Three Common Factors

Elena Andreou and Patrick Gagliardini and Eric Ghysels and Mirco Rubin*

Abstract

We present a set of novel theoretical results providing a solution that addresses the perils of beta dynamics misspecification in the estimation of conditional asset pricing models. We show empirically relevant conditions under which the true factors are those common between conditional asset pricing models for individual stocks and models estimated from sorted portfolios. We find that there are at least three common factors. These are not entirely spanned by the three or five Fama-French factors, and are also related to some proxies of idiosyncratic risk and liquidity. The three factors feature superior out-of-sample pricing performance compared to standard asset pricing models.

^{*}Elena Andreou, elena.andreou@ucy.ac.cy, University of Cyprus. Patrick Gagliardini, patrick.gagliardini@usi.ch, Università della Svizzera italiana, Lugano and Swiss Finance Institute. Eric Ghysels, eghysels@unc.edu, University of North Carolina - Chapel Hill, Kenan-Flagler Business School and CEPR. Mirco Rubin, mirco.rubin@edhec.edu, ED-HEC Business School, Nice. We would like to thank the Editor and two referees for their comments which helped us improve our paper. We would also like to thank Laurent Calvet, Mike Chernov, Jianqing Fan, Andrei Gonçalves, Cam Harvey, Abraham Lioui, Raymond Kan (discussant), Daniele Massacci, Filippos Papakonstantinou, Olivier Scaillet, Raman Uppal and Dacheng Xiu (discussant) for insightful comments, as well as seminar participants at EDHEC Business School, Nice, Kenan-Flagler Finance and UNC Economics seminars, LUISS Business School, King's College London, University of Gothenburg, and participants at the 2021 Annual SoFiE Conference, the 2021 Annual Conference of the IAAE, the 2022 WFA Meeting, the 2022 American and European meetings of the Econometric Society, the 2022 (EC)² conference at ESSEC Business School, the 2022 Workshop on "Advances in alternative data and machine learning for macroeconomics and finance" at the Institute Luis Bachelier in Paris, and the 2023 AFA Annual Meeting. The first author would like to acknowledge that this work was funded by the Republic of Cyprus through the Research and Innovation Foundation (Project: INTERNATIONAL/USA/0118/0043), and the European Commission Recovery and Resilience Plan in Cyprus for the Project "Economic Modelling for Economic Policy".

I. Introduction

Currently, the dominant research theme in empirical asset pricing is the low dimensional factor representation of a large set of asset returns. Ideally any high dimensional set of asset returns should contain the information necessary to recover the factors. In practice, the literature has taken two different approaches. Jensen, Black, and Scholes (1972) and Fama and MacBeth (1973), among many others, have advocated to collect stocks into portfolios and subsequently run cross-sectional regressions using portfolios as test assets. An alternative approach is to estimate cross-sectional risk premia using the entire universe of stocks as advocated by Litzenberger and Ramaswamy (1979), among others.

The message of our paper is the following: if one is interested in estimating the latent factors in conditional asset pricing models using statistical techniques such us Principal Component Analysis (PCA), or recent extensions allowing for time varying betas such as Instrumented Principal Component Analysis (IPCA), and one has concerns about potential specification errors in the assumed beta dynamics (necessary to implement these statistical techniques for factor extraction), then one benefits from estimating models from both panels of individual stocks *and* sorted portfolios. Hence, contrary to the conventional wisdom in the literature, it is not about *either* individual stocks *or* sorted portfolios. Instead, estimating two potentially misspecified models from the two datasets helps in extracting a smaller set of factors which are not prone to beta dynamics misspecification, and explain well both the variability and the risk premia of both panels of returns.

Our paper starts with a set of novel theoretical results regarding conditional asset pricing models for individual stocks. Such models have two critical inputs: (a) the specification of risk

factors, and (b) the time series behavior of exposures to those risk factors, i.e. the specification of the time-varying betas. We establish a connection between beta dynamics prone to specification errors and the discovery of risk factors. One can think of many reasons why specification errors might occur, either related to the functional form specification of the betas and/or the economic variables driving those betas. In particular, misspecification errors in the dynamic exposure to risk spill over into additional "spurious" risk factors with the assumed beta dynamics. Our theory also predicts that potentially misspecified conditional beta models using individual stocks will lead to an overstatement of the number of risk factors. The situation is similar with portfolios where additional spurious factors will appear due to the interactions among the specific types of portfolio sorting and underlying and "true", but unknown (or only partially known), beta dynamics.

More importantly, we provide a solution to account for the perils of beta dynamics misspecification. We show that under a set of mild regularity conditions the true factors are those common between conditional asset pricing models for individual stocks and models estimated from sorted portfolios. Put differently, the factors that are common between the panels of (a) individual stocks and (b) sorted portfolios, identify the true factors, even though the time-varying exposures are misspecified.

The task for finding the pervasive factors that are common between two large panels is not trivial and requires theoretical insights so far not explored in the empirical asset pricing literature. To achieve the task set forth we need to expand the theory underpinning a procedure proposed by Andreou, Gagliardini, Ghysels, and Rubin (2019) (henceforth AGGR). They study a situation where (latent) factors $h_{1,\tau}$ and $h_{2,\tau}$, are estimated from two separate panels of data by classical PCA, assume constant loadings, and one is interested in testing how many factors are common

between them. However, AGGR show that the common factor space is identified by examining how many linear combinations of respectively $h_{1,\tau}$ and $h_{2,\tau}$ are perfectly correlated. Equivalently, they introduce a test for the number of canonical correlations between $h_{1,\tau}$ and $h_{2,\tau}$ equal to one and derive its asymptotic distribution. The AGGR setting does not directly translate into a procedure suitable for asset pricing applications. In particular, we need to generalize the AGGR asymptotic theory to (a) estimators for conditional asset pricing models, i.e. linear models with time-varying betas, and (b) estimators where the means of the latent factors represent risk premia and the risk premia of the test assets enter explicitly in the factor's estimation objective function. Indeed, one of the contributions of this paper is the extension of the theory of AGGR to allow factors to be estimated by means of Instrumented PCA (or IPCA) of Kelly, Pruitt, and Su (2019), which allows betas to be linear functions of stock characteristics and Risk Premium PCA (RP_PCA) advocated by Lettau and Pelger (2020b), a version of PCA with a penalty term accounting also for the cross-sectional pricing error in expected returns.

The novel testing procedure identifies at least 3 factors at the intersection of larger factor spaces extracted by potentially misspecified IPCA and RP_PCA from individual US stock returns and a large set of sorted portfolios of the same stocks during the historical period 1966M1-2020M12. Surprisingly, we find that neither the Fama and French 3 (FF3) nor 5 (FF5) factor models, both with or without a momentum factor (hence up to 6 factors), span those 3 factors, i.e. span the factor space common between individual stocks and sorted portfolios. In fact out of the 6 factors considered, only the excess market returns factor seems to be the most related to the common factors, while all the other 5 factors are only partially spanned by the common

¹To sort out genuine risk factors Pukthuanthong, Roll, and Subrahmanyam (2019) also rely on canonical correlations, but do not present a formal statistical procedure.

factors, and a large part of their variability is specific to portfolio sorting. For convenience we will call the 3 common factors 3CF. As our theory predicts, Kelly et al. (2019) find more factors, namely 5, applying IPCA to individual stocks, and Lettau and Pelger (2020b) applying RP_PCA to sorted portfolios find between 5 and 10 factors depending on which penalty is used.

The search for factors has been on steroids with literally hundreds of potential additional candidate factors beyond FF3 suggested in the literature. The endeavor has been dubbed the factor zoo by Cochrane (2011) and terms such as p-hacking (meaning data-snooping or data-mining) have been used to describe the hunt for factors.² The literature started of with the pretty tame single factor model, i.e. the CAPM. It is perhaps more appropriate to say that we moved from a petting zoo to a jungle. For example, Harvey and Liu (2019) have documented over 400 factors published in top journals. In our empirical application we use a data set of over a thousand portfolios associated with 205 characteristics. It takes up to 10 PCs from this factor zoo to span the space of 3CF.

Using multiple in-sample performance evaluation measures we find that 3CF perform better than a large collection of observable and latent factor models in pricing individual stock and sorted portfolios assets. Turning to the out-of-sample (OOS) analysis the results yield several interesting empirical findings. The 3CF yield the highest total, pricing and predictive OOS R^2 s with respect to the same benchmark models. For the individual stocks as well as sorted portfolios the OOS predictive R^2 s gains using 3CF can be 80% and 50% vis-à-vis for example the Fama-French factors.

²See Harvey, Liu, and Zhu (2016), McLean and Pontiff (2016), Chorida, Goyal, and Saretto (2020) Hou, Xue, and Zhang (2020), Feng, Giglio, and Xiu (2020), Chen (2019), among others.

In order to provide economic interpretation, we regress the observable factors from the zoo one by one onto 3CF and document which ones yield the best fit. We find two of the three Fama and French factors (CAPM Beta and Size) partially explain 3CF. Other factors in the zoo also explain part of the variation of 3CF. These portfolios are mostly different proxies of idiosyncratic risk and liquidity or uncertainty, such as Bid-ask Spread, Cash-flow to price variance, Volume to market equity, EPS Forecast Dispersion, Days with zero trades, Volume Variance, and Price delay R-squared. Our findings corroborate those of Dello-Preite, Uppal, Zaffaroni, and Zviadaze (2024), who use a different methodology from ours involving only sorted portfolios as test assets. They also find that a model with only three factors, one being the market and another one being related to different idiosyncratic/unsystematic risk factors, prices well a large cross-section of sorted portfolio returns similar to ours.

The rest of the paper is organized as follows. Section II contains novel theoretical results on the relationship between conditional asset pricing models for individual stocks, beta dynamics, and sorted portfolios. Section III introduces the spanning test and details the data on the various cross-sections of asset returns used to estimate common latent factor spaces. Section IV covers the empirical implementation of the testing procedure, followed by Section V where we report the results of an extensive empirical study comparing the asset pricing performance of the common factors with widely used factors in the asset pricing literature. Section VI revisits the topic of the factor zoo. Conclusions appear in Section VII. Technical details and additional empirical results appear online as *Supplementary Material*.

II. Conditional Asset Pricing Models and Common Factors

We present a number of theoretical results pertaining to the estimation of conditional asset pricing models where the dynamic specification of latent factor exposure is potentially misspecified. In a first subsection we present results for individual stocks, whereas in a second we cover sorted portfolios. The final subsection deals with group-factor models.

A. Factors and Time-Varying Betas - Individual Stocks

A researcher is interested in extracting factors from the cross-section of individual stock expected returns, but faces potential model specification issues. We start with the IPCA estimator of Kelly et al. (2019) for the conditional linear factor asset pricing model for individual stocks:

(1)
$$y_{i,\tau}^{ind} = \beta'_{i,\tau-1} f_{\tau}^c + \varepsilon_{i,\tau},$$

with latent factors f_{τ}^{c} (the reason for the superscript 'c' will become clear later) and where the betas are linear functions of so called instruments:

(2)
$$\beta_{i,\tau-1} = \Gamma \mathscr{Z}_{i,\tau-1} + \eta_{i,\tau-1}.$$

We can interpret (2) as a cross-sectional (multivariate) linear regression of the betas onto the instruments for any τ , with residual $\eta_{i,\tau-1}$. Because such a projection is always possible, the only restriction is that the matrix Γ is time-independent (see Gagliardini and Ma (2022) for further discussion). Intuitively, a constant Γ means that the time-variation in instruments $\mathcal{Z}_{i,\tau-1}$ is rich enough to guarantee that the cross-sectional relationship between betas and instruments $\mathcal{Z}_{i,\tau-1}$ is

stable across time. Substituting (2) into (1) we get:

$$y_{i,\tau} = [\mathscr{C}\mathscr{Z}_{i,\tau-1}]' f_{\tau}^c + e_{i,\tau},$$

where
$$\mathscr{C} = \Gamma$$
 and $e_{i,\tau} = \eta'_{i,\tau-1} f_{\tau} + \varepsilon_{i,\tau}$.

