as in (6). Precisely because we place virtually no restriction on the joint distribution of (θ_i, s_i) , the characterization of the estimand is of a nonparametric nature.

Uniformly valid inference. Let p be the number of lags of $\hat{s}_i X_t$ included in the controls W_{it} . Both p and h are fixed as $T, N \to \infty$ while $T/N \to 0$.²³ Our main result is that $\hat{C}_{\alpha}(h)$ has correct coverage for $\beta(h)$ uniformly over κ so long as $h \le p$:

Proposition 1. *Under Assumptions* 1, 2 and 3, for $h \le p$,

$$\lim_{T,N\to\infty} \sup_{\kappa} \left| P_{\kappa} \Big(\beta(h) \in \hat{C}_{\alpha}(h) \Big) - (1-\alpha) \right| = 0.$$
 (14)

²²For example, the best linear approximation $E^*[\beta_{ih} \mid s_i] = E[\beta_{ih}] + (\text{Cov}(s_i, \beta_{ih}) / \text{Var}(s_i))(s_i - E[s_i])$ requires both estimands or, alternatively, the interaction of X_t with $(1, s_i)$ rather than s_i alone (omitting time effects). If s_i is multivariate, a confidence region constructed on the basis of a time-clustered heteroskedasticity-robust variance estimate enjoys the same uniform validity property of Proposition 1. We illustrate this in our empirical calculations in Section 5.

²³We regard $T/N \to 0$ as a mild requirement for the empirical applications of reference. It follows from the proof of Proposition 1 that if T/N is not asymptotically negligible (as if taking N as fixed), (14) holds with $\beta(h)$ replaced by the *finite-population* estimand $\tilde{\beta}(h) = (\sum_{i=1}^{N} \hat{s}_i^2)^{-1} \sum_{i=1}^{N} \hat{s}_i \beta_{ih}$.

Proposition 1 states that valid inference follows from clustering standard errors at the time level, which accounts for cross-sectional dependence induced by omitted aggregate shocks, and from ex-ante including lags of $\hat{s}_i X_t$ as controls, which renders the regression scores uncorrelated. We refer to this strategy as time-clustered lag-augmented heteroskedasticity-robust (t-LAHR) inference. As in Section 2 and as explained below, it is closely linked to inference in time series LPs.

Despite the general error dynamics in (6), the regression score $\sum_{i=1}^{N} X_i \hat{s}_i \xi_{it}(h, \kappa)$, with $\xi_{it}(h, \kappa)$ the population counterpart to $\hat{\xi}_{it}(h)$ defined in (19), has limited serial correlation. It is an MA(h) process with the first p autocovariances set to zero. Thus, it becomes uncorrelated when $p \ge h$ which is why t-LAHR works. Besides, when p < h, the autocovariances stem only from leftover leads of X_t and not from the unobserved macro error Z_t or micro error u_{it} . In fact, they will tend to be small compared to the variance of the score in low-signal (large κ) DGPs or if $\beta_{i\ell}$ decays quickly. We therefore expect t-LAHR inference to have small coverage distortions even for p < h; we provide affirmative evidence via simulations in Section 4.

A striking implication of Proposition 1 is that t-LAHR inference remains valid even in the low-signal setting $\kappa/\sqrt{N} \to \infty$ where there is scarcity of information about aggregate shocks in the sample and $\hat{\beta}(h)$ is inconsistent. The uniformity over DGPs with different macro-micro signal-noise obviates the need to take a stand on the κ -regime, which is important because κ is not always consistently estimable.

In contrast, inference based on unit-level clustering of the regression score is not uniformly valid as it tends to severely undercover $\beta(h)$ in high- and moderate-signal regimes. Similarly to Section 2, provided lags of $\hat{s}_i X_t$ are included, unit-level clustering is asymptotically equivalent to not clustering at all, whereas two-way clustering is equivalent to time-level clustering. That is, unit-level clustering is neither necessary nor sufficient for valid inference — yet another implication of X_t being a shock that has no counterpart in a more generic time series setup.

Remark 3 (Proof steps). To establish (14), we decompose the problem into showing (A) asymptotic normality of the score, (B) consistency of the standard error, and (C) negligibility of some remainder terms. We obtain uniformity via the drifting parameter sequence approach (see Andrews, Cheng, and Guggenberger (2020)).

In (A), although the regression score is serially uncorrelated, it contains leads and lags of macro and micro errors. This makes the reverse martingale technique of Montiel Olea

and Plagborg-Møller (2021) inapplicable. Instead, using a similar insight to Xu (2023), we produce a martingale approximation by rearranging the score so that leads at time t become lags at a time in the future of t (Lemma ?? in Supplemental Appendix ??).

In (B) and (C), we rely on direct calculation of uniform bounds. The presence of heterogeneity poses a challenge which has no counterpart in the time series case. Because of Assumption 3, we can derive many of the bounds by first conditioning on $\{\theta_i, s_i\}_{i=1}^N$, exploiting the connection between conditional and unconditional convergence.

Remark 4 (Synthetic time series). A useful device to interpret panel LPs is the following representation. The residual $\hat{x}_{it}(h)$ in (8) can be written as $\hat{x}_{it}(h) = \hat{s}_i \hat{X}_t(h)$, where $\hat{X}_t(h)$ is the residual from regressing X_t on X_{t-1}, \ldots, X_{t-p} and an intercept (on T-h observations).

Then, the panel LP estimator in (9) can be written as

$$\hat{\beta}(h) = \frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{s}_{i} \hat{X}_{t}(h) Y_{i,t+h}}{\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{s}_{i} \hat{X}_{t}(h)^{2}} = \frac{\sum_{t=1}^{T-h} \hat{X}_{t}(h) \hat{Y}_{t+h}}{\sum_{t=1}^{T-h} \hat{X}_{t}(h)^{2}},$$

i.e., the time series LP estimator that regresses cross-sectional regression coefficients $\hat{Y}_{t+h} = (\sum_{i=1}^N \hat{s}_i^2)^{-1} \sum_{i=1}^N \hat{s}_i Y_{i,t+h}$ on X_t controlling for $1, X_{t-1}, \ldots, X_{t-p}$. The standard error $\hat{\sigma}(h)$ in (10) is also the Eicker–Huber–White standard error calculated on the time series LP residuals $\hat{\xi}_t(h) = (\sum_{i=1}^N \hat{s}_i^2)^{-1} \sum_{i=1}^N \hat{s}_i \hat{\xi}_{it}(h)$. Hence, t-LAHR inference for panel LPs and lag-augmented heteroskedasticity-robust inference for time series LPs are intimately related.

Remark 5 (s_i and precision). This representation is also useful to illuminate the fact that estimation error is of order $T^{-1/2}$ in environments with $\kappa \propto \sqrt{N}$, despite what otherwise looks like a standard panel regression with potentially very rich microdata. We can give interpretable conditions under which variation in s_i affords faster convergence rates. These are akin to s_i being a cross-sectional instrument: we require s_i to correlate with β_{ih} — that is, be relevant for heterogeneity in transmission of X_t at horizon h — but to be orthogonal to all other exposures to aggregate shocks, ($\{\beta_{i\ell}\}_{\ell \neq h}$, γ_i). These conditions seem particularly hard to meet: for each horizon h, a source of variation that is orthogonal to responses at all other horizons is required. (Assumption 3(v) rules this out in our formulation.) In some sense, this reveals an intrinsic trade-off between documenting interesting transmission mechanisms and finding valid instruments for precision.

To see this, note that $\hat{x}_{it}(h)$ is $\hat{s}_i X_t$ minus a linear combination of $\hat{s}_i X_{t-1}, \dots, \hat{s}_i X_{t-p}$ and unit and possibly time indicators which is orthogonal to all of the latter. When W_{it} includes additional controls, the synthetic time series representation is asymptotically but not numerically equivalent. Arkhangelsky and Korovkin (2023) derive a similar connection in a regional-exposure design context.

Remark 6 (*t***-HAR).** In principle, time-clustered HAR inference is a valid alternative to *t*-LAHR. An analogue to Proposition 1 can be shown for a confidence interval that replaces $\hat{V}(h)$ in (11) with the Hansen and Hodrick (1980) variance estimator $\hat{V}(h) + 2\sum_{\ell=p+1}^{h} \tilde{V}_{\ell}(h)$ where

$$\tilde{V}_{\ell}(h) = \frac{1}{(T-h)} \sum_{t=\ell+1}^{T-h} \left(\frac{1}{N} \sum_{i=1}^{N} \hat{x}_{it}(h) \hat{\xi}_{it}(h) \right) \left(\frac{1}{N} \sum_{i=1}^{N} \hat{x}_{i,t-\ell}(h) \hat{\xi}_{i,t-\ell}(h) \right),$$

This boils down to $\hat{V}(h)$ for $p \ge h$. Unlike $\hat{V}(h)$, this alternative variance estimator is not guaranteed to be positive semidefinite. Also, *t*-LAHR inference is simpler to implement and refine, remains tractable over moderate horizons under VAR DGPs (Section 3.3), and performs better in small samples (Section 4).

Remark 7 (State-dependence). In some applications, interest is in the differential pass-through of shocks to responses along an observable (time-varying) state, denoted now s_{it} . Formalizing this requires extending (6) to allow for time-varying impulse responses:

$$Y_{it} = \mu_i + \sum_{\ell=0}^{\infty} \beta_{it\ell} X_{t-\ell} + v_{it}, \quad v_{it} = \sum_{\ell=0}^{\infty} \gamma_{it\ell} Z_{t-\ell} + \kappa \sum_{\ell=0}^{\infty} \delta_{it\ell} u_{i,t-\ell}.$$

Letting $\hat{s}_{it} = s_{it} - N^{-1} \sum_{j=1}^{N} s_{jt}$, the corresponding panel LP estimator on $\hat{s}_{it} X_t$ retains its interpretation as the slope coefficient of the linear projection $E^*[\beta_{ith} \mid s_{it}]$ as long as s_{it} and impulse responses are exogenous with respect to X_t . Although a more detailed exploration is beyond the scope of our paper, the treatment of s_{it} is analogous to that of s_i , and all the results above carry over with little modification. We revisit this in simulations in Section 4 and in our empirical illustration in Section 5.²⁵

3.3 Panel VAR model

It is not uncommon in applications that the researcher is interested in responses at an horizon h which is a non-negligible fraction of T. Proposition 1 guarantees exact coverage for short horizons depending on the number of lags of the outcome and shock used as controls. There is, however, one important class of DGPs for which our uniformity result extends to $h \propto T$: the VAR class.

²⁵Rambachan and Shephard (2021, Section 3.4) offer a nonparametric characterization of local projection estimands when states are endogenous in a time-series potential outcomes framework; see also Gonçalves, Herrera, Kilian, and Pesavento (forthcoming) for the case where $s_t = 1\{X_t > c\}$.

We now assume a panel VAR(p) model for (Y_{it} , X_t) (with $p < \infty$):

$$Y_{it} = m_i + \sum_{\ell=1}^{p} A_{\ell} Y_{i,t-\ell} + \sum_{\ell=0}^{p} B_{i\ell} X_{t-\ell} + C_{i0} Z_t + \kappa D_{i0} u_{it}.$$
 (15)

If $\sum_{\ell=1}^{p} A_{\ell} < 1$, as implied by Assumption 3(iv), we can recover unit-specific parameters μ_{i} , $\{\beta_{i\ell}\}$, $\{\gamma_{i\ell}\}$, $\{\delta_{i\ell}\}$ from m_{i} , $\{A_{\ell}\}$, $\{B_{i\ell}\}$, C_{i0} , D_{i0} inverting the polynomial $A(L) = 1 - \sum_{\ell=1}^{p} A_{\ell} L^{\ell}$. That is, the VAR model (15) is a special case of (6).

Assuming that p is known and that W_{it} contains p lags of Y_{it} and s_iX_t , the t-LAHR confidence interval $\hat{C}_{\alpha}(h)$ defined in (12) has uniform validity even for moderately long horizons h exceeding p:

Proposition 2. Under Assumptions 1, 2 and 3, for some positive constant $\phi < 1$,

$$\lim_{T,N\to\infty} \sup_{0\le h\le \phi T} \sup_{\kappa} \left| P_{\kappa} \Big(\beta(h) \in \hat{C}_{\alpha}(h) \Big) - (1-\alpha) \right| = 0.$$
 (16)

Proof. See Appendix A.

The intuition and proof for Proposition 2 mirror that of Proposition 1. In the VAR model (15), the regression score $\sum_{i=1}^{N} X_t \hat{s}_i \xi_{it}(h, \kappa)$, with $\xi_{it}(h, \kappa)$ now defined in (20), is serially uncorrelated not just for $h \leq p$ but for any h. The basic consequence is that if a low-order VAR model is a reasonable approximation, the t-LAHR inference approach that relies on controlling for a small number of lags of the outcome and shock is robust over long horizons and regardless of the amount of micro noise. ²⁶

Remark 8 (LP inference when the shock comes from SVAR). Proposition 2 can be read as the panel data counterpart to the result in Montiel Olea and Plagborg-Møller (2021) under stationarity when the shock is observed. That parallel implies that if instead we observe a serially correlated aggregate $X_t^* = \sum_{\ell=1}^p \alpha_\ell X_{t-\ell}^* + X_t$ and we run a lag-augmented local projection of $Y_{i,t+h}$ on $s_i X_t^*$ including p lags of Y_{it} and $s_i X_t^*$ in the control vector W_{it} , t-LAHR inference is again uniformly valid over h and κ .

²⁶The results in Montiel Olea et al. (2024) suggest that for a *fixed* horizon h, t-LAHR inference would also remain valid if the VAR model (15) were contaminated by moving averages of Z_t and u_{it} in a $T^{-1/4}$ -neighborhood of zero — that is, if the VAR model holds only approximately. The simulation evidence in Section 4 based on DGPs which are not VARs is consistent with this idea.

²⁷The connection with Montiel Olea and Plagborg-Møller (2021) also suggests that $\hat{C}_{\alpha}(h)$ is uniformly valid over the VAR parameter space (including unit roots) if a certain condition on uniform non-singularity of the least squares denominator matrix (Assumption 3 in their paper) holds.

An important extension is recursive identification in a macro system such as (7), where the researcher observes $\mathbf{R}_t = (R_{1t}, \dots, R_{nt})'$ and X_t is the j-th entry of ε_t in the SVAR model, so that $R_{jt} = \sum_{k=1}^{j-1} \alpha_{0,jk} R_{kt} + \sum_{\ell=1}^{p_R} \alpha_{\ell,j\bullet} \mathbf{R}_{t-\ell} + X_t$. Consider the panel local projections

$$Y_{i,t+h} = \hat{\beta}_{R}(h)s_{i}R_{jt} + \hat{\eta}_{R}(h)'W_{it} + \hat{\xi}_{R,it}(h) \text{ and } Y_{i,t+h} = \hat{\beta}_{X}(h)s_{i}\hat{X}_{t} + \hat{\eta}_{X}(h)'W_{it} + \hat{\xi}_{X,it}(h),$$

where W_{it} includes $s_i R_{1t}, \ldots, s_i R_{j-1,t}$, p_R lags of $s_i R_t$, p lags of Y_{it} , and unit and time effects, and where \hat{X}_t is the estimate of X_t from Section 3.1. Then, the analogue to $\hat{C}_{\alpha}(h)$ that uses $\hat{\beta}_R(h)$ or $\hat{\beta}_X(h)$ with t-LAHR standard errors also enjoys the uniform validity result (16).²⁸

Remark 9 (Heterogeneity in VAR coefficients). Model (15) assumes homogeneous coefficients $\{A_\ell\}$. This is common in the microeconometric literature on panel VARs (Arellano, 2003, Chapter 6) but it is not necessary for (16). We can prove Proposition 2 in a moderate heterogeneity environment that replaces A_ℓ with $A_{i\ell}$ where $\sup_{1 \le i \le N} |A_{i\ell} - A_\ell| = O_p(T^{-1/2})$. Proposition 2 can also be established (under different regularity conditions) if we allow for heterogeneity in $\{A_\ell\}$ but include p unit-specific lags of Y_{it} as controls in W_{it} .

3.4 Panel LP-IV and proxy shocks

The most common implementation of panel LPs in empirical work treats the shock of interest as observed. Nevertheless, it is sometimes more realistic to assume there is measurement error in the shock elicitation process. This creates an endogeneity problem that can be dealt with by using the shock measures as instruments for the actual underlying shock (Ramey, 2016; Stock and Watson, 2018).

The researcher observes the outcome Y_{it} and characteristics s_i , but instead of the actual shock X_t she observes an endogenous aggregate state variable \tilde{X}_t and a proxy shock X_t^* . In the context of our empirical analysis, \tilde{X}_t denotes the Fed Funds rate and X_t^* an imperfect measurement of monetary policy surprises X_t . In addition to (6), we assume

$$\tilde{X}_{t} = \sum_{\ell=0}^{\infty} b_{\ell} X_{t-\ell} + \sum_{\ell=0}^{\infty} c_{\ell} Z_{t-\ell}, \tag{17}$$

$$X_t^* = a_0 X_t + \nu_t, (18)$$

In the local projection on $s_i\hat{X}_t$, controlling for $s_iR_{1t},\ldots,s_iR_{j-1,t}$ and lags of s_iR_t ensures that the generated regressor error $X_t - \hat{X}_t$ does not affect inference asymptotically by orthogonalizing the residual $\hat{\xi}_{X,it}(h)$ with respect to $s_i(X_t - \hat{X}_t)$. The same applies to the other macro identification methods discussed in Section 3.1.

where v_t is measurement error. We normalize $b_0 = 1$ to fix the scale of the estimand as only relative impulse responses are identified.²⁹ We also adopt the following:

Assumption 4 (LP-IV).