The perils of specification errors in conditional betas are well known, see e.g. Ghysels (1998). It is therefore reasonable to suppose that the researcher faces model specification challenges and uses instruments $Z_{i,\tau-1}$ not necessarily equal to $\mathcal{Z}_{i,\tau-1}$. One can think of many reasons for potential specification errors, including omitted instruments or the assumed linear functional form - where sieve estimators and neural net approximations have been suggested (see e.g. Chen, Roussanov, and Wang (2023) and references therein). Polynomial approximations in particular would expand the set of instruments and therefore the sources of specification errors. In addition, the specification error causes the cross-sectional relationship with betas to be time varying:

(4)
$$\beta_{i,\tau-1} = \Gamma_{\tau-1} Z_{i,\tau-1} + \eta_{i,\tau-1}^Z,$$

³Note that the error term $e_{i,\tau}$ contains the component $\eta'_{i,\tau-1}f_{\tau}$ which may be cross-sectionally correlated. Kelly, Pruitt, and Su (2020) argue that this is not an issue for the consistent estimation of the factor values and the matrix Γ , because the error component $\eta'_{i,\tau-1}f_{\tau}$ is cross-sectionally orthogonal to the instruments, namely $1/N\sum_{i=1}^{N}\mathscr{Z}_{i,\tau-1}\eta'_{i,\tau-1}f_{\tau}\to 0$ in probability as $N\to\infty$ (see their Assumption A). Such an orthogonality is of course a consequence of the error term $\eta_{i,\tau-1}$ being the residual of the cross-sectional regression of betas onto instruments. That error term contributes to the asymptotic distribution of the IPCA estimates.

yielding:

(5)
$$y_{i,\tau}^{ind} = Z'_{i,\tau-1}g_{\tau} + e_{i,\tau}^{Z},$$

where $g_{\tau}:=\Gamma_{\tau-1}'f_{\tau}^c$ and $e_{i,\tau}^Z=(\eta_{i,\tau-1}^Z)'f_{\tau}^c+\varepsilon_{i,\tau}$. To proceed we make the following assumption:

ASSUMPTION 1 It is assumed that the linear dynamics appearing in equation (4) has the following decomposition: $\Gamma_{\tau} = \Gamma^0 + \sum_{j=1}^J \Gamma^j \zeta_{\tau}^j$, with ζ_{τ}^j latent scalar stochastic processes and Γ^j for j=0,1,...,J deterministic matrices, such that the $m \times k^1$ compound matrix $C':=[\Gamma^{0'}:\Gamma^{1'}:\cdots:\Gamma^{J'}]$ is full column rank, where $k_1:=(J+1)k^c \leq m$.

The full rank condition on the matrix C is used to exclude that one can rewrite the model in terms of a smaller number of latent factors upon rotation. Then the following Proposition characterizes how specification errors impact the time-varying beta model, when a researcher uses instruments $Z_{i,\tau-1}$ which feature cross-sectional variation $\Gamma_{\tau-1}$, instead of $\mathscr{Z}_{i,\tau-1}$ with fixed Γ .

PROPOSITION 1 Let Assumption 1 hold with the data generating process appearing in equations (1)-(2). Then under the model specification appearing in (5) we obtain: $g_{\tau} = \Gamma'_{\tau-1} f_{\tau}^c = \Gamma'_{\tau-1}$

(6)
$$y_{i,\tau}^{ind} = (CZ_{i,\tau-1})' f_{1,\tau} + e_{i,\tau}^Z.$$

The proof of the proposition appears in Appendix Section A. A few remarks about model specification are in order. First, so far we have not included a constant vector in the time-varying beta specification. This is easy to do and in fact it is equivalent to including a constant (across

assets and dates) in the instrument vector $Z_{i,\tau-1}$. If we make the constant explicit (without including it in the instrument set), equation (3) becomes: $y_{i,\tau}^{ind} = [\mathcal{B} + \mathcal{C}\mathcal{Z}_{i,\tau-1}]' f_{\tau}^c + e_{i,\tau}$. Then, a result similar to equation (6) in Proposition 1 is as follows: $y_{i,\tau}^{ind} = (B + CZ_{i,\tau-1})' f_{1,\tau} + e_{i,\tau}^Z$, where the lower block of B is nil. Second, what about having the effect of macro-variables in the time-varying betas? This would be something like: $y_{i,\tau}^{ind} = [\mathcal{B}^0 + \mathcal{B}^1\mathcal{Z}_{\tau-1} + \mathcal{C}\mathcal{Z}_{i,\tau-1}]' f_{\tau}^c + e_{i,\tau}$. Because the part $(\mathcal{B}^0 + \mathcal{B}^1\mathcal{Z}_{\tau-1})' f_{\tau}^c$ is stock-independent we can get rid of it by cross-sectional demeaning of the returns and then estimate the rest of the model. In our empirical investigation we do not consider macroeconomic instruments, and therefore do not apply such a demeaning to individual stock returns. Finally, we do not include stock-indexed parameter matrices \mathcal{B}_i^0 , \mathcal{B}_i^1 and \mathcal{C}_i , since this would be a significant departure from the IPCA framework.

B. Factors and Time-Varying Betas - Sorted Portfolios

Let us now consider portfolio formation. Assuming that the weight of individual assets i at time τ in portfolio j equals $\alpha_{j,i,\tau-1}$ and the model for the individual stock is: $y_{i,\tau}^{ind} = [\mathscr{C}\mathscr{Z}_{i,\tau-1}]'f_{\tau}^c + e_{i,\tau}, \text{ then the portfolio returns are}$

(7)
$$y_{j,\tau}^p = \frac{1}{N_1} \sum_{i=1}^{N_1} \alpha_{j,i,\tau-1} y_{i,\tau}^{ind} = [\mathscr{C}\mathscr{Z}_{j,\tau-1}^p]' f_{\tau}^c + e_{j,\tau}^p,$$

where $\mathscr{Z}^p_{j,\tau-1}:=\frac{1}{N_1}\sum_{i=1}^{N_1}\alpha_{j,i,\tau-1}\mathscr{Z}_{i,\tau-1}$ and $e^p_{j,\tau}=\frac{1}{N_1}\sum_{i=1}^{N_1}\alpha_{j,i,\tau-1}e_{i,\tau}$. Next, we suppose the following assumption about the large cross-sectional limit holds.

ASSUMPTION 2 It is assumed $\mathscr{Z}_{j,\tau-1}^p \to W_j^0 + W_j Z_{\tau-1}^*$ as $N_1 \to \infty$, for a latent vector of macro-variables Z_{τ}^* .

Furthermore, we want that the variables driving the conditional betas, and portfolios formation, do not mask the common factor f_{τ}^{c} . Hence, the following assumption.

ASSUMPTION 3 The elements of $\mathscr{C}W_j^0$ are non-zero for a non-vanishing fraction of portfolios.

It is worth noting that the above assumption puts few restrictions on how sorted portfolios are formed. The sorted portfolio panel in our empirical analysis is based on standard portfolio-formation methods used to construct the Chen and Zimmermann (2022) data set. But this is not key to the analysis in our paper. For example, Bryzgalova, Pelger, and Zhu (2024) suggest to used random forest models to construct portfolios. That would work perfectly fine for us as well as long as the portfolios do not mask any of the factors f_{τ}^c .

PROPOSITION 2 *Under Assumptions 2 and 3 the portfolios excess returns are such that:*

(8)
$$y_{j,\tau}^p = \lambda_j' f_{2,\tau} + e_{j,\tau}^p,$$

where $f_{2,\tau}=(f_{\tau}^{c\prime},f_{2,\tau}^{s\prime})'$ with the factor spaces spanned by f_{τ}^c and $f_{2,\tau}^s=Z_{\tau-1}^*\otimes f_{\tau}^c$, do not intersect.

From equation (6) we learned that the cross-section of individual stock returns features a set of factors $f_{1,\tau}$ when a researchers uses a IPCA estimator with $Z_{i,\tau-1}$, not necessarily the proper set of instruments. From equation (8) we observed that sorted portfolio returns features factors $f_{2,\tau}$ when estimated with PCA. Vectors $f_{1\tau}$ and $f_{2\tau}$ consist of a common subgroup of latent factors f_{τ}^c as well as factors $f_{1\tau}^s$ and $f_{2\tau}^s$ that are specific to the panel of respectively individual stocks and sorted portfolios. This is essentially a group-factor setup we discuss next.

C. Group-factor models and (I)PCA

A group-factor model, studied by AGGR, applies to a situation where one has two panel data sets - in our application one will be a panel of individual stocks and the other of sorted portfolios. AGGR only considered fixed loadings PCA applications in group-factor models. Here we extend their framework to a setting where one uses IPCA in (at least) one of the groups. The panel data in group ℓ is denoted by $y_{\ell,\tau}$ with:

(9)
$$\begin{bmatrix} y_{1,\tau} \\ y_{2,\tau} \end{bmatrix} = \begin{bmatrix} \Lambda_{1,\tau}^c & \Lambda_{1,\tau}^s & 0 \\ \Lambda_{2}^c & 0 & \Lambda_{2}^s \end{bmatrix} \begin{bmatrix} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{bmatrix} + \begin{bmatrix} \varepsilon_{1,\tau} \\ \varepsilon_{2,\tau} \end{bmatrix},$$

where without loss of generality, the group-specific factors $f_{1,\tau}^s$ and $f_{2,\tau}^s$ are unconditionally orthogonal to the common factor f_{τ}^c . Since the unobservable factors can be standardized, we have:

$$(10) \qquad \mathbb{E} \begin{bmatrix} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{bmatrix} = \begin{bmatrix} \mu^c \\ \mu_1^s \\ \mu_2^s \end{bmatrix}, \qquad \text{and} \qquad \Sigma_F := \mathbb{V} \begin{bmatrix} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{bmatrix} = \begin{bmatrix} I_{k^c} & 0 & 0 \\ 0 & I_{k_1^s} & \Upsilon \\ 0 & \Upsilon' & I_{k_2^s} \end{bmatrix},$$

where the expected values of the factors are finite, and matrix Σ_F is positive-definite. Note that the group-specific factors can be correlated with a non-zero covariance Υ . Indeed, in our application we expect that $f_{1,\tau}^s = \zeta_{\tau-1} \otimes f_{\tau}^c$ and $f_{2,\tau}^s = Z_{\tau-1}^* \otimes f_{\tau}^c$ (see Proof of Proposition 2) are correlated. The factor f_{τ}^c is common to the two panels (groups). For our methodology to apply we need that this factor is indeed pervasive in both panels, and is the only one having this property (up to normalization). For this to hold, we allow the group-specific factors to be correlated but

they do not share common linear transformations that are *perfectly* correlated, as stated in the next assumption.

ASSUMPTION 4 The canonical correlations between $\zeta_{\tau-1} \otimes f_{\tau}^c$ and $Z_{\tau-1}^* \otimes f_{\tau}^c$ are all strictly smaller than 1.

Then we have the following result:

PROPOSITION 3 Under Assumptions 1 through 4 the group-factor model appearing in equations (9) - (10) will yield estimates of f_{τ}^{c} unlike IPCA alone applied to individual stocks when proper characterization of instruments is in doubt.

Hence the combination of (misspecified) IPCA on the panel of single stocks, and PCA on the panel of portfolios, yields a group-factor model with the fundamental factor f_{τ}^c that is common among the two.