- (*i*) $a_0 \neq 0$.
- (ii) Assumptions 1, 2 and 3 hold with Z_t replaced by (Z_t, v_t) .
- (iii) For the same constants C_{ℓ} and C of Assumption 3,

$$|b_{\ell}| \le C_{\ell}, \quad |c_{\ell}| \le C_{\ell}, \quad \sum_{\ell=0}^{\infty} b_{\ell}^2 \ge \underline{C}, \quad \sum_{\ell=0}^{\infty} c_{\ell}^2 \ge \underline{C}.$$

Assumption 4(i) is needed for instrument relevance, and we restrict our attention to the strong instrument case where we keep a_0 fixed as $N, T \to \infty$. On the other hand, Assumption 4(ii) implies that ν_t is orthogonal to $\{X_\tau, Z_\tau\}$. This embodies the key lead-lag exogeneity condition requiring X_t^* to be contemporaneously correlated only with X_t , a well-known condition in the time series LP-IV context.³⁰ Finally, Assumption 4(iii) imposes regularity on the endogenous variable \tilde{X}_t .

LP-IV estimation and inference. LP-IV regresses $Y_{i,t+h}$ on $\tilde{X}_t = (\tilde{X}_t, \tilde{X}_{t-1}, \dots, \tilde{X}_{t-p})'$ using $X_t^* = (X_t^*, X_{t-1}^*, \dots, X_{t-p}^*)'$ as instruments (both interacted with s_i), controlling for unit and time effects (W_{it} denotes controls). The residualized instrument is

$$\hat{x}_{it}(h) = \hat{s}_i X_t^* - \hat{\pi}(h)' W_{it} = \hat{s}_i \hat{X}_t^*(h),$$

where $\hat{X}_t^*(h) = X_t^* - (T - h)^{-1} \sum_{t=1}^{T-h} X_t^*$. The panel LP-IV estimator $\hat{\beta}^{\text{IV}}(h)$ is then

$$\hat{\boldsymbol{\beta}}^{\text{IV}}(h) = \left(\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{\boldsymbol{x}}_{it}(h) \hat{s}_{i} \tilde{\boldsymbol{X}}_{t}'\right)^{-1} \sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{\boldsymbol{x}}_{it}(h) Y_{i,t+h} = \left(\sum_{t=1}^{T-h} \hat{\boldsymbol{X}}_{t}^{*}(h) \tilde{\boldsymbol{X}}_{t}'\right)^{-1} \sum_{t=1}^{T-h} \hat{\boldsymbol{X}}_{t}^{*}(h) \hat{Y}_{i,t+h}'$$

²⁹It is straightforward to include intercepts in both (17) and (18). Additionally, as in Section 3.3, we can derive uniformity results with respect to the horizon h by assuming a VAR model in (6) and (17).

³⁰See, for instance, Stock and Watson (2018, p. 924) and Plagborg-Møller and Wolf (2021, p. 970). The setup can be extended to allow v_t to be serially correlated and to the case where X_t^* is valid only after conditioning on a set of controls.

where $\hat{Y}_{i,t+h}$ is the synthetic outcome defined in Remark 4. Put another way, panel LP-IV admits a synthetic time series LP-IV representation.

The only entry of $\hat{\beta}^{IV}(h)$ that has interpretation as an estimate of a relative impulse response is $\hat{\beta}_0^{IV}(h) = e_1'\hat{\beta}^{IV}(h)$ where e_1 is the first column of I_{p+1} . The remaining entries are necessary for t-LAHR inference to be valid. Given residuals

$$\hat{\xi}_{it}^{\mathrm{IV}}(h) = Y_{i,t+h} - \hat{s}_i \tilde{\boldsymbol{X}}_t' \hat{\boldsymbol{\beta}}^{\mathrm{IV}}(h) - \hat{\eta}^{\mathrm{IV}}(h)' W_{it},$$

we define

$$\hat{J}^{\text{IV}}(h) = \frac{1}{N(T-h)} \sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{x}_{it}(h) \hat{s}_i \tilde{X}'_t, \quad \hat{V}^{\text{IV}}(h) = \frac{1}{(T-h)} \sum_{t=1}^{T-h} \left(\frac{1}{N} \sum_{i=1}^{N} \hat{x}_{it}(h) \hat{\xi}_{it}^{\text{IV}}(h) \right)^2.$$

The time-clustered heteroskedasticity-robust standard error for $\hat{\beta}_0^{\text{IV}}(h)$ is

$$\hat{\sigma}_0^{\text{IV}}(h) = \left[\frac{1}{(T-h)} \cdot \left(e_1' \hat{\boldsymbol{J}}^{\text{IV}}(h)^{-1}\right) \hat{\boldsymbol{V}}^{\text{IV}}(h) \left(e_1' \hat{\boldsymbol{J}}^{\text{IV}}(h)^{-1}\right)'\right]^{1/2}$$

and the $(1 - \alpha)$ confidence interval, $\hat{C}_{\alpha}^{IV}(h) = \left[\hat{\beta}_{0}^{IV}(h) \pm z_{1-\alpha/2}\hat{\sigma}_{0}^{IV}(h)\right]$. Then:

Proposition 3. *Under Assumption 4, for* $h \le p$ *,*

$$\lim_{T,N\to\infty}\sup_{\kappa}\left|P_{\kappa}\left(\beta(h)\in\hat{C}_{\alpha}^{IV}(h)\right)-(1-\alpha)\right|=0.$$

Proof. See Appendix A.

Remark 10 (Absence of first-stage heterogeneity). The LP-IV estimand coincides (under the normalization $b_0 = 1$) with the LP estimand (13) despite the presence of heterogeneity. This is far from obvious: under treatment effect heterogeneity, IV estimands are generally (weighted averages of) local average treatment effects (Angrist and Imbens, 1995; Angrist, Imbens, and Graddy, 2000). This is caused by the aggregate-only nature of the first-stage model, yet another illustration of the unique setting that we study in this paper.

4 Simulation study

We ran a comprehensive simulation study to verify the finite-sample robustness of the inference procedures analyzed in Section 3. Here we provide a summary and defer

additional detail and results to Supplemental Appendix ??.

Designs. Our study relies on two different DGPs. The first is the general setup (6) supplemented with (17)–(18) to cover the endogenous case. We begin by simulating shocks $\{X_t, Z_t, v_t, \{u_{it}\}_{i=1}^N\}$ as mutually and serially independent N(0,1) random variables, and by drawing $\{\theta_i, s_i\}_{i=1}^N$ independently across units. To ensure correlation between observed and unobserved heterogeneity we use a technique described in Supplemental Appendix **??**. We calibrate the distribution of $\{\beta_{i\ell}, \gamma_{i\ell}, \delta_{i\ell}\}$ and the value of $\{b_\ell, c_\ell\}$ to produce realistic degrees of shock persistence.

Given these elements, we generate the inputs for panel LP and LP-IV procedures, namely Y_{it} , X_t , S_t , X_t^* . We also simulate the time-varying covariate $S_{it} = S_t + \zeta_{it}$ (where ζ_{it} is such that S_{it} remains strictly exogenous) to compare panel LPs on $S_t X_t$ and $S_{it} X_t$ —this illustrates the point we made in Remark 7.

The second DGP is the VAR model (15). Again we generate shocks as i.i.d. N(0,1) and we simulate the heterogeneity as detailed in Supplemental Appendix ??. When calibrating the VAR parameters $\{A_\ell\}$ we allow the largest AR root to be 1-c/T (we use c=5) to capture the essence of a near non-stationary environment.³¹

The results below are based on $n_{\rm MC}=5,000$ Monte Carlo samples. Motivated by our survey of the empirical literature, we look at designs with T=30 and T=100. We set N=1,000 (although we also considered experiments with larger N) and we let κ take values consistent with $\bar{R}^2(\kappa) \in \{0.99,0.66,0.33\}$ as defined in (4). As a reference, $\bar{R}^2(\kappa)=0.66$ corresponds to the one-third of aggregate fluctuations explained by micro shocks suggested by Gabaix (2011) for GDP growth, which we take as moderate signal-to-noise.

Inference procedures. We compare t-LAHR inference with one-way (1W), two-way (2W), and Driscoll-Kraay (DK98) inferences. These are implemented without lag augmentation, as is common practice. For illustrative purposes, we also include t-HR (the non-lag-augmented counterpart to t-LAHR) and t-HAR alternatives.

For *t*-LAHR inference we use the simple lag selection rule $p = \min\{h, (T-h)^{1/3}\}$ (except in the VAR DGP where *p* is known) and we apply the finite-sample refinement advocated

³¹We also considered experiments where (a) in the first DGP shocks are conditionally heteroskedastic, and (b) in the VAR DGP we have unit-specific VAR parameters $\{A_{i\ell}\}$. We did not find any major difference with what we report here.

by Imbens and Kolesár (2016). The lag selection rule is motivated by Xu (2023, Section 3.3) for fixed h and provides fairly generous lag augmentation. For t-HAR inference we use the equally-weighted cosine approach (Müller, 2004) with the choice of tuning parameter recommended in Lazarus et al. (2018).

Results. In Figure 1, we report pointwise coverage rates for horizons $0 \le h \le 0.25T$ with T = 100. These correspond to 90% confidence intervals for panel LP and LP-IV using s_i to interact the aggregate shock. Panels (a)-to-(c) display LP while (d)-to-(f) display LP-IV in the general DGP; panels (g)-to-(i) display LP in the VAR DGP.

Figure 1 suggests four takeaways. First, t-LAHR performs best in all scenarios, with coverage close to the nominal rate even in low-signal cases and for horizons h well beyond p. Its mean absolute coverage distortion never exceeds 2%, whereas it is between 4% and 7% for the second best option (t-HAR) under high signal.

Second, estimating the long-run variance of the score (instead of lag augmenting) can be challenging with small *T*. This is particularly true for DK98 which relies on Newey–West. Interestingly, these approaches do better in low-signal DGPs where, as mentioned before, there is less to gain from doing HAC.

Third, one-way clustering is very sensitive to $\bar{R}^2(\kappa)$, suffering severe distortions in intermediate- and high- $\bar{R}^2(\kappa)$ cases. What is more, it is outperformed by *t*-LAHR even if micro shocks explain the majority of aggregate variation. This is consistent with the view that 1W guards against the wrong type of correlation in the score.

Finally, two-way clustering is usually close to *t*-HR, its non-*i*-clustered version; another indication that there is no clear advantage in clustering by units. In fact, in certain occasions (mainly low-signal and near non-stationary designs), 2W gives worse inferences than *t*-HR or 1W alone. This is possibly due to the non-standard behavior of variance estimators when there are micro (near) unit roots.

Identical takeaways emerge in experiments where we substitute s_i with either 1 or s_{it} (Supplemental Appendix ??), and with a sample size T = 30 (Figure 2).

In sum, the small-sample evidence reinforces many of our theoretical results. It shows that the large-sample approximations of Section 3 provide reliable guidance for understanding estimation and inference with aggregate shocks. Furthermore, it illustrates the practical relevance of achieving uniformity with respect to κ , and it delivers a clear methodological prescription: t-LAHR inference.

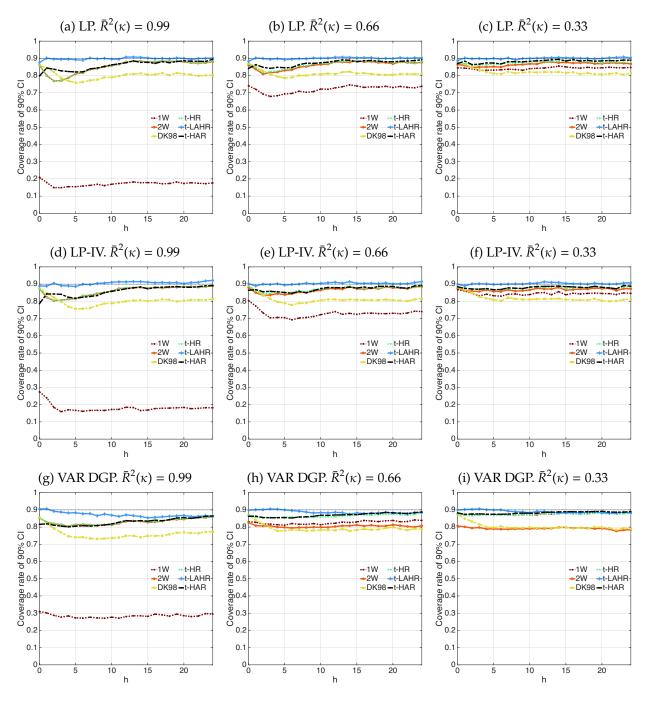


FIGURE 1. Coverage rates of 90% confidence intervals for T = 100.

Note: 1W refers to one-way (unit-level) clustering, 2W to two-way clustering, DK98 to Driscoll–Kraay, and *t*-HR/*t*-LAHR/*t*-HAR to the time-level clustering approaches discussed in the text.

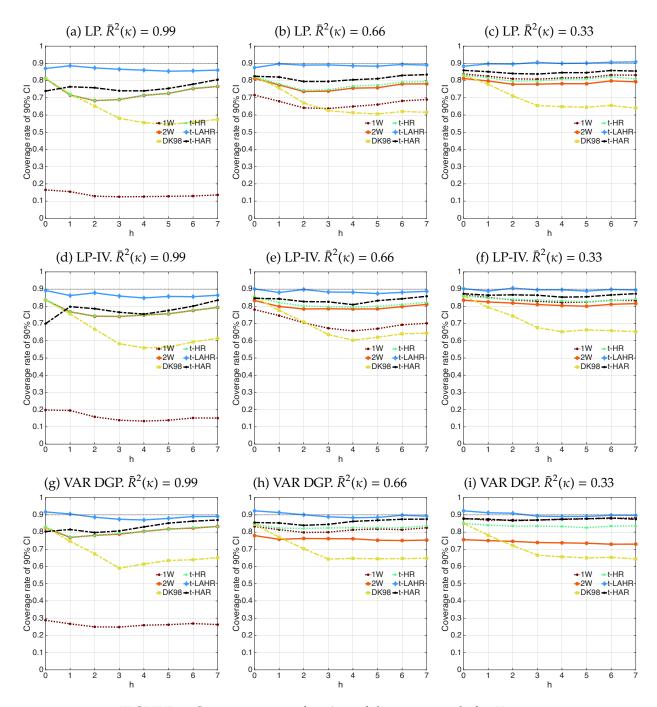


FIGURE 2. Coverage rates of 90% confidence intervals for T = 30.

Note: 1W refers to one-way (unit-level) clustering, 2W to two-way clustering, DK98 to Driscoll–Kraay, and *t*-HR/*t*-LAHR/*t*-HAR to the time-level clustering approaches discussed in the text.

5 Empirical illustration

We now discuss an empirical exercise that demonstrates the applicability of our methods in a setup featuring time-varying s_{it} and unbalanced panels, and compares our practical recommendation to popular alternatives. The exercise is motivated by the burgeoning literature on the role played by firm heterogeneity and financial frictions in the propagation of monetary policy.

Data and background. Quantifying firm-level responses to exogenous changes in policy is a key empirical goal as it is informative on the mechanisms through which monetary policy operates. For instance, Crouzet and Mehrotra (2020) focus on the role of firm size for investment response heterogeneity, finding larger (albeit noisy) responses for smaller firms; Ottonello and Winberry (2020) instead focus on default risk, finding larger responses for less risky companies.

For our empirical analysis, we construct a dataset similar to the latter based on Compustat and high-frequency identified monetary policy shocks (Gürkaynak, Sack, and Swanson, 2005; Gorodnichenko and Weber, 2016). This results in an unbalanced panel for the period 1990Q1–2010Q4 with observations on firm-level investment, size, and leverage.³² In total, there are T = 80 quarters and N = 4,187 individual companies which, net of missing data, amount to 235,233 observations.

We consider regressions of cumulative investment changes $Y_{i,t+h} = \log(k_{i,t+h}/k_{i,t-1})$ (k_{it} being the capital stock) on policy shocks X_t interacted with s_{it} , a vector containing size, leverage, and their product. From Section 3, we know that under unrestricted heterogeneity the population counterpart is the linear projection of firm-level impulse responses on s_{it} . Thus, including size and leverage together (as well as their interaction) in s_{it} is a way to enrich the linear approximation.