III. Factor Space Spanning Test

Summarizing the results in the previous section, we have shown that if individual stocks were used, a researcher would conclude that f^c_{τ} and $f^s_{1,\tau}$ are risk factors, if we cannot rule out the possibility of specification errors in the dynamics of loadings. Conversely, a researcher starting from sorted portfolio returns would conclude that the risk factors are instead f^c_{τ} and $f^s_{2,\tau}$. In the former case IPCA or some variation thereof is used, in the latter IPCA, PCA or some variation thereof. Whilst to the best of our knowledge, other methods cannot address this, our testing procedure allows us to identify and estimate f^c_{τ} , despite the fact that a researcher faces specification errors in the formulation of conditional asset pricing models. In this section we will

focus on a succinct presentation of the econometric methods, while Appendix Sections A.1 and A.2 and *Supplementary Material* Section OA.3 cover all the technical details, including the definitions of the RP_PCA and IPCA estimators. There are four steps, namely:

Step 1: We start with extracting factors from each panel separately. This means extracting factors from individual stock returns and doing the same for sorted portfolios. As is typically done (see e.g. Lehmann and Modest (1988) and recently Kim and Korajczyk (2024), among many others), we will study non-overlapping sample blocks of generic length w of the panels covering data from t-w+1 to t.⁴

Step 2: Let $k_j = k^c + k_j^s$, for j = 1, 2, be the dimensions of the factor spaces for the two panels, and define $\underline{k} = \min(k_1, k_2)$. We collect the factors of each group in the k_j -dimensional vectors $h_{j,\tau}$ then $h_{1,\tau} = \mathcal{H}_1 \left[f_{\tau}^{c\prime} , \ f_{1,\tau}^{s\prime} \right]'$ and $h_{2,\tau} = \mathcal{H}_2 \left[f_{\tau}^{c\prime} , \ f_{2,\tau}^{s\prime} \right]'$, meaning that the factors we extract from each group are some linear transformation \mathcal{H}_j of the underlying factors, with \mathcal{H}_j full rank $k_j \times k_j$ matrices, for j = 1, 2. This means that some linear combinations of $h_{1,\tau}$ - namely those corresponding to f_{τ}^c - are *perfectly* correlated with linear combinations of $h_{2,\tau}$ and vice versa. Let us recall at this point the purpose of canonical correlation analysis. In general canonical correlation applies to a setting where we have two random vectors, in our application $h_{1,\tau}$ and $h_{2,\tau}$, and one finds linear combinations of respectively $h_{1,\tau}$ and $h_{2,\tau}$ which have maximum correlation with each other. Therefore we are interested in finding how many of these linear

⁴As noted by Kelly et al. (2019) there is no need to look at block sampling for IPCA using individual stocks. We will implement IPCA on the full sample as well as on block samples. We use a block sampling scheme to avoid look ahead biases in full sample factor extraction as well as survivorship biases for individual firms: see Section OA.1 in the *Supplementary Material* for further discussion.

combinations, also known as canonical variables, are perfectly correlated, i.e. have canonical correlation equal to one.

Step 3: Proposition A.1 in the Appendix tells us that the dimension k^c is the number of unitary canonical correlations between $h_{1,\tau}$ and $h_{2,\tau}$. The largest possible number of common factors is $\underline{k} = \min(k_1, k_2)$. We develop a test for $k^c : H_0(r) : k^c = r$ against $H_1(r) : k^c < r$, for any given $r = \underline{k}, \underline{k} - 1, \ldots, 1$. More precisely, we sort the canonical correlations from high to low and let $\hat{\rho}_{\ell}$ be the ℓ -th sorted sample canonical correlation between the factors $\hat{h}_{1,\tau}$ and $\hat{h}_{2,\tau}$ estimated on a sample of length w, and let:

(11)
$$\hat{\xi}(r) = \sum_{\ell=1}^{r} \hat{\rho}_{\ell},$$

be the sum of the r largest sample canonical correlations. We reject the null for $r=k^c$ common factors $H_0=H(k^c)$ when $\hat{\xi}(k^c)-k^c$ is negative and large - namely the sum of the largest k^c estimated canonical correlations is substantially less than k^c .⁵ The test statistic is a re-centered and standardized version of $\hat{\xi}(r)$, namely with $N=\min\{N_1,N_2\}$. The terms $\hat{\Omega}_{U,1}^{ipc}$ and $\hat{\Sigma}_{U}^{ipc}$ are defined in Appendix A.2 and are highly non-linear functions of the estimation errors of the factors in each group. As N and T grow large, under the generic null hypothesis $H_0(r):k^c=r$ we have, $\hat{\xi}(r)\stackrel{d}{\longrightarrow} N\left(0,1\right)$, while under the alternative hypothesis $H_1(r):k^c< r$, we have $\hat{\xi}(r)\stackrel{p}{\longrightarrow} -\infty$.⁶ Step 4: Once the dimension k^c is identified, we can recover the common factors f_{τ}^c via the

⁵When we reject the null $H(\underline{k})$ we look at the null hypothesis: $H(\underline{k}-1) = \{\rho_1 = ... = \rho_{\underline{k}-1} = 1\}$, and so forth until we identify the dimension of the common factor space. Sequential testing issues are addressed in AGGR.

⁶See Theorem 2 of AGGR, and its extension Theorem A.2 in Appendix A.2. The asymptotic distribution and rate of convergence of the test statistic $\tilde{\xi}(k^c)$ in Theorem A.2 are unchanged when the true numbers of factors k_1 and k_2 are unknown, and are estimated by some consistent empirical selection method.

canonical directions - i.e. the weights of the linear combinations yielding unitary canonical correlations - applied to the factors estimated from each of the separate panels.

It is worth highlighting a number of theoretical contributions of the paper. The theory in AGGR only covers panel data centered at zero and the PCA estimator, which is appropriate for a model with constant loadings. Appendix A.2 extends the estimators and theoretical results of AGGR to the models where i) factors are allowed to have any finite mean, and ii) loadings are time-varying. In particular, we cover RP_PCA and IPCA estimators for respectively the sorted portfolios and the individual stock panels. This general set-up is more relevant for asset pricing applications.

IV. Testing Test Assets

We implement the four-step procedure described in Section III for each of the non-overlapping sample blocks from Jan. 1966 until Dec. 2020 with 5-year increments. These are 5-year panels $y_{1,\tau}$ of individual stocks, and panels $y_{2,\tau}$ of test asset portfolios monthly returns. Broadly speaking we consider four monthly data sets in our analysis, namely (i) individual US stock returns from CRSP, (ii) the panel of test asset portfolios returns from the April 2021 release of the database "Open Source Cross-Sectional Asset Pricing" created by Chen and Zimmermann (2022), CZ21 hereafter, (iii) the panel of factors from the zoo considered by CZ21, and (iv) the

⁷More precisely, we have 11 panels ending at t = Dec. 1970, Dec. 1975, ..., Dec. 2020 with observations $y_{1,\tau}$ and $y_{2,\tau}$ for $\tau = t - 59, t - 58, \ldots, t$. We also report findings with full-sample estimates, although most of the discussion here will focus on the 5-year block samples. In *Supplementary Material* Section OA.1 we provide a detailed description of the data.

values of the 35 individual stock characteristics from Freyberger, Neuhierl, and Weber (2020) constructed using the new characteristics dataset of Jensen, Kelly, and Pedersen (2023).

The factors are estimated using IPCA of Kelly et al. (2019) for individual stocks implemented with the linear instruments for the time-varying loadings using the characteristics from Freyberger et al. (2020) and the Jensen et al. (2023) update. The second panel consists of quantile portfolios from Chen and Zimmermann (2022) where we use Lettau and Pelger's RP-PCs , fixing $\gamma_{RP}=-1$ to estimate factors: henceforth, we simply refer to them as PCs. Similar to Pukthuanthong et al. (2019) we decide to fix a priori the maximum number of pervasive factors (i.e. IPCs ans PCs) in each panel to 10.

A. Common Factors

Panel (a) of Figure 1 displays the first seven canonical correlations between 10 PCs from each panel. There is a noticeable gap between the top three canonical correlations and the remaining four. Applying our formal test we reject the null hypothesis that the number k_{τ}^{c} of common factors between the first 10 PCs of individual stocks and the first 10 PCs of portfolio test

⁸In the main body of the paper we report empirical results where models for sorted portfolios (resp. individual stocks) are estimated by RP_PCA with tuning parameter $\gamma_{RP} = -1$ (resp. IPCA). In the *Supplementary Material* we also report results of the same analyses performed by estimating the latent factors for CZ21 test assets by RP_PCA with different values of the tuning parameter, namely $\gamma_{RP} = +1, +5$, and +10, i.e. progressively increasing the relative importance of fitting the risk premia (i.e the realized means) instead of time-series variation of the returns of the CZ21 test assets, when estimating their latent factors. In unreported empirical results, appeared in previous version of this paper and available upon request, we have considered factors estimated by PCA on balanced non-overlapping panels of individual stock returns and obtained qualitatively similar results to those reported in the following sections.

assets is more than 3 for all 5-year blocks. While there is some variation we will proceed with the number of factors being equal to 3 and then investigate whether a fourth common factor has any significant impact on our findings. Henceforth we will refer to these three factors as the "common" factors and use the acronym 3CF.⁹

Panel (b) of Figure 1 displays the sum of the canonical correlations of the three common factors with (1) the 3 Fama and French factors (henceforth FF3 - blue circles, thick dotted line), (2) the 5 Fama and French factors (henceforth FF5 - black stars, thick dashed line). In addition, it also displays FF3/FF5 augmented with the momentum factor (FF3+Mom - blue circles, thin dotted line, FF5+Mom - black stars, thin dashed line). The red line across the plot marks the 3-factor benchmark common factor space.

[Insert Figure 1 approximately here]

The results in Panel (b) of Figure 1 convey a surprisingly simple and clear message. We observe that over the entire sample the canonical correlations between FF3 and common factors (blue circles, thick dotted line) are well below 3. This implies that over this sample period FF3 does not span the common factor space. What happens if we move from FF3 to FF5, i.e. we add RMW-operating profitability and CMA-investment style? In the same figure, using the same approach, the black thick dashed line with stars shows that adding two FF factors falls again short of spanning the common factor space. In fact, in most years the improvements of the two additional factors appears to be only minor. The same analysis is repeated with observable factors

⁹We also conducted the same exercise with the full sample where we only find one common factor. More specifically, the top three canonical correlations for the full sample are: 0.98, 0.72, and 0.59. While this would favor the time-varying beta CAPM, it is based on only roughly two thirds of the sorted portfolios that are available throughout the full sample (see *Supplementary Material Section OA.1* for further details).

FF3 and FF5 plus momentum (FF3+Mom - blue circles, thin dotted line, FF5+Mom - black stars, thin dashed line). While the higher number of observable factors increases mechanically the value of the sum of non-zero canonical correlations, it remains the case that adding the momentum factor is not enough to span the common factors. This means that the popular factor models fall short of capturing the three common factors. Andreou, Gagliardini, Ghysels, and Rubin (2024) derive a formal test between latent and observed factors similar in spirit to the test of AGGR covered in the previous section. Applying an extension of such a test, we reject the null that the FF5 and momentum factors span the $k^c = 3$ common factors - put differently all the lines in Panel (b) of Figure 1 are significantly below 3. Finally, this also begs the question whether members of the zoo might help us out in recovering the common factor space, a topic addressed in Section VI.

B. Comparing 3CF with Observable Factors

So far we established the existence of common factors between individual stocks and sorted portfolios, and shown that they are only partially spanned by FF5 and momentum. We now study how each of the FF5 and momentum observable factors taken one-by-one are related to (a) the common factors and (b) the factors which are specific respectively to the sorted portfolios or to the CRSP individual stocks. This analysis allows us to check (a) whether and to what degree the common factors, i.e. f_{τ}^c , are spanned by some of the observable factors, and (b) to understand the nature of the panel-specific factors: $f_{2,\tau}^s$ portfolio-specific and $f_{1,\tau}^s$ individual stock-specific factors as they were referred to in Section II.

To achieve the task at hand, we regress each of the 6 observable factors on (i) the 3CF

¹⁰In Section VI we further discuss the plausible interpretation of the three common factors in terms of their correlations with hundreds of observable factors collected by CZ21, i.e. the factor zoo, as well as macro factors.

factors and (ii) all the panel-specific factors, and report the R^2 s of these regression for FF3, RMW, CMA and momentum in *Supplementary Material* Figures OA.2 and OA.3. Not surprisingly, typically over 90% of the variability of the market factor is explained by the common factors, and $f_{2,\tau}^s$ portfolio-specific factors tend to explain almost all the remaining part of its variability, with their R^2 s ranging between 2% and 20%, depending on the time period. We also find that the factors specific to individual stocks are not able to capture the same amount of the variability unexplained by the common factors, as their R^2 s are typically below 10%.

The fact that CRSP stock-specific factors explain much less of the variability of observable factors compared to the $f_{2,\tau}^s$ portfolio-specific factors is an empirical regularity that we observe across all the FF5 and momentum factors. We also note that on average only about 50% of the variability of SMB is explained by the 3CF, and the remaining 50% is mostly portfolio sorted specific factors that explain SMB. For HML a similar pattern appears and analogous conclusions can be drawn for RMW, CMA and momentum as detailed in the *Supplementary Material*.

The finding implies that out of the 6 factors considered, only the market is predominantly related to the common factors, while all the other FF and momentum factors are only partially spanned by 3CF, and a large part of their variability is due to a risk dimension which is specific to portfolio sorting.

We can examine the same question from a different angle. Recall that we started with 10 IPCs for the panel of individual stocks, 10 PCs for the panel of test asset portfolios and found three common factors. Therefore, we have 7 remaining group-specific factors in each panel. Figure 2 displays the sum of the canonical correlations of the 7 group-specific factors (Panel (a) the CZ21- $f_{2,\tau}^s$ portfolio-specific and in Panel (b) CRSP-specific $f_{1,\tau}^s$) with FF3 (blue circles, thick dotted line), FF5 (black stars, thick dashed line), FF3 factors and momentum (blue circles, thin

dotted line), and FF5 factors and momentum (black stars, thin dashed line). Recall that the group-specific factors reflect (1) specification errors in the time variation of betas and (2) characteristic-sorting mixed with beta dynamics. Since we showed that the market factor is basically one of the 3CF, we expect for FF3 at most a sum of canonical correlations equal to 2, and for FF5 equal to 4, adding momentum increases both numbers by one. Panels (a) and (b) of Figure 2 show that the FF factors with or without momentum are well below those numbers. The sum of canonical correlations is equal to about one for FF3, two for FF3 plus momentum and FF5 and finally the sum equals about three for FF5 plus momentum. The $f_{2,\tau}^s$ portfolio-specific factors appearing in Panel (a) yield higher sums of canonical correlations, particularly when momentum is added. This finding is perhaps not surprising, since the FF factors are constructed by sorting.

If we combine the findings in Figures 1 and 2 in the it appears that one linear combination of FF3 (putting most of the weight on the market) is perfectly correlated with one common factors. Adding the evidence from Figure OA.2 in the *Supplementary Material* it appears that another linear combination of FF3 highly correlates with a sorting-specific factor (Panels (c) through (f) in Figure OA.2 suggest this is a combination of SMB and HML).

V. Asset Pricing Performance of Common Factors

How do the common factors perform as predictors? How does their performance compare with widely used factors in the asset pricing literature? These are questions we address in this section.

A. Empirical Models

We model the excess returns of CRSP individual stocks and test asset sorted portfolios as linear functions of different sets of K factors. In particular, we consider the following sets of factors:

- (i) FF + mom: Under this header we have a set of models starting with the market factor only, defined as the value-weighted index of all CRSP stocks minus the risk-free (K = 1); FF3 and FF5 only (K = 3, 5); FF3 + Momentum, FF5 + Momentum, (K = 4, 6);
- (ii) 3CF: $K = k^c = 3$ common factors;
- (iii) 3CF + CRSP-spec.: $k^c = 3$ common factors and $k_1^s = 1, 2, 3$ group-specific factors from the panel of CRSP stocks (K = 4, 5, 6);
- (iv) 3CF + CZ21-spec.: $k^c = 3$ common factors and $k_2^s = 1, 2, 3$ group-specific factors from the panel of CZ21 test assets portfolios (K = 4, 5, 6);
- (v) *PCA on CZ21*: factors estimated as the first K=1,3,4,5,6 PCs on CZ21 test assets portfolios;
- (vi) IPCA on CRSP: factors estimated as the first K = 1, 3, 4, 5, 6 IPCs on CRSP individual stocks.

All the models, factors and betas/loadings are estimated from the 60 monthly returns in each of the eleven non-overlapping 5-years blocks ending in year t, with t=1970,...,2020. Therefore, loadings/betas of PCA, and the matrix Γ mapping stock-specific characteristics to IPCA loadings, are constant for all the dates τ within each non-overlapping block, but are allowed to change across different blocks.¹¹

¹¹See definition of the IPCA estimator and its parametrization in *Supplementary Material* Section OA.3.

B. Performance Evaluation Measures

We describe the in-sample and out-of-sample performance evaluation measures.¹²

In-sample performance evaluation

We compute the following in-sample performance measures across the entire sample, that is across all B blocks:

- $Total\ R^2$ of Kelly et al. (2019) which represents the fraction of return variance for all the assets explained by both the dynamic behavior of the loadings (for IPCA) and the contemporaneous factor realizations across different blocks, aggregated over all assets and all time periods.
- $Pricing\ error\ R^2$ of Kelly, Palhares, and Pruitt (2023) which pertains to the fraction of the squared unconditional mean excess returns, i.e. risk premium, that is described by factors and betas. This metric is close to (one minus) the numerator of formal tests (like the GRS test) for the null hypothesis that the pricing errors for the test assets are zero.
- Predictive R^2 from Kelly et al. (2019) which represents the fraction of realized return variation explained by the model's description of conditional expected returns, and summarizes the model's ability to describe risk compensation only through exposure to systematic risk.

Additionally, we test the significance of the CZ21 group-specific factors in explaining the cross-section of returns for the CZ21 test assets when added to the 3 common factors in

¹²The technical details appear in *Supplementary Material* Section OA.6.

Fama-MacBeth regressions.¹³ Failure to reject the null, implies that the 3CF suffice in explaining the cross-section of risk premia of the CZ21 portfolios. As we re-estimate the loadings and the factors in each one of the eleven five-years blocks, this implies that the regressors in the Fama-MacBeth cross-sectional regressions change across blocks, and therefore we perform the test in each block.

Out-of-sample performance evaluation

We implement the out-of-sample version of the $Total\ R^2$, $Pricing\ R^2$ and $Predictive\ R^2$ with betas and factor loadings computed using information from block ending in year t-5 to price month τ assets in the following 5-years block ending in year t, with $t=1975,\ldots,2020$. Analogously to Lettau and Pelger (2020b) we also compute the annualized Sharpe ratio (SR) of the "Maximum Sharpe-ratio portfolio" that can be obtained by an optimal (in a mean-variance sense) linear combination of the factors, which are ultimately portfolios of individual stocks. The

¹³ Using estimated factors instead of true unobserved factors as regressors does not affect asymptotically the R^2 measures because the factor estimates are consistent. Indeed, the estimation errors for factors obtained by PCA or IPCA on panel ℓ are of the order $O_p(1/\sqrt{N_\ell}+1/T)$ and shrink to zero when both panel dimensions grow, see e.g. Bai and Ng (2006). One can show that the same holds for common factors obtained from canonical analysis applied to PCs or IPCs; we refer to AGGR and the OA for details. In fact, the use of estimated factors from large cross-sections in the first pass of the Fama-MacBeth procedure is also inconsequential for inference on risk premia asymptotically. Indeed, the estimation errors on factor values induce a supplementary term at order $O_p(1/\sqrt{N_\ell}+1/T)$ in beta estimates. However, when $N_\ell \gg T$, this term is negligible with respect to the usual term at order $O_p(1/\sqrt{T})$ from first-pass time-series regression, which yields the well-known Error-in-Variable (EIV) problem in the second-pass; we refer to Gagliardini, Ossola, and Scaillet (2016) for risk premia estimation with large N_ℓ, T .

out-of-sample performance measures are defined as: (a) OOS Total \mathbb{R}^2 , (b) OOS Pricing \mathbb{R}^2 , (c) OOS Predictive R^2 , and (d) Maximum Sharpe-ratio, OOS SR.

In- and Out-of-sample Performance Evaluation of Factor Models

The objective of this subsection is to compare the role of the factor model specifications in explaining the variation of returns for individual CRSP stocks as well as CZ21 test assets, both inand out-of-sample, during the period 1966-2020. Panels A - C in Table 1 present respectively the Total, Pricing and Predictive \mathbb{R}^2 evaluation measures. The first two rows in each panel pertain to the benchmark models which consist of the FF and momentum factors, starting with the one-factor market (CAPM), then FF3, FF3 with momentum, FF5 and finally FF5 plus momentum. The columns are therefore labeled 1, 3, 4, 5 and 6 corresponding to the number of factors K in each model. For comparison, the next two rows refer to the corresponding R^2 s of the three common factors (3CF). These are followed by the models which consider both the three common factors as well as the panel-specific factors (individual CRSP stocks and CZ21 test assets) in order to evaluate whether the latter have additional explanatory power for the variation of returns beyond the three common factors. Finally, the last four rows in each panel of Table 1 pertain to the R^2 s for models with factors based on the PCs from each of the two panels as well as the PCs across the two panels (namely IPCs from CRSP are used to price CZ21 assets, assuming constant loadings estimated by OLS regressions of the CZ21 test assets on IPCs, and vice versa).

[Insert Table 1 approximately here]

From the in-sample analysis reported in Table 1 we can draw the following two important observations. First, the three common factors typically yield better or comparable in-sample Total, Pricing and Predictive R^2 s vis-à-vis the benchmark models (i.e. CAPM K=1, FF3 K=3, FF3 plus momentum, K = 4, FF5 K = 5 and FF5 plus momentum K = 6). There are a number of cases, particularly involving CRSP individual stocks data, where the traditional models do better in-sample. Second, adding the corresponding group-specific factors from the two panels (of individual stock and test assets) leads only to marginal improvements compared to models with 3CF. Of greater interest are the out-of-sample results in Table 1. They yield the following key empirical findings:

- the three common factors yield the highest OOS Total, Pricing and Predictive R^2 s compared to *any* FF (plus momentum) benchmark model (with up to K = 6). The relative gains of the OOS Total, Pricing and Predictive R^2 s are often orders of magnitude larger particularly for individual stocks;
- the 3CF model does much better than IPCA on CRSP with any number of factors and any metric (with the exclusion of K=6 for Pricing R^2 only) in explaining CZ21 returns and does as good as or better than IPCA on CRSP with 3 IPCs according to any metric. The differences are quite substantial in many cases. The FF type factors applied to CRSP reach a max of roughly 16% Total R^2 whereas 3CF yields almost 25% and applying IPCA on CRSP gets similar results only with 5 factors. More dramatic is the OOS Pricing R^2 , where variations on FF reach a max of 8% while 3CF yields almost 35%. Applying IPCA in this case is comparable to 3CF;
- when we use (RP_PCA) PCs from portfolios (CZ21 assets) to explain the risk premia of the CRSP stocks out of sample the OOS performance is extremely poor: the Pricing R^2 in Table 1 (and *Supplementary Material* Table OA.3) is among the lowest among all the

models, and the lowest compared to other models with latent factors. In fact, the OOS Pricing R^2 decreases when more PCs from the CZ21 portfolios are added;

- adding to 3CF the corresponding panel-specific factors (from individual stock and CZ21 assets) leads sometimes to only marginal improvements;
- RP_PCA on CZ21 yields slightly better results than 3CF for OOS Total and Pricing (but not Predictive) R^2 s for CZ21 returns. However, those factors poorly predict individual stocks out-of-sample according to the three types of R^2 s considered.