Synthetic time series representation. A fundamental insight of our paper is that the synthetic time series form of the microdata is a sufficient statistic for the panel LP; a low

³²We use the paper's replication code to build the data and verify that we can replicate the original results, with minor numerical differences that can be attributed to input data revisions. Firm size is measured by the value of total assets held by a company while leverage is its debt-to-assets ratio. We have also tried the distance-to-default measure in Ottonello and Winberry (2020) with qualitatively similar results.

dimensional representation of a highly complex, unbalanced dataset.³³

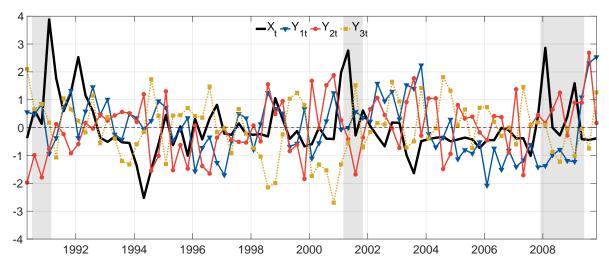


FIGURE 3. Synthetic time series representations.

Note: Grey areas are NBER-dated recessions. s_1 is size, s_2 is leverage and s_3 is the interaction. X_t and \hat{Y}_t are standardized to zero mean and unit variance; $X_t > 0$ indicates a surprise cut in the Fed Funds rate.

Figure 3 displays it for the three components of s_{it} . It is clear that movements in synthetic outcomes concurrent with surprise cuts in policy rates, mostly around recessions, are the main source of identification. There is also substantial variation in synthetic outcomes unrelated to X_t , indicating the presence of omitted aggregate or non-negligible idiosyncratic shocks — the central premises of our paper.

Estimation and inference method comparison. Figure 4 reports point estimates and 90% confidence intervals for the coefficient on each entry of $s_{it}X_t$ at different horizons.³⁴ According to the *t*-LAHR intervals, the evidence favors the hypothesis that larger and less

$$\hat{\beta}(h) = \frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} d_{it} s_{it} X_t Y_{i,t+h}}{\sum_{t=1}^{T-h} \sum_{i=1}^{N} d_{it} s_{it}^2 X_t^2} = \frac{\sum_{t=1}^{T-h} \omega_t X_t \hat{Y}_{t+h}}{\sum_{t=1}^{T-h} \omega_t X_t^2},$$

where $\omega_t = \sum_{i=1}^N d_{it} s_{it}^2$ and $\hat{Y}_{t+h} = (\sum_{i=1}^N s_{it}^2)^{-1} \sum_{i=1}^N s_{it} Y_{i,t+h}$. This is a weighted least squares regression of slope coefficients \hat{Y}_{t+h} on X_t . Note that if $s_{it} = 1$ the weights boil down to the number of non-missing observations $\omega_t = \sum_{i=1}^N d_{it}$, as intuition suggests. Our theory applies with data missing at random.

³³Remark 4 generalizes as follows. Let $d_{it} = 0$ indicate a missing observation with $d_{it} = 1$ otherwise. Abstracting from controls, the panel local projection estimator with a time-varying s_{it} is

 $[\]omega_t = \sum_{i=1}^N d_{it}$, as intuition suggests. Our theory applies with data missing at random.

³⁴The panel local projections include as controls unit and time effects, lagged firm-level sales growth, and both lagged GDP growth and lagged unemployment interacted with s_{it} . For t-LAHR inference we include p lags of Y_{it} and $s_{it}X_t$. We limit $p = \min\{h, 2\}$ to discipline the number of regressors in view of the dimension of s_{it} . One-way, two-way and Driscoll-Kraay are implemented without lag-augmentation.

indebted firms respond less to monetary policy shocks, with the size effect more persistent and not much interaction between the two.

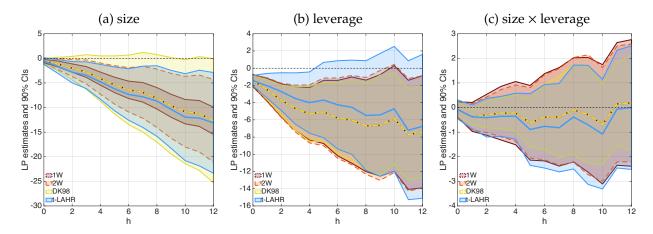


FIGURE 4. Point estimates and 90% confidence intervals.

Note: The procedures are one-way unit-level clustering (1W, dotted line), two-way clustering (2W, dashed line), Driscoll-Kraay (DK98, dash dotted line), and *t*-LAHR (solid line).

From an applied point of view, the main message is that popular methods can deviate significantly from the (asymptotically robust) t-LAHR method. For example, one-way clustering produces intervals that are too short (panel (a)) and too wide (panel (c)). Two-way clustering is close to t-LAHR but produces different conclusions in panel (b) and may be unreliable in high-persistence, low-signal setups. Finally, Driscoll-Kraay intervals can be misleading with T = 80. In fact, they lead to exactly the opposite conclusions about the role of size and leverage.

6 Conclusion

The use of microdata to answer macro questions offers an exciting avenue to study how agents respond to economy-wide policies. Possibilities include a better understanding of the dynamic transmission of shocks and the nature of heterogeneity.

Challenges are ubiquitous too. We propose a disciplined approach to uncertainty quantification when both aggregate and idiosyncratic shocks coexist and interest lies in parameters identified principally by macro shocks. The size of aggregate shocks relative

 $^{^{35}}$ This can happen even in the same exercise because the estimation errors of different coefficients load differently on the macro and micro components of the regression score. Figure 4 suggests the size coefficient is driven by the macro component and the other coefficients by the micro component.

to micro-level noise determines both the strength of identifying variation and the nature of estimation error. Since this signal-noise ratio is unknown to the researcher and context-specific, it is important to design inference procedures that are uniformly valid across all such macro signal regimes.

One such scenario is the estimation of impulse responses to macro shocks via panel local projections when rich microdata and a measurement of the shock of interest are available. This includes cases where the shock of interest is directly observed, retrieved from a macroeconomic system, or contaminated with measurement error. Despite the complex environment, inference is simple and robust: it involves including lags as controls and clustering at the time level, and is valid regardless of the relative signal of macro shocks in the microdata.

Our basic framework generalizes beyond the empirical applications we have focused on. Other, related literatures where identification comes from randomness in group level shocks include regional-exposure and shift-share designs. In fact, impulse responses are sometimes an object of interest too — see, for instance, the literature on cross-sectional fiscal multipliers (Chodorow-Reich, 2019).

We also leave some interesting dimensions for future research. Quantifying signal-to-noise (perhaps a lower bound) seems relevant in settings where uniform inference is not possible; we expect that these issues become more salient as macroeconomists embrace the use of microdata to sharpen identification (Nakamura and Steinsson, 2018). On a different note, strong persistence of micro-level shocks is likely a feature of many datasets, and this is only captured in an indirect sense by our signal-to-noise device. Formalizing the idea of (possibly heterogeneous) non-stationarities along these lines seems promising and full of empirical content. Finally, extensions to simultaneous inference over impulse response horizons could be made building on the techniques in Jordà (2009) and Montiel Olea and Plagborg-Møller (2019).

A Proofs

Proposition 1

Let $\tilde{\beta}(h) = \left(\sum_{i=1}^{N} \hat{s}_{i}^{2}\right)^{-1} \sum_{i=1}^{N} \hat{s}_{i} \beta_{ih}$ be the coefficient in the (infeasible) regression of β_{ih} on \hat{s}_{i} —the finite-population counterpart to $\beta(h)$. Also, define

$$\xi_{it}(h,\kappa) = \sum_{\ell=0}^{\infty} \left(\iota_{\ell}(h) \beta_{i\ell} X_{t+h-\ell} + \gamma_{i\ell} Z_{t+h-\ell} + \kappa \delta_{i\ell} u_{i,t+h-\ell} \right), \tag{19}$$

$$\xi_{t}(h,\kappa) = \frac{1}{N} \sum_{i=1}^{N} \hat{s}_{i} \xi_{it}(h,\kappa) = \sum_{\ell=0}^{\infty} \left(\iota_{\ell}(h) \bar{\beta}_{\ell} X_{t+h-\ell} + \bar{\gamma}_{\ell} Z_{t+h-\ell} + \frac{\kappa}{N} \sum_{i=1}^{N} \hat{s}_{i} \delta_{i\ell} u_{i,t+h-\ell} \right)$$

where $\iota_{\ell}(h) = 1 - 1\{h \leq \ell \leq h + p\}, \ \bar{\beta}_{\ell} = N^{-1} \sum_{i=1}^{N} \hat{s}_{i} \beta_{i\ell} \text{ and } \bar{\gamma}_{\ell} = N^{-1} \sum_{i=1}^{N} \hat{s}_{i} \gamma_{i\ell}.$ Finally, let $V(h, \kappa) = \operatorname{Var}_{\kappa} \left(X_{t} \xi_{t}(h, \kappa) \mid \{\theta_{i}, s_{i}\}_{i=1}^{N} \right).$