[Insert Table 2 approximately here]

Next we test, at least in-sample, using the Fama-MacBeth procedure whether the CZ21 group-specific factors explain beyond 3CF, the risk premia of the CZ21 portfolios. The results appear in Panel A of Table 2.¹⁴ Testing for the joint significance of the risk premia - estimated in cross-sectional regressions - of the three CZ21-specific factors when added to the 3CF, is an

¹⁴We conduct the test for each one of the 5-years blocks for which we estimate our model, by computing the risk premia of the factors with Weighted Least Squares (WLS) cross-sectional regressions, and estimate their covariance matrix by applying the Fama-MacBeth procedure. The variance covariance matrix of estimated factors' risk premia in Table 2 are computed using the Fama-MacBeth procedure, with 3-lags (optimally selected) Newey-West weighting. In unreported results available upon request, we verify that our results are robust to different choices of the number of lags for the Newey-West weighting. As a further robustness check, we also consider ordinary least squares cross-sectional regressions, instead of WLS, and obtain similar results, although for one block out of eleven the null hypothesis of the insignificance of the risk premia of the CZ21-specific factors is not rejected at 1% significance level (p-value = 0.003), and in another block the null is not rejected at the 5% significant level (p-value = 0.036). For all the other nine 5-years blocks the joint test of significance of the risk premia is not rejected at 5% significance level.

indirect way to test for the significance of the slight increase from 93.0 to 95.7 observed for the Pricing R^2 in Panel B of Table 1.¹⁵

Interestingly, the risk premia of the three CZ21-specific factors considered in Table 1, are never significant, indicating the CZ21-specific factors seem not to be relevant to explaining the (large) cross-section of risk-premia of the CZ21 portfolios that we consider, when added to the 3 common factors. The results in Panel A of Table 2 confirm that across all subsamples, the addition of the CZ21-specific factors has no statistically significant impact on risk premia.

D. Macro Factors and 3CF

So far we focused on factors extracted from equity returns. Macroeconomic factors are often cited as drivers of stock prices and in Panel B of Table 2 we examine whether macro indicators/factors contain relevant pricing information beyond that embed in the three common factors. We consider a wide range of macro indicators/factors. This includes fixed income market factors, such as term spreads 10y-1y, 10y-3m, 1y-FEDfund, default spread BAA-AAA and the Cochrane-Piazzesi factor. It also covers various key real economy series such as consumption, labor income growth, various measures of industrial production and principal components extracted from panels of macroeconomic indicators. Finally, different inflation measures (including expected and unexpected inflation), the Baker, Bloom, and Davis (2016) economic policy uncertainty and Jurado, Ludvigson, and Ng (2015) macroeconomic uncertainty measures are also included. There are a total of 219 statistics in Panel B of Table 2 covering all

 $^{^{15}}$ To the best of our knowledge no formal test exists, yet, for the increase of the $R^{2\prime}$ s displayed in Table 1. The construction of such tests, although being on our research agenda, is beyond the scope of this paper.

¹⁶The details appearing in *Supplementary Material* Table OA.1.

the series considered across the 11 non-overlapping 5-year samples. Under the null that none of the macro factors considered is significant, we expect to find (in large independent samples) 2 (and there is one in Panel B), 10 (and there are two)and 20 (and there are seven) rejections of the null at respectively the 1, 5 and 10 percent level. Of course, these are not large and independent samples, but the few significant tests suggest that overall the macroeconomic factors do not provide incremental pricing gains beyond what is embedded in the 3CF. Put differently, although the three common factors are extracted from the cross-section of asset returns - individual stocks and sorted portfolios - they do incorporate the macroeconomic drivers of stock returns.

E. Out-of-sample Portfolio Performance

Last but not least, Table 3 presents the out-of-sample annualized maximum Sharpe ratio for the different factor model specifications. Interestingly, we find that the three common factors perform relatively better producing a SR of 0.66, vis-à-vis the CAPM and FF3 models which have SRs of 0.39 and 0.30, respectively. Nevertheless, the FF5 factor model plus momentum (K = 6) produces the highest SR of 0.81 among the traditional benchmark models. If we limit ourselves to K = 3, then the common factor model outperforms all other specifications in Table 3 except for IPCA on individual stocks. Adding factors beyond K = 3 does increase SR, however, and the best model is IPCA on CRSP with K = 6.

[Insert Table 3 approximately here]

When compared with models with a larger number of factors, the 3CF model is not consistently the best according to the different metrics we consider across both panels of test assets. In particular, an IPCA model with K=6 factors produces a 7.6% (resp. 0.6%) higher out-of-sample Pricing \mathbb{R}^2 , when compared to the 3CF for the panel of individual stocks (resp. CZ21 portfolios).

Nevertheless, when CRSP-specific factors are added to the 3CF some metrics, such as the OOS Pricing and Predictive R^2 , show that the model with K=6 IPCs is slightly inferior to explain CZ21 portfolio returns. These results, together with the Total R^2 of a model for individual stocks with K=6 IPCs being systematically inferior to one of other models with either 3 or 6 factors, and the high SR of the 6 IPCs in Table 3, could be explained by the fact that the higher order IPCs - especially the 6^{th} - are able to explain the returns of a relatively small set of stocks with very high realized risk premia, and possibly low return variability.

We repeat the empirical analyses summarized in Tables 1 - 3 with one notable difference, namely we estimate the latent factors for CZ21 test assets by RP_PCA with the tuning parameter γ_{RP} taking values +1, +5, and +10, and report the results in Supplementary Material Section OA.5.1. We find that, overall, by progressively increasing the relative importance of fitting the risk premia when estimating the latent factors on the CZ21 test assets, the results reported in Tables 1 - 3 (produced using classical PCA, i.e. RP_PCA with $\gamma_{RP}=-1$) are mostly robust to different values of γ_{RP} , with some exceptions that do not change the main message. As expected, by increasing γ_{RP} RP-PCs explain better the risk premia of CZ21 test assets as reflected by the slight increases in their in-sample Pricing R^2 of between 1% and 1.5%, depending on the number of RP-PCs considered and values of γ_{RP} . More surprisingly, the RP-PCs from the CZ21 assets improve the in-sample Pricing R^2 for the CRSP individual stocks from 2.4% to 11.3%. Nevertheless, these improvements in pricing abilities of the RP-PCs almost completely disappear in the out-of-sample analysis, where RP-PCs seem to perform very similarly to PCs. The models with 3CF (with or without group-specific factors) have improvements of the Pricing R^2 ranging from -0.2%, i.e. a deterioration, to 0.5%. On the other hand, the predictive R^2 generated by the models with RP-PCs decreases slightly for almost all specifications. Additionally, Fama

MachBeth regressions of CZ21-specific factors obtained by RP-PCs confirm the results obtained by PCs, that the 3CF are enough to explaining the risk premia of CZ21 test assets, with the exception of only one 5 year block (different for different values of γ_{RP}), where the p-value of the test is 4%. Finally, the out-of-sample Sharpe Ratios of the 3CF for the first 3 RP-PCs do not improve when RP-PCs are considered instead of classical PCs, but improvements are noticed when adding the higher order PCs from the fourth onward and CZ-21 -specific factors based on RP-PCs, but not for the CRSP-specific factors. Overall, we conclude that the nature and the properties of 3CF do not change when we consider RP-PCs instead of PCs in the estimation of latent factors in the CZ21 test assets.

The results in Table 2 also inform us about whether augmenting the 3 common factors by a fourth one (4CF) improves the pricing of risk. Starting with individual stocks, we can read these results as follows: the third row, CRSP/3CF K=3 needs to be compared with the fifth row K=4 as we consider one group-specific factor as common. This means that for example the out-of-sample Total R^2 increases from 25.5 to 26.4. Similarly, for sorted portfolios we need to compare row four, CZ21/3CF with the sixth where for example the out-of-sample Total R^2 moves from 91.2 to 92.0. The results in Table 2 can be used here again to inform us of the fact that moving from 3CF to 4CF has no significant impact on pricing. Hence, while for some of the subsamples there might be evidence of more than 3 common factors, it does not seem to matter in terms of pricing the cross-section of stocks and sorted portfolios.

Finally, as mentioned above, it is worth noting that the empirical results reported in Sections IV and V, but also those in the next Section VI remain qualitatively similar when PCA is used to extract factors on balanced panels of individual stocks returns, instead of IPCA (results appeared in previous versions of this paper and are available upon request). Therefore, it is

unlikely that our identification of the 3CF is coming simply from the differences of the smaller set of characteristics used to estimate the loadings in the IPCA model, and those used by CZ21 to construct the sorted portfolios used in our analysis as second panel of test assets.

VI. Plausible Interpretations of 3CF and the Factor Zoo

We address two related questions in this section. First we explore the relationship between 3CF and the factor zoo. Second, we try to understand which factors from the zoo might help us with finding an economic interpretation of 3CF. The factor zoo is represented by all the factors collected by CZ21.¹⁷

A. Common Factors and the Zoo

We investigate first the ability of PCs extracted from the factor zoo, which we call zoo PCs, to price and explain the variability of the panels of individual stocks and the CZ21 portfolios. Figure 3 shows the sum of the canonical correlations of the three common factors with the first 3, 5, 10 and 15 zoo PCs. As before, all PCs and common factors are estimated from non-overlapping 5-year balanced panels of monthly data over the period 1966-2020. The figure allows us to understand whether the PCs from the zoo panel span the space of the common factors. The first observation emerging from Figure 3 is that 3 zoo PCs yield a sum of canonical correlations between 2 and 2.5 as if there is constantly a missing factor. Going to 5 zoo PCs gets us above 2.5, but it takes up to 10 PCs from the factor zoo to approximately span the set of 3 common factors.

¹⁷As detailed in *Supplementary Material* Section OA.1, when we refer to the factor zoo we use a data set of over a thousand portfolios associated with 205 characteristics.

[Insert Figure 3 approximately here]

[Insert Table 4 approximately here]

Table 4 documents which factors from the zoo are the most related (a) to 3CF, and (b) to 3CF augmented with the first three group-specific factors of the CZ21 portfolios. Panel A reports the twenty factors with the largest average R^2 s - across all non-overlapping windows - when regressed on the 3CF factors. Among them, we find two of the three Fama-French factors (CAPM Beta and Size) along with portfolios based on market beta put forward by Frazzini and Pedersen (2014) and Price (which can be considered a cruder alternative to the Size factor). ¹⁸ In addition, among the zoo factors we also find different measures of idiosyncratic risk and liquidity or uncertainty, such as bid-ask Spread, cash-flow to price variance, volume to market equity, trading volume, EPS Forecast Dispersion, days with zero trades, Volume Variance, and Price delay R-squared. Interestingly, book-to market using the most recent market equity also appears among one of the most correlated factors with the 3CF (with an average R^2 of 77.1%), nevertheless it is worth mentioning that the "classical" book-to-market factor of Fama and French (1992) - which uses as denominator the previous year December value of the equity - is not on the list, although across the entire sample it has a correlation of 0.58. This implies that the difference between two similar, but not identical factors, matters in their ability to explain the 3CF.¹⁹

It is also interesting to note that idiosyncratic risk measures (i.e. those suggested by Ang, Hodrick, Yuhang, and Zhang (2006) and Ali, Lee-Seok, and Trombley (2003)), a factor reflecting market frictions such as the betting-against beta of Frazzini and Pedersen (2014),

¹⁸The Price and the Size factors have a correlation of 0.78 computed over the entire sample in our analysis.