Proof of Propositions 1. Let $\sum_{i,t}$ denote summation over $1 \le t \le T - h$ and $1 \le i \le N$. For any $\psi \in \mathbb{R}^d$,

$$\begin{split} \left(\sum_{i,t} \hat{x}_{it}(h)^2\right) & \left(\hat{\beta}(h) - \tilde{\beta}(h)\right) = \sum_{i,t} \hat{x}_{it}(h) \left(Y_{i,t+h} - \tilde{\beta}(h)\hat{s}_i X_t - \psi' W_{it}\right) \\ & = \sum_{i,t} \hat{s}_i X_t \left(Y_{i,t+h} - \beta_{ih} X_t - \psi' W_{it}\right) \\ & - \sum_{i,t} (\hat{\pi}(h)' W_{it}) \left(Y_{i,t+h} - \tilde{\beta}(h)\hat{s}_i X_t - \psi' W_{it}\right). \end{split}$$

The first line uses $\sum_{i,t} \hat{x}_{it}(h)^2 = \sum_{i,t} \hat{x}_{it}(h)\hat{s}_i X_t$ and $\sum_{i,t} \hat{x}_{it}(h) W_{it} = 0_{d \times 1}$ (to introduce ψ). The second line uses $\hat{x}_{it}(h) = \hat{s}_i X_t - \hat{\pi}(h)' W_{it}$ and $\sum_{i,t} \hat{s}_i X_t (\tilde{\beta}(h) \hat{s}_i X_t - \beta_{ih} X_t) = 0$.

We can choose ψ so that

$$\sum_{i,t} \hat{s}_i X_t \left(Y_{i,t+h} - \beta_{ih} X_t - \psi' W_{it} \right) = \sum_{i,t} \hat{s}_i X_t \xi_{it}(h,\kappa) = N \sum_{t=1}^{T-h} X_t \xi_t(h,\kappa).$$

Here, W_{it} consists of p lags of $\hat{s}_i X_t$, unit indicators, and (possibly) time indicators (so that d = p + N + T). To choose ψ , we set the coefficient on $\hat{s}_i X_{t-\ell}$ to $\tilde{\beta}(h+\ell) = \left(\sum_{i=1}^N \hat{s}_i^2\right)^{-1} \sum_{i=1}^N \hat{s}_i \beta_{i,h+\ell}$, the coefficient on the unit-i indicator to μ_i , and the coefficients on time indicators to zero. Moreover, $\hat{\pi}(h)'W_{it} = \hat{s}_i(X_t - \hat{X}_t(h))$ with $\hat{X}_t(h)$ the residual from a regression of X_t on

 X_{t-1}, \ldots, X_{t-p} and an intercept. Then,

$$\sum_{i,t} (\hat{\pi}(h)'W_{it}) \left(Y_{i,t+h} - \tilde{\beta}(h)\hat{s}_i X_t - \psi'W_{it} \right) = \sum_{i,t} (\hat{\pi}(h)'W_{it}) \xi_{it}(h,\kappa).$$

It follows that the standardized estimation error can be written as

$$\begin{split} \frac{\hat{\beta}(h) - \tilde{\beta}(h)}{\hat{\sigma}(h)} &= \frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} \hat{x}_{it}(h) (Y_{i,t+h} - \tilde{\beta}(h) \hat{x}_{it}(h))}{N \sqrt{(T-h)\hat{V}(h)}} \\ &= \sqrt{\frac{V(h,\kappa)}{\hat{V}(h)}} \times \left(\frac{\sum_{t=1}^{T-h} X_{t} \xi_{t}(h,\kappa)}{\sqrt{(T-h)V(h,\kappa)}} + R_{T}(h,\kappa)\right) \end{split}$$

where the remainder term is

$$R_T(h,\kappa) = -\frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} (\hat{\pi}(h)' W_{it}) \, \xi_{it}(h,\kappa)}{N \, \sqrt{(T-h)V(h,\kappa)}}.$$

To establish our uniform approximation we exploit drifting parameter sequences (see Andrews et al. (2020) for formal results connecting the two). For simplicity we index everything to T, including $N = N_T$. We show that for any $\{\kappa_T\}$, as $T \to \infty$,

(A)
$$\{(T-h)V(h,\kappa_T)\}^{-1/2} \sum_{t=1}^{T-h} X_t \xi_t(h,\kappa_T) \xrightarrow{d} N(0,1),$$

(B)
$$\hat{V}(h)/V(h, \kappa_T) \xrightarrow{p} 1$$
,

(C)
$$R_T(h, \kappa_T) \xrightarrow{p \ P_{\kappa_T}} 0.$$

Hence, for any such $\{\kappa_T\}$,

$$\frac{\hat{\beta}(h) - \tilde{\beta}(h)}{\hat{\sigma}(h)} \xrightarrow{P_{\kappa_T}} N(0,1).$$

We establish (A), (B) and (C) in Lemmas ??, ?? and ?? in Supplemental Appendix ??. Now, Assumptions 1(ii) and 3(iv) imply $\tilde{\beta}(h) - \beta(h) = O_{P_{\kappa_T}}(N^{-1/2})$ whereas Lemma ?? implies $\min\{1, \kappa_T^{-1}\}\hat{\sigma}(h) = O_{P_{\kappa_T}}((T-h)^{-1/2})$. Since $T/N \to 0$,

$$\frac{(\hat{\beta}(h) - \beta(h))}{\hat{\sigma}(h)} = \frac{(\hat{\beta}(h) - \tilde{\beta}(h))}{\hat{\sigma}(h)} + o_{P_{\kappa_T}}(1)$$

and the result follows.

Proposition 2

Define

$$\xi_{it}(h,\kappa) = \sum_{\ell=0}^{h} \left(\iota_{\ell}(h) \beta_{i\ell} X_{t+h-\ell} + \gamma_{i\ell} Z_{t+h-\ell} + \kappa \delta_{i\ell} u_{i,t+h-\ell} \right), \tag{20}$$

$$\xi_{t}(h,\kappa) = \frac{1}{N} \sum_{i=1}^{N} \hat{s}_{i} \xi_{it}(h,\kappa) = \sum_{\ell=0}^{h} \left(\iota_{\ell}(h) \bar{\beta}_{\ell} X_{t+h-\ell} + \bar{\gamma}_{\ell} Z_{t+h-\ell} + \frac{\kappa}{N} \sum_{i=1}^{N} \hat{s}_{i} \delta_{i\ell} u_{i,t+h-\ell} \right),$$

and, as before, let $V(h, \kappa) = \operatorname{Var}_{\kappa} \left(X_t \xi_t(h, \kappa) \mid \{\theta_i, s_i\}_{i=1}^N \right)$. By recursive substitution,

$$Y_{i,t+h} = m_i(h) + \sum_{\ell=1}^{p} (A_{\ell}(h)Y_{i,t-\ell} + B_{i\ell}(h)X_{t-\ell}) + \beta_{ih}X_t + \xi_{it}(h,\kappa),$$

for some $m_i(h)$, $\{A_{\ell}(h)\}$, $\{B_{i\ell}(h)\}$ that depend on the VAR parameters m_i , $\{A_{\ell}\}$, $\{B_{i\ell}\}$.

Proof of Proposition 2. We follow exactly the same steps as for Proposition 1. The control vector W_{it} includes p lags of Y_{it} and $\hat{s}_i X_t$ in addition to unit and time effects. In the step where we choose ψ , we set the coefficient on $Y_{i,t-\ell}$ to $A_{\ell}(h)$, the coefficient on $\hat{s}_i X_{t-\ell}$ to $\tilde{B}_{\ell}(h) = \left(\sum_{i=1}^N \hat{s}_i^2\right)^{-1} \sum_{i=1}^N \hat{s}_i B_{i\ell}(h)$, the coefficient on the unit-i indicator to $m_i(h)$, and the coefficients on time indicators to zero.