¹⁹And makes the results compatible with those in Panel (e) of *Supplementary Material* Figure OA.2.

i.e. Frazzini-Pedersen Beta factor in Table 4, and a factor reflecting behavioral biases, in particular the factor reflecting expectations of future earning growth by La Porta (1996), i.e. Long-term EPS forecast in Table 4, are among the portfolio strategies that Dello-Preite et al. (2024) find to be related, although not perfectly, to their "unsystematic risk component", i.e. the third component of their SDF which added to two other systematic factors prices well the returns of a large cross-section of sorted portfolios. Turning to Panel B, we report the factors in the zoo showing the highest increase in the \mathbb{R}^2 - averaged across all non-overlapping years windows - when the first 3 CZ21-specific factors are added as regressors to 3CF. Interestingly, the majority of the factors which are the most related to the three CZ21 group-specific factors are associated with momentum.

Overall, these findings complement the results of *Supplementary Material* Figures OA.2 and OA.3 confirming that the market and size are the factors most correlated with 3CF (for size this is especially true for the first half of our sample), while the majority of the variability of momentum is mostly explained by CZ21-specific factors.

Table 5 reports the performance evaluation measures for the zoo PCs. For convenience of comparison we repeat the results for 3CF from Panels A - C in Table 1. The zoo PCs have positive in- and out-of-sample Total R^2 s for individual stocks which is better than the observable factors, but worse than all the other latent factor models (comparing with results in Table 1). Moreover, they perform worse than the other set of factors in explaining CZ21 portfolio returns. The pricing performance of the zoo PCs is also the worst, as evident from their low or negative Pricing R^2 appearing in Table 5.

[Insert Table 5 approximately here]

The results of Table 5 are obtained by estimating factors from the returns of the factors in the Zoo

assets by classical PCA. In Tables OA.7 - OA.9 in Supplementary Material Section OA.2 we repeat the same analyses by estimating the latent factors in the Zoo by RP_PCA setting the tuning parameter $\gamma_{RP} = +1$, +5 and +10, and the results are slightly worse that those reported for PCA, with the exception of the out-of-sample pricing R^2 which increases slightly for CZ21 test assets when explained by the 3CF only.

Overall, these results for (RP-)PCs obtained from the Zoo of factors are not surprising. In fact, RP-PCs are meant to capture most of the variation and, (implicitly or explicitly) the risk premia of the assets from which they are estimated. The factors in the Zoo are, by construction, very special linear combinations of the CZ21 test assets (which are themselves special linear combinations of returns of individual stocks). Therefore, fitting the risk premia of a special combination (portfolio) of assets does not guarantee that the risk premia of all the components of the linear combinations, or assets not included in the linear combination themselves.

Finally, we also report on an exercise where we pick two factors from the zoo in addition to the market and find the highest canonical correlation among any two combined with the market with 3CF. To illustrate for 200 (say) factors in the zoo (other than the market) there are 19,900 = (200!/(2!198!) possible three observable factor models with the market and two other zoo members. For these 19900 combinations we compute the sum of the canonical correlations with the three common (latent) factor and look at the max among all 19,900. The analysis is repeated for each of the five-years non-overlapping blocks and reported in Table 6. The size factor is the most prominent (again) featured in the table, particularly in the early part of the sample. Other than that there does not seem to be a consistent pattern. This may explain why there are so many factors in the zoo. While the 3CF are parsimonious, they are latent, and do not seem to consistently relate to observable factors. Also noteworthy are the numbers in the last column.

They reflect the maximal sum of correlations attained by each of the selected combinations. Applying the test of Andreou et al. (2024) of observable versus latent factors rejects the null hypothesis that these sums are equal to three. This means we need to search for at least four observable factors to match 3CF.

[Insert Table 6 approximately here]

B. On the Macroeconomic Factors Embedded in 3CF

Recall that in Panel B of Table 1 we examined whether macro factors contain relevant pricing information beyond that embedded in the three common factors. We found that although the three common factors are extracted from the cross-section of asset returns - individual stocks and sorted portfolios - they do incorporate the macroeconomic drivers of stock returns. This begs the question which macroeconomic factors are mostly related to 3CF. In Table 7 we report an exercise similar to that in Table 6 but instead look ad the maximum canonical correlations between macroeconomic indicators/factors and 3CF. The table reports the highest canonical correlation among any three combined macro factors (Panel A) or two macro factors and market (Panel B) with 3CF for each of the 11 five-years non-overlapping blocks. The numbers reflect the maximal sum of correlations attained by each of the selected combinations. In Panel A, where we only look at macro factors, we note that most canonical correlations are barely above one, with on exception the last subsample where it reaches 1.75. The macro PCs innovations (namely, the MacCracken and Ng (M&N) MacroFRED MD PCs and the Ludvigson and Ng (L&N) Macro PCs VAR innovations) PCs appear most prominently across the subsamples. Panel B, where the market return is paired with two macro factors, we see that the canonical correlations

increase by about 0.50 on average, which is still substantially below the maximum canonical correlations between observed zoo factors and 3CF reported in Table 6.

[Insert Table 7 approximately here]

The existing (factor zoo) literature has added 200+ observable factors to explain the cross-section of returns. From the evidence in Tables 6 and 7, we must conclude that the 3CF embed macroeconomic signals as well as information from the cross-section of returns, but combinations of a small number of observable factors from either stock returns or the real economy do not fully capture the 3CF risk factors.

VII. Conclusions

For a number of decades we have been debating the selection of test assets, with arguments pro and con individual stocks versus sorted portfolios. In addition, we also have been debating about the drivers of time-varying betas. The fact many studies resort to sieve and neural network type searches for functional specifications for time-varying loadings is indicative of the complexity of the task. Our paper proposes an alternative to account for the well-known challenges in specifying beta dynamics of single stocks, different from the 'machine-learning' trend (which uses complex beta dynamics, for which however there may be issues of overfitting, robustness, etc.) or recent approaches that use short time series panel methods (which however focus on short rolling time-windows instead of a global picture of the factor structure).

The argument put forward in our paper is that we should not look at *either* individual stocks *or* sorted portfolios separately. We should consider the information from both of these panels and find the factor space that is common among the two. Doing so, shields us from the

thorny issue of specifications errors in beta dynamics. The task of finding the common factor space is non-trivial and one of the key contributions of our paper. We extract factors from both panels and find the common factor space between the two panels, yielding factors which price both individual stocks and sorted portfolios. We labeled these three common factors 3CF. We show that this provides a path toward extracting factors neither affected by sorting characteristics nor by varying risk exposures and recalcitrant features of individual stocks.

The projection arguments put forward in Hansen and Jagannathan (1991) imply, as noted by Kozak, Nagel, and Santosh (2018), that there exists a factor representation of the stochastic discount factor (SDF). Moreover, there is practically no disagreement that the space of factors spanning the SDF is low-dimensional. In this paper we found 3 factors which were selected via a novel procedure addressing a longstanding debate in the empirical asset pricing literature. These three factors are related to the first two FF factors, namely Market and Size, but also to an idiosyncratic/unsystematic risk factor, similarly to Dello-Preite et al. (2024), liquidity and uncertainty proxies. Importantly, 3 PCs from a large panel of sorted CZ21 portfolios, or 3 IPCs extracted from individual stock returns cannot explain both panels of stock returns and the 3CF. Adding more panel-specific factors to either the 3CF or the 3 PCs or IPCs increases only marginally the ability of the models to explain both panels of returns.

Last but not least, it should be noted that the testing procedure introduced in our paper can be applied in many other asset pricing settings. A few examples are: comparing panels of private equity and publicly traded companies, international asset pricing comparing systematic risks in stock returns in different countries, among others.

References

- Ali, A.; W. Lee-Seok; and M. A. Trombley. "Arbitrage Risk and the Book-to-Market Anomaly." *Journal of Financial Economics*, 69 (2003), 355–373.
- Andreou, E.; P. Gagliardini; E. Ghysels; and M. Rubin. "Inference in Group Factor Models with an Application to Mixed-frequency Data." *Econometrica*, 87 (2019), 1267–1305.
- Andreou, E.; P. Gagliardini; E. Ghysels; and M. Rubin. "Spanning Latent and Observable Factors." (2024). *Journal of Econometrics*, (forthcoming).
- Ang, A.; R. J. Hodrick; X. Yuhang; and X. Zhang. "The Cross-Section of Volatitlity and Expected Returns." *Journal of Finance*, 61 (2006), 259–299.
- Bai, J., and S. Ng. "Confidence Intervals for Diffusion Index Forecasts and Inference for Factor-Augmented Regressions." *Econometrica*, 74 (2006), 1133–1150.
- Baker, S. R.; N. Bloom; and S. J. Davis. "Measuring Economic Policy Uncertainty." *Quarterly Journal of Economics*, 131 (2016), 1593–1636.
- Bryzgalova, S.; M. Pelger; and J. Zhu. "Forest Through the Trees: Building Cross-Sections of Stock Returns." (2024). *Journal of Finance* (forthcoming).
- Chen, A. Y. "The Limits of p-Hacking: A Thought Experiment." (2019). Available at SSRN https://ssrn.com/abstract=3272572.
- Chen, A. Y., and T. Zimmermann. "Open Source Cross-Sectional Asset Pricing." *Critical Finance Review*, 11 (2022), 207–264.

- Chen, Q.; N. Roussanov; and X. Wang. "Semiparametric Conditional Factor Models: Estimation and Inference." (2023). NBER Working Paper 31817.
- Chorida, T.; A. Goyal; and A. Saretto. "Anomalies and False Rejections." *Review of Financial Studies*, 35 (2020), 2134–2179.
- Cochrane, J. H. "Presidential Address: Discount Rates." *Journal of Finance*, 66 (2011), 1047–1108.
- Dello-Preite, M.; R. Uppal; P. Zaffaroni; and I. Zviadaze. "Cross-Sectional Asset Pricing with Unsystematic Risk." (2024). Available at SSRN

 https://ssrn.com/abstract=4135146.
- Fama, E. F., and J. D. MacBeth. "Risk, Return, and Equilibrium: Empirical Tests." *Journal of Political Economy*, 81 (1973), 607–636.
- Feng, G.; S. Giglio; and D. Xiu. "Taming the Factor Zoo." *Journal of Finance*, 70 (2020), 1327–1370.
- Frazzini, A., and L. H. Pedersen. "Betting against Beta." *Journal of Financial Economics*, 111 (2014), 1–25.
- Freyberger, J.; A. Neuhierl; and M. Weber. "Dissecting Characteristics Nonparametrically." *Review of Financial Studies*, 33 (2020), 2326–2377.
- Gagliardini, P., and H. Ma. "Extracting Statistical Factors when Betas are Time-Varying." (2022).

 Swiss Finance Institute Research Paper.

- Gagliardini, P.; E. Ossola; and O. Scaillet. "Time-Varying Risk Premium in Large Cross-Sectional Equity Data Sets." *Econometrica*, 84 (2016), 985–1046.
- Ghysels, E. "On Stable Factor Structures in the Pricing of Risk: Do Time-Varying Betas Help or Hurt?" *Journal of Finance*, 53 (1998), 549–573.
- Hansen, L. P., and R. Jagannathan. "Implications of Security Market Data for Models of Dynamic Economies." *Journal of Political Economy*, 99 (1991), 225–262.
- Harvey, C. R., and Y. Liu. "A Census of the Factor Zoo." (2019). Available at SSRN https://ssrn.com/abstract=3341728.
- Harvey, C. R.; Y. Liu; and H. Zhu. "... and the Cross-Section of Expected Returns." *Review of Financial Studies*, 29 (2016), 5–68.
- Hou, K.; C. Xue; and L. Zhang. "Replicating Anomalies." *Review of Financial Studies*, 35 (2020), 2019–2133.
- Jensen, I.; B. Kelly; and L. Pedersen. "Is There a Replication Crisis in Finance." *Journal of Finance*, 78 (2023), 2465–2518.
- Jensen, M. C.; F. Black; and M. S. Scholes. In *Studies in the Theory of Capital Markets*. "The Capital Asset Pricing Model: Some Empirical Tests.", M. C. Jensen, ed. Praeger Publishers Inc (1972).
- Jurado, K.; S. C. Ludvigson; and S. Ng. "Measuring Uncertainty." *American Economic Review*, 105 (2015), 1177–1216.

- Kelly, B.; D. Palhares; and S. Pruitt. "Modeling corporate bond returns." *Journal of Finance*, 78 (2023), 1967–2008.
- Kelly, B.; S. Pruitt; and Y. Su. "Instrumented Principal Component Analysis." Working Paper.
- Kelly, B. T.; S. Pruitt; and Y. Su. "Characteristics are Covariances: A Unified Model of Risk and Return." *Journal of Financial Economics*, 134 (2019), 501–524.
- Kim, S., and R. A. Korajczyk. "Large sample estimators of the stochastic discount factor." *Journal of Financial Econometrics*, 22 (2024), 1672–1713.
- Kozak, S.; S. Nagel; and S. Santosh. "Interpreting Factor Models." *Journal of Finance*, 73 (2018), 1183–1223.
- La Porta, R. "Expectations and the Cross-Section of Stock Returns." *Journal of Finance*, 51 (1996), 1715–1742.
- Lehmann, B., and D. Modest. "The Empirical Foundations of the Arbitrage Pricing Theory." *Journal of Financial Economics*, 21 (1988), 213–254.
- Lettau, M., and M. Pelger. "Estimating Latent Asset-Pricing Factors." *Journal of Econometrics*, 218 (2020a), 1–31.
- Lettau, M., and M. Pelger. "Factors That Fit the Time Series and Cross-Section of Stock Returns." *Review of Financial Studies*, 33 (2020b), 2274–2325.
- Litzenberger, R. H., and K. Ramaswamy. "The Effect of Personal Taxes and Dividends on Capital Asset Prices: Theory and Empirical Evidence." *Journal of Financial Economics*, 7 (1979), 163–195.

Magnus, J. R., and H. Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley and Sons: Chichester/New York. (2007).

McLean, R. D., and J. Pontiff. "Does Academic Research Destroy Stock Return Predictability?" *Journal of Finance*, 71 (2016), 5–32.

Pukthuanthong, K.; R. Roll; and A. Subrahmanyam. "A Protocol for Factor Identification." *Review of Financial Studies*, 32 (2019), 1573–1607.

 $\label{eq:TABLE 1}$ In- and Out-of-sample Performance Evaluation of Factor Models

Panels A - C of the table report Total, Pricing and Predictive R^2 s in percent for observable factor models (lines 1-2 in each panel), a latent factor model with 3 factors common between individual stocks and CZ21 portfolios (lines 3-4 in each panel), the same 3 common factors together with 1, 2, or 3 CRSP-specific factors (line 5 in each panel), again the same 3 common factors together with 1, 2, or 3 CZ21-specific factors (line 6 in each panel), a latent factor model where the factors are K PCs extracted from the CZ21 portfolios only (lines 7-8 in each panel), and a latent factor model where the factors are K IPCs extracted from the CRSP individual stocks only (lines 9-10). Observable factor model specifications are CAPM, FF3, FF3 + Momentum, FF5, and FF5 + Momentum in the K = 1, 3, 4, 5, 6 columns, respectively. The models (IPCA on individual stocks and PCs on CZ21 portfolios, and group-factor model based on the previous two models) are estimated on non-overlapping windows starting in year t - 4 and ending in year t, for each t = 1970, ..., 2020. The R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r : CRSP) or CZ21 portfolios (r : CZ21) as described in Section B.

	In-Sample			Ou	t-of-San	nple				
N. of factors, K	1	3	4	5	6	1	3	4	5	6
Panel A: Total R ²										
r: CRSP, f : FF + mom	21.5	29.4	32.2	33.1	35.8	16.3	15.9	13.1	11.9	9.4
r: CZ21, f : FF + mom	73.9	89.8	91.5	90.7	92.3	70.8	83.6	84.7	83.6	84.8
r: CRSP, f : 3CF		23.5					25.5			
r: CZ21, f: 3CF		93.6					91.2			
r: CRSP, f : 3CF + CRSP spec.			24.1	24.8	25.4			26.4	27.6	28.6
r: CZ21, f : 3CF + CZ21 spec.			94.3	94.9	95.4			92.0	92.4	92.7
r: CRSP, f : PCA on CZ21	24.5	31.8	34.2	36.8	38.8	18.7	18.4	17.9	17.4	16.9
r: CZ21, f : PCA on CZ21	89.3	94.5	95.2	95.7	96.1	87.3	91.9	92.5	92.9	93.1
r: CRSP, f : IPCA on CRSP	5.4	17.8	22.2	24.2	25.1	5.3	18.7	22.1	24.8	27.1
r: CZ21, f : IPCA on CRSP	21.4	65.6	80.3	86.5	88.8	9.0	52.9	60.7	68.9	73.4
Panel B: Pricing \mathbb{R}^2										
r: CRSP, f : FF + mom	32.5	51.0	56.3	54.3	59.1	7.5	8.1	6.5	4.2	1.2
r: CZ21, f: FF + mom	75.2	89.0	88.4	90.0	89.7	78.8	85.5	84.4	86.7	86.4
r: CRSP, f: 3CF		45.1					34.8			
r: CZ21, f: 3CF		93.0					92.9			
r: CRSP, f : 3CF + CRSP spec.			49.3	49.1	49.6			36.4	37.8	36.8
r: CZ21, f : 3CF + CZ21 spec.			94.7	95.2	95.7			93.5	93.9	94.5
r: CRSP, f : PCA on CZ21	43.8	57.4	59.2	61.8	64.8	8.9	6.3	7.0	5.5	4.4
r: CZ21, f : PCA on CZ21	89.4	93.9	95.0	95.9	96.4	90.7	93.2	94.0	94.7	95.0
r: CRSP, f : IPCA on CRSP	17.1	33.9	40.8	45.8	45.8	6.9	33.3	35.6	39.9	42.4
r: CZ21, f : IPCA on CRSP	66.1	94.0	91.6	94.7	95.5	32.0	81.1	85.1	91.9	93.5
Panel C: Predictive R ²										
r: CRSP, f : FF + mom	0.51	0.89	1.03	0.96	1.07	0.35	0.40	0.36	0.17	0.17
r: CZ21, f: FF + mom	2.60	4.00	4.03	4.07	4.12	2.63	3.92	4.09	3.91	4.11
r: CRSP, f: 3CF		0.80					0.79			
r: CZ21, f: 3CF		4.18					4.75			
r: CRSP, f : 3CF + CRSP spec.			0.80	0.86	0.91			0.82	0.84	0.89
r: CZ21, f: 3CF + CZ21 spec.			4.30	4.33	4.36			4.79	4.81	4.83
r: CRSP, f : PCA on CZ21	0.77	1.04	1.08	1.15	1.20	0.47	0.40	0.38	0.29	0.28
r: CZ21, f : PCA on CZ21	4.02	4.30	4.36	4.40	4.43	4.49	4.71	4.73	4.75	4.77
r: CRSP, f : IPCA on CRSP	0.39	0.85	0.95	1.05	1.02	0.36	0.79	0.94	0.97	0.94
r: CZ21, f : IPCA on CRSP	1.90	2.54	3.79	4.15	4.22	1.94	3.75	4.30	4.36	4.14

TABLE 2

Wald Tests

The table reports the values of the Wald test statistics (first row) and associated p-values (second row) for the joint significance of the risk premia of the three CZ21-specific factors estimated by cross-sectional regressions of the risk premia of the CZ21 portfolios on the betas of 3CF, and the three CZ21-specific factors considered in Table 1. Details and data sources of the macro indicators/factors are provided in Supplementary Material Table OA.1. Cross-sectional regressions are estimated by Weighted Least Squares, where each cross-sectional observation is scaled by the square root of the pricing error from the regression itself. P-values are obtained using the asymptotic distribution of the Wald test for the joint significance of the three coefficients, namely a $\chi^2(3)$, with variance covariance matrix of estimated factors' risk premia computed using the Fama-MacBeth procedure and 3-lags Newey-West weighting. Stars are related to p-values as follows: ***p < 0.01, **p < 0.05, *p < 0.10. The model (PCs on CZ21 portfolios, and group-factor model) and cross-sectional WLS regressions (of risk premia of CZ21 quantile portfolios on factors' betas) are estimated on non-overlapping windows of 60 months starting in year t - 4 and ending in year t, for each t = 1970, ..., 2020.

Sample	66-70	71-75	76-80	81-85	86-90	90-95	96-00	01-05	06-10	11-15	16-20
Panel A: 3 CZ21-specific factors, con	ntrols: 3 (CFs									
3 CZ21-spec.	5.012	1.999	2.409	6.246	2.991	2.100	0.297	4.577	2.810	0.784	1.731
Panel B: observable macro factors, c	controls: 3	3 CFs									
Term Spread 10y-1y	0.27	0.37	0.06	1.78	0.12	0.08	0.33	2.06	3.31*	0.12	0.33
Term Spread 10y-3m	0.14	0.41	0.04	1.14	0.15	0.07	2.81*	1.43	1.16	0.03	0.66
Term Spread 1y-FEDfund	0.25	0.06	0.07	1.26	0.00	0.00	0.56	1.10	0.64	0.08	1.37
Default Spread BAA-AAA	0.14	0.25	0.22	0.38	0.63	0.17	0.41	5.41**	1.95	1.22	1.33
Cochrane Piazzesi factor	0.39	0.00	0.05	0.05	0.26	0.67	1.61	5.75**	0.24	0.00	0.69
Consumption growth	0.97	0.62	0.07	0.36	1.03	0.32	0.01				•
Labor Income growth	0.02	0.06	1.51	0.10	0.00	0.05	0.03	0.96	0.00	0.14	0.03
IP growth	0.49	0.28	0.63	0.89	0.02	0.00	0.05	0.11	0.66	0.04	0.03
IP growth innov	0.77	0.76	0.71	1.45	0.01	0.29	0.13	0.41	0.48	0.01	0.00
L&N Macro PCs (1:3) VAR innov	1.18	0.82	3.13	3.20	0.85	0.45	1.72	0.40	1.01	0.41	2.70
M&N Macro FRED MD PCs (1:4)	4.17	1.53	2.38	4.57	5.41	1.43	3.98	19.27***	5.87	1.30	9.03*
Inflation	0.18	0.42	0.01	0.41	0.46	3.33*	0.06	0.21	0.10	0.28	0.00
Lagged Inflation	0.01	0.12	0.02	1.11	1.70	0.35	0.69	2.01	0.69	0.03	0.00
Expected Inflation	0.03	0.16	0.01	2.99*	1.79	0.91	0.33	1.36	1.28	0.19	0.94
Unexpected Inflation	0.29	0.12	0.08	0.59	0.00	0.90	0.11	0.09	0.54	0.72	0.10
Inflation innov	0.13	1.23	0.31	1.97	0.05	0.24	0.25	0.00	0.89	0.33	0.00
EPU					0.03	0.00	2.88*	2.48	0.77	0.04	2.54
EPU growth					0.33	0.13	0.12	0.52	2.08	0.02	0.16
Macro Uncertainty (h=1m)	0.11	0.13	0.14	0.69	0.13	0.12	2.06	0.41	1.60	0.00	2.19
Macro Uncertainty (h=3m)	0.19	0.01	0.29	1.02	0.28	0.13	2.27	0.40	2.00	0.01	1.69
Macro Uncertainty (h=12m)	0.02	0.22	0.29	0.90	0.01	0.07	3.03*	0.55	0.59	0.02	1.65

TABLE 3

Out-of-sample Sharpe ratios of factor portfolios

Table 1 for details on the different factor models. Factors and loadings for each model are first estimated on non-overlapping windows starting in year t-4 and ending in year t, for each t=1970,...,2020. On the same windows, mean-variance efficient (i.e. maximum Sharpe ratio) portfolio weights are computed from the first K factors. Then, the loadings estimated in the previous block are used to reconstruct the factor returns on the next 5-years block, jointly with the returns of that block, and finally the OOS returns are combined together using mean-variance efficient portfolio weights form the previous block, to form the OOS maximum Sharpe ratio portfolio returns, from which the reported ex-post Sharpe Ratios are computed.