The standardized estimation error can then be written as

$$\frac{\hat{\beta}(h) - \tilde{\beta}(h)}{\hat{\sigma}(h)} = \sqrt{\frac{V(h,\kappa)}{\hat{V}(h)}} \times \left(\frac{\sum_{t=1}^{T-h} X_t \xi_t(h,\kappa)}{\sqrt{(T-h)V(h,\kappa)}} + R_T(h,\kappa)\right)$$

where the remainder term is now

$$R_{T}(h,\kappa) = -\frac{\sum_{t=1}^{T-h} \sum_{i=1}^{N} (\hat{\pi}(h)'W_{it}) \left[(\beta_{ih} - \tilde{\beta}(h)\hat{s}_{i})X_{t} + \sum_{\ell=1}^{p} (B_{i\ell}(h) - \tilde{B}_{\ell}(h)\hat{s}_{i})X_{t-\ell} + \xi_{it}(h,\kappa) \right]}{N\sqrt{(T-h)V(h,\kappa)}}$$

Let ϕ < 1. In contrast to Proposition 1, instead of a single drifting parameter we now have two. We show that for any $\{h_T, \kappa_T\}$ such that $h_T \leq \phi T$,

(A)
$$\{(T - h_T)V(h_T, \kappa_T)\}^{-1/2} \sum_{t=1}^{T - h_T} X_t \xi_t(h_T, \kappa_T) \xrightarrow{P_{\kappa_T}} N(0, 1),$$

(B)
$$\hat{V}(h_T)/V(h_T, \kappa_T) \xrightarrow{p} 1$$
,

(C)
$$R_T(h_T, \kappa_T) \xrightarrow{p} 0.$$

We prove (A), (B) and (C) in Lemmas \ref{Lemmas} and \ref{Lemmas} in Supplemental Appendix \ref{Lemmas} . The rest of the argument is identical to that of Proposition 1.

Proposition 3

Using (17), substitute $\tilde{X}_t, \tilde{X}_{t-1}, \dots, \tilde{X}_{t-p}$ in succession into (6) to obtain

$$\begin{split} Y_{i,t+h} &= \mu_i + \beta_{ih} \tilde{X}_t + \sum_{\ell=1}^p \tilde{\eta}_{i\ell} \tilde{X}_{t-\ell} + \xi_{it}(h,\kappa), \\ \xi_{it}(h,\kappa) &= \sum_{\ell=0}^\infty \left(\iota_\ell(h) \tilde{\beta}_{i\ell} X_{t+h-\ell} + \tilde{\gamma}_{i\ell} Z_{t+h-\ell} + \kappa \delta_{i\ell} u_{i,t+h-\ell} \right), \end{split}$$

for some coefficients $\{\tilde{\eta}_{i\ell}\}$, $\{\tilde{\beta}_{i\ell}\}$, $\{\tilde{\gamma}_{i\ell}\}$ that depend on $\{\beta_{i\ell}\}$, $\{\gamma_{i\ell}\}$, $\{b_{\ell}\}$, $\{c_{\ell}\}$ and satisfy the bound conditions in Assumption 3 for a suitable choice of C_{ℓ} and \underline{C} . Also define $\tilde{\boldsymbol{\beta}}(h) = \left(\sum_{i=1}^{N}\hat{s}_{i}^{2}\right)^{-1}\sum_{i=1}^{N}\hat{s}_{i}\boldsymbol{\beta}_{ih}$ with $\boldsymbol{\beta}_{ih} = (\beta_{ih}, \tilde{\eta}_{i1}, \ldots, \tilde{\eta}_{ip})'$, $\xi_{t}(h, \kappa) = N^{-1}\sum_{i=1}^{N}\hat{s}_{i}\xi_{it}(h, \kappa)$ and $\boldsymbol{V}(h, \kappa) = \operatorname{Var}_{\kappa}(\boldsymbol{X}_{t}^{*}\xi_{t}(h, \kappa) \mid \{\theta_{i}, s_{i}\}_{i=1}^{N}\}$.

Proof of Proposition 3. Following similar steps to the derivation in Proposition 1, let $\sum_{i,t}$ denote summation over $1 \le t \le T - h$ and $1 \le i \le N$. For any ψ ,

$$\begin{split} \left(\sum_{i,t} \hat{\boldsymbol{x}}_{it}(h) \hat{\boldsymbol{s}}_{i} \tilde{\boldsymbol{X}}_{t}'\right) \left(\hat{\boldsymbol{\beta}}^{\text{IV}}(h) - \tilde{\boldsymbol{\beta}}(h)\right) &= \sum_{i,t} \hat{\boldsymbol{s}}_{i} \boldsymbol{X}_{t}^{*} \left(\boldsymbol{Y}_{i,t+h} - \tilde{\boldsymbol{X}}_{t}' \tilde{\boldsymbol{\beta}}_{ih} - \boldsymbol{\psi}' \boldsymbol{W}_{it}\right) \\ &- \left(\frac{\sum_{t=1}^{T-h} \boldsymbol{X}_{t}^{*}}{(T-h)}\right) \sum_{i,t} \hat{\boldsymbol{s}}_{i} \left(\boldsymbol{Y}_{i,t+h} - \hat{\boldsymbol{s}}_{i} \tilde{\boldsymbol{X}}_{t}' \tilde{\boldsymbol{\beta}}(h) - \boldsymbol{\psi}' \boldsymbol{W}_{it}\right). \end{split}$$

Note W_{it} includes unit and (possibly) time effects. To choose ψ , set the coefficient on the unit-i indicator to μ_i and the coefficients on time indicators to zero, so that

$$\sum_{i,t} \hat{s}_i X_t^* \left(Y_{i,t+h} - \tilde{X}_t' \tilde{\boldsymbol{\beta}}_{ih} - \psi' W_{it} \right) = N \sum_{t=1}^{T-h} X_t^* \xi_t(h,\kappa),$$

$$\left(\frac{\sum_{t=1}^{T-h} X_t^*}{(T-h)} \right) \sum_{i,t} \hat{s}_i \left(Y_{i,t+h} - \hat{s}_i \tilde{X}_t' \tilde{\boldsymbol{\beta}}(h) - \psi' W_{it} \right) = \left(\frac{\sum_{t=1}^{T-h} X_t^*}{(T-h)} \right) \sum_{t=1}^{T-h} \xi_t(h,\kappa).$$

Thus, the standardized estimation error can be written as

$$\frac{\hat{\beta}_{0}^{\mathrm{IV}}(h) - \tilde{\beta}(h)}{\hat{\sigma}_{0}^{\mathrm{IV}}(h)} = \sqrt{\frac{\left(e_{1}^{\prime}\boldsymbol{J}^{-1}\right)\boldsymbol{V}(h,\kappa)\left(e_{1}^{\prime}\boldsymbol{J}^{-1}\right)^{\prime}}{\left(e_{1}^{\prime}\hat{\boldsymbol{J}}^{\mathrm{IV}}(h)^{-1}\right)\hat{\boldsymbol{V}}^{\mathrm{IV}}(h)\left(e_{1}^{\prime}\hat{\boldsymbol{J}}^{\mathrm{IV}}(h)^{-1}\right)^{\prime}}}} \times \left(\frac{\left(e_{1}^{\prime}\hat{\boldsymbol{J}}^{\mathrm{IV}}(h)^{-1}\right)\hat{\boldsymbol{V}}^{\mathrm{IV}}(h)\left(e_{1}^{\prime}\hat{\boldsymbol{J}}^{\mathrm{IV}}(h)^{-1}\right)^{\prime}}}{\sqrt{\left(T - h\right)\left(e_{1}^{\prime}\boldsymbol{J}^{-1}\right)\boldsymbol{V}(h,\kappa)\left(e_{1}^{\prime}\boldsymbol{J}^{-1}\right)^{\prime}}} + R_{T}(h,\kappa)}\right)$$

where $J = (N^{-1} \sum_{i=1}^{N} \hat{s}_i^2) E[X_t^* \tilde{X}_t']$ and the remainder term is

$$R_T(h,\kappa) = -\frac{\left\{ (T-h)^{-1} \left(e_1' \hat{\boldsymbol{J}}^{\text{IV}}(h)^{-1} \right) \sum_{t=1}^{T-h} X_t^* \right\} \sum_{t=1}^{T-h} \xi_t(h,\kappa)}{\sqrt{(T-h) \left(e_1' \boldsymbol{J}^{-1} \right) V(h,\kappa) \left(e_1' \boldsymbol{J}^{-1} \right)'}}.$$

As in Proposition 1, we show that for any $\{\kappa_T\}$ and $\lambda \neq 0_{(p+1)\times 1}$

(A)
$$\{(T-h)\lambda'V(h,\kappa_T)\lambda\}^{-1/2}\sum_{t=1}^{T-h}\lambda'X_t^*\xi_t(h,\kappa_T) \xrightarrow{\mathrm{d}} N(0,1),$$

(B)
$$\left(\lambda' \hat{V}^{\text{IV}}(h)\lambda\right)/\left(\lambda' V(h,\kappa_T)\lambda\right) \xrightarrow{p} 1 \text{ and } \hat{J}^{\text{IV}}(h) \xrightarrow{p} J$$
,

(C)
$$R_T(h, \kappa_T) \xrightarrow{p} 0$$
.

The technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ and $\ref{lem:state}$ and $\ref{lem:state}$ are technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (A), (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (B), and (C) are stated in Lemmas $\ref{lem:state}$, $\ref{lem:state}$ and $\ref{lem:state}$ are the technical steps for (B), and (C) are stated in Lemmas $\ref{lem:state}$ and $\ref{lem:stated}$ are the technical steps for (B), and (C) are stated in Lemmas $\ref{lem:stated}$.

References

Adão, R., M. Kolesár, and E. Morales (2019): "Shift-Share Designs: Theory and Inference." *Quarterly Journal of Economics*, 134, 1949–2010.

Anderson, G. and A. Cesa-Bianchi (2024): "Crossing the Credit Channel: Credit Spreads and Firm Heterogeneity." *American Economic Journal: Macroeconomics*, 16, 417–446.

Andrews, D., X. Cheng, and P. Guggenberger (2020): "Generic results for establishing the asymptotic size of confidence sets and tests," *Journal of Econometrics*, 218, 496–531.

- Andrews, D. K. (2005): "Cross-section regression with common shocks." *Econometrica*, 73, 1551–1585.
- Angrist, J. D. and G. W. Imbens (1995): "Two-Stage Least Squares Estimation of Average Causal Effects in Models with Variable Treatment Intensity." *Journal of the American Statistical Association*, 90, 431–442.
- Angrist, J. D., G. W. Imbens, and K. Graddy (2000): "The Interpretation of Instrumental Variables Estimators in Simultaneous Equations Models with an Application to the Demand for Fish." *Review of Economic Studies*, 67, 499–527.
- Arellano, M. (2003): Panel Data Econometrics, Oxford University Press.
- Arkhangelsky, D. and V. Korovkin (2023): "On Policy Evaluation with Aggregate Time-Series Shocks." Working Paper arXiv:1905.13660.
- Chang, M., X. Chen, and F. Schorfheide (2024): "Heterogeneity and Aggregate Fluctuations." *Journal of Political Economy*, 132.
- Chodorow-Reich, G. (2019): "Geographic Cross-Sectional Fiscal Multipliers: What Have We Learned?" *American Economic Journal: Economic Policy*, 11, 1–34.
- CROUZET, N. AND N. R. MEHROTRA (2020): "Small and Large Firms over the Business Cycle." American Economic Review, 110, 3549–3601.
- Drechsel, T. (2023): "Earnings-Based Borrowing Constraints and Macroeconomic Fluctuations." American Economic Journal: Macroeconomics, 15, 1–34.
- Driscoll, J. and A. Kraay (1998): "Consistent Covariance Matrix Estimation with Spatially Dependent Panel Data." *Review of Economics and Statistics*, 80, 549–560.
- Fukui, M., E. Nakamura, and J. Steinsson (2023): "The Macroeconomic Consequences of Exchange Rate Depreciations." NBER Working Paper w31279.
- Gabaix, X. (2011): "The Granular Origins of Aggregate Fluctuations." Econometrica, 79, 733–772.
- Gonçalves, S. (2011): "The moving blocks bootstrap for panel linear regression models with individual fixed effects." *Econometric Theory*, 27, 1048–1082.
- Gonçalves, S., A. M. Herrera, L. Kilian, and E. Pesavento (forthcoming): "State-dependent local projections." *Journal of Econometrics*.

- GORODNICHENKO, Y. AND M. WEBER (2016): "Are Sticky Prices Costly? Evidence from the Stock Market." *American Economic Journal: Macroeconomics*, 8, 160–199.
- GÜRKAYNAK, R. S., B. SACK, AND E. T. SWANSON (2005): "Do Actions Speak Louder Than Words? The Response of Asset Prices to Monetary Policy Actions and Statements." *International Journal of Central Banking*, 1, 55–93.
- Hahn, J., G. Kuersteiner, and M. Mazzocco (2020): "Estimation with Aggregate Shocks." *Review of Economic Studies*, 87, 1365–1398.
- Hahn, J., G. Kuersteiner, A. Santos, and W. Willigrod (2024): "Overidentification in Shift-Share Designs," Working Paper arXiv:2404.17049.
- Hansen, L. P. and R. Hodrick (1980): "Forward Exchange Rates as Optimal Predictors of Future Spot Rates: An Econometric Analysis." *Journal of Political Economy*, 88, 829–853.
- HERBST, E. P. AND B. K. JOHANSENN (2023): "Bias in Local Projections." Working paper.
- Imbens, G. W. and M. Kolesár (2016): "Robust Standard Errors in Small Samples: Some Practical Advice." *Review of Economics and Statistics*, 98, 701–712.
- JEENAS, P. AND R. LAGOS (2024): "Q-Monetary Transmission." Journal of Political Economy, 132.
- JORDÀ, O. (2009): "Simultaneous Confidence Regions for Impulse Responses." *The Review of Economics and Statistics*, 91, 629–647.
- JORDÀ, O. (2005): "Estimation and inference of impulse responses by local projections." *American Economic Review*, 95, 161–182.
- Jovanovic, B. (1987): "Micro Shocks and Aggregate Risk." Quarterly Journal of Economics, 102, 395–409.
- Känzig, D. (2021): "The Macroeconomic Effects of Oil Supply News: Evidence from OPEC Announcements." *American Economic Review*, 111, 1092–1125.
- LAZARUS, E., D. J. LEWIS, J. H. STOCK, AND M. W. WATSON (2018): "HAR Inference: Recommendations for Practice." *Journal of Business and Economic Statistics*, 36, 541–559.
- LEEPER, E. M., C. A. SIMS, AND T. ZHA (1996): "What Does Monetary Policy Do?." Brookings Papers on Economic Activity, 27, 1–78.
- Lusompa, A. (2023): "Local Projections, Autocorrelation, and Efficiency." *Quantitative Economics*, 14, 1199–1220.

- Majerovitz, J. and K. A. Sastry (2023): "How Much Should We Trust Regional-Exposure Designs?." Working paper.
- Montiel Olea, J. and M. Plagborg-Møller (2021): "Local Projection Inference is Simpler and More Robust Than You Think." *Econometrica*, 89, 1789–1823.
- Montiel Olea, J. L. and M. Plagborg-Møller (2019): "Simultaneous Confidence Bands: Theory, Implementation, and an Application to SVARs." *Journal of Applied Econometrics*, 34, 1–17.
- Montiel Olea, J. L., M. Plagborg-Møller, E. Qian, and C. K. Wolf (2024): "Double Robustness of Local Projections and Some Unpleasant VARithmetic." Working Paper NBER Working Paper 32495.
- Müller, U. K. (2004): "A theory of robust long-run variance estimation." Working paper.
- NAKAMURA, E. AND J. STEINSSON (2014): "Fiscal Stimulus in a Monetary Union: Evidence from US Regions." *American Economic Review*, 104, 753–792.
- ——— (2018): "Identification in Macroeconomics." *Journal of Economic Perspectives*, 32, 59–86.
- Nunn, N. and N. Qian (2014): "US Food Aid and Civil Conflict." *American Economic Review*, 104, 1630–1666.
- Ottonello, P. and T. Winberry (2020): "Financial Heterogeneity and the Investment Channel of Monetary Policy." *Econometrica*, 88, 2473–2502.
- Pakel, C. (2019): "Bias reduction in nonlinear and dynamic panels in the presence of cross-section dependence." *Journal of Econometrics*, 213, 459–492.
- PESARAN, M. H. (2006): "Estimation and inference in large heterogeneous panels with a multifactor structure." *Econometrica*, 74, 967–1012.
- Plagborg-Møller, M. and C. K. Wolf (2021): "Local projections and VARs estimate the same impulse responses." *Econometrica*, 89, 955–980.
- Rambachan, A. and N. Shephard (2021): "When do common time series estimands have nonparametric causal meaning?" Working paper.
- RAMEY, V. (2016): "Macroeconomic shocks and their propagation." in *Handbook of Macroeconomics*, ed. by J. B. Taylor and H. Uhlig, Elsevier, vol. 2, chap. 2.
- STAIGER, D. AND J. H. STOCK (1997): "Instrumental variables regression with weak instruments." *Econometrica*, 65, 557–586.

STOCK, J. H. AND M. W. WATSON (2016): "Dynamic factor models, factor-augmented vector autoregressions, and structural vector autoregressions in macroeconomics." in *Handbook of Macroeconomics*, ed. by J. B. Taylor and H. Uhlig, Elsevier, vol. 2, chap. 8.

——— (2018): "Identification and estimation of dynamic causal effects in macroeconomics using external instruments." *Economic Journal*, 128, 917–948.

Xu, K.-L. (2023): "Local Projection Based Inference under General Conditions." Working paper.