N. of factors, K	1	3	4	5	6
FF + mom	0.39	0.30	0.62	0.63	0.81
3CF		0.66			
3CF + CZ21 spec.			0.79	0.89	0.93
3CF + CRSP spec.			0.24	0.37	0.61
PCA on CZ21	0.49	0.54	0.81	0.98	0.84
IPCA on CRSP	0.31	0.77	1.09	1.10	1.44

TABLE 4

3CF and the Zoo

The table documents which factors from the zoo are the most related (a) to 3CF, and (b) to 3CF augmented with the first three group-specific factors of the CZ21 portfolios. Panel A reports the twenty factors with the largest average R^2 s - across all our 5 non-overlapping windows - when regressed on the 3CF factors. Panel B reports the factors in the zoo showing the highest increase in the R^2 - averaged across all non-overlapping 5 years windows - when the first 3 CZ21-specific factors are added as regressors to 3CF.

Panel A: 3 Common factors only		Panel B: 3 CZ21-specif. added to 3 common	factors
Factor	R^2	Factor	ΔR^2
CAPM beta (1973)	90.9	Momentum and LT Reversal (2006)	31.7
Price (1972)	87.7	Junk Stock Momentum (2007)	28.6
Bid-ask spread (1986)	87.6	Momentum in high volume stocks (2000)	27.7
Frazzini-Pedersen Beta (2014)	86.5	Option to stock volume (2012)	26.1
Cash-flow to price variance (1996)	85.4	gross profits / total assets (2013)	25.1
Volume to market equity (1996)	84.6	Momentum based on FF3 residuals (2011)	24.4
Past trading volume (1998)	81.2	Composite equity issuance (2006)	24.3
52 week high (2004)	79.4	Conglomerate return (2012)	23.8
Size (1981)	79.2	Intangible return using EP (2006)	23.7
EPS Forecast Dispersion (2002)	77.3	Firm Age - Momentum (2004)	23.6
Book to market, most recent ME (1985)	77.1	Inst Own and Forecast Dispersion (2005)	23.1
Idiosyncratic risk (AHT) (2003)	76.7	Taxable income to income (2004)	23.0
Days with zero trades (2006)	76.6	Momentum without the seasonal part (2008)	22.1
Days with zero trades (2006)	75.4	Volatility smirk near the money (2010)	21.7
Price delay R-squared (2005)	73.7	O Score (1998)	21.4
Intangible return using Sale2P (2006)	72.8	Total assets to market (1992)	21.1
Days with zero trades (2006)	72.5	Market leverage (1988)	21.0
Idiosyncratic risk (3 factor) (2006)	72.1	Industry Momentum (1999)	20.6
Idiosyncratic risk (2006)	71.5	Earnings Surprise (1984)	20.5
Long-term EPS forecast (1996)	71.0	Breadth of ownership (2002)	20.4

 ${\it TABLE~5}$ ${\it Total, Pricing~and~Predictive~} {\it R}^2 {\it s} \mbox{ - Common factors versus zoo~PCs}$

The table reports Total, Pricing and Predictive R^2 s in percent for a latent factor model with only 3 common factors (lines 1-2 in each of the three panels) - a repeat of lines 3-4 in Panels A - C of Table 1, and a latent factor model where the factors are K PCs extracted from the factors in the zoo only (lines 3-4). The models are estimated on the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020. R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

	In-Sample			Out-of-Sample						
N. of factors, K	1	3	4	5	6	1	3	4	5	6
Panel A: Total R ²										
r: CRSP, f: 3CF		23.5					25.5			
r: CZ21, f: 3CF		93.6					91.2			
r: CRSP, f : PCA on Zoo	15.8	26.1	29.5	32.3	35.1	6.6	6.6	7.4	6.9	8.3
r: CZ21, f : PCA on Zoo	42.2	68.6	72.6	75.6	77.7	29.4	49.0	53.5	56.0	59.0
Panel B: Pricing \mathbb{R}^2										
r: CRSP, f : 3CF		45.1					34.8			
r: CZ21, f: 3CF		93.0					92.9			
r: CRSP, f : PCA on Zoo	<0	<0	10.4	21.7	22.4	5.3	<0	<0	<0	<0
r: CZ21, f: PCA on Zoo	<0	<0	<0	9.2	7.7	<0	<0	<0	<0	<0
Panel C: Predictive R^2										
r: CRSP, f: 3CF		0.80					1.13			
r: CZ21, f: 3CF		4.18					4.77			
r: CRSP, f : PCA on Zoo	<0	<0	0.14	0.39	0.41	<0	<0	<0	<0	<0
r: CZ21, f : PCA on Zoo	<0	<0	0.07	0.85	1.04	0.11	<0	<0	0.05	0.25

 $\label{eq:TABLE 6} \mbox{Maximum Canonical Correlations Observed Zoo Factors and 3CF}$

The table reports the highest canonical correlation among any two combined zoo factors with the market with 3CF. The analysis is repeated for each of the 11 five-years non-overlapping blocks. The numbers in the last column reflect the maximal sum of correlations attained by each of the selected combinations.

1966-1970	Gross profits/total assets + Size	2.61
1971-1975	Earnings-to-Price Ratio + Share issuance (5 year)	2.54
1976-1980	Book leverage (annual) + Size	2.46
1981-1985	Operating leverage + Size	2.68
1986-1990	Analyst earnings per share + Size	2.51
1991-1995	Accruals + Past trading volume	2.32
1996-2000	Change in PPE and inv/assets + O Score	2.49
2001-2005	Industry concentration (sales) + Short term reversal	2.50
2006-2010	BTM using most recent ME + Intangible return using EP	2.58
2011-2015	Off season reversal years 6 to 10 + Size	2.49
2016-2020	Predicted div yield next month + Net payout yield	2.46

TABLE 7

Maximum Canonical Correlations between 3CF and observed factors

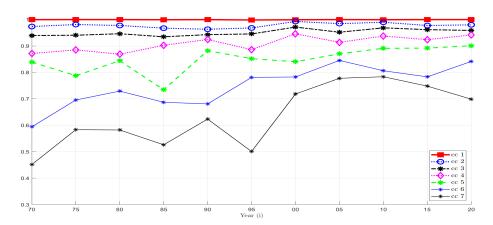
The table reports the highest canonical correlation among any three combined macro factors (Panel A) or two macro factors and market (Panel B) with 3CF. The analysis is repeated for each of the 11 five-years non-overlapping blocks. The numbers in the last column reflect the maximal sum of correlations attained by each of the selected combinations.

Panel A. Observed factors: macro 1966-1970 Term Spread 1y-FEDfund + Default Spread BAA-AAA + Expected Inflation 1.22 1971-1975 Default Spread BAA-AAA +L&N Macro PC2 VAR innov + Unexpected Inflation 1.07 L&N Macro PC1 VAR innov + Unexpected Inflation + M&N Macro FRED MD PC2 1976-1980 1.16 L&N Macro PC1 VAR innov + M&N Macro FRED MD PC1 + M&N Macro FRED MD PC2 1981-1985 1.16 1986-1990 L&N Macro PC1 VAR innov + EPU growth + M&N Macro FRED MD PC1 1.42 IP growth + Macro PC var innov 1 + Expected Inflation 1991-1995 1.30 1996-2000 Consumption growth + Macro Uncertainty (h=1m) + M&N Macro FRED MD PC2 1.05 2001-2005 Term Spread 1y-FEDfund + EPU + EPU growth 1.17 L&N Macro PC2 VAR innov + L&N Macro PC3 VAR innov + CP 2006-2010 1.10 2011-2015 IP growth innov + L&N Macro PC3 VAR innov + EPU growth 0.81 2016-2020 IP growth innov + L&N Macro PC2 VAR innov + Macro Uncertainty (h=1m)1.75 Panel B. Observed factors: market and two macro 1966-1970 Default Spread BAA-AAA + L&N Macro PC1 VAR innov 1.67 1971-1975 Unexpected Inflation + M&N Macro FRED MD PC3 1.59 1976-1980 L&N Macro PC1 VAR innov + M&N Macro FRED MD PC2 1.64 1981-1985 Macro Uncertainty (h=1m) + M&N Macro FRED MD PC1 1.68 1986-1990 L&N Macro PC1 VAR innov + M&N Macro FRED MD PC1 1.85 1991-1995 L&N Macro PC1 VAR innov + EPU growth 1.69 1996-2000 Consumption growth + Macro Uncertainty (h=1m)1.54 2001-2005 Term Spread 1y-FEDfund + EPU growth 1.70 2006-2010 Default Spread BAA-AAA + L&N Macro PC2 VAR innov 1.65 2011-2015 Labor Income growth + IP growth innov 1.42 2016-2020 Labor Income growth + Macro Uncertainty (h=1m) 2.21

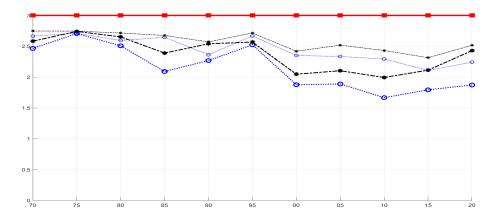
FIGURE 1

Canonical correlations

Panel (a) displays the first seven canonical correlations between 10 IPCA factors (implemented using individual stocks from CRSP and the characteristics used in Freyberger et al. (2020), updated by Jensen et al. (2023)) and 10 PCs from Chen and Zimmermann (2022) quantile portfolios. Panel (b) displays the sum of the canonical correlations of the three factors common across CRSP and CZ21 test assets with the 3 Fama and French factors (FF3): Market, SMB and HML factors (blue circles, thick dotted line), and the sum of the canonical correlations of the three common factors with the 5 Fama and French factors (FF5) - adding RMW-operating profitability, CMA-investment style (black stars, thick dashed line). In addition, it also displays FF3/FF5 augmented with the momentum factor (FF3+Mom - blue circles, thin dotted line, FF5+Mom - black stars, thin dashed line). The red line across the plot marks the 3-factor benchmark common factor space. For each year the 10 IPCs from individual stocks, the 10 PCs from the quintile portfolios, the three common factors among them and the different sets of canonical correlations in both panels are computed on the 5-years block of monthly data starting in year t - 4, for each t = 1970, 1975, ..., 2020.



(a) Canonical correlations between 10 IPCs on individual stocks and 10 PCs sorted portfolios

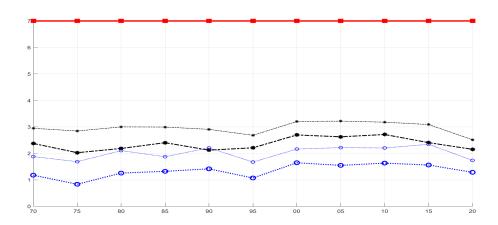


(b) Sum of canonical correlations between 3 common factors (3CF) and observable factors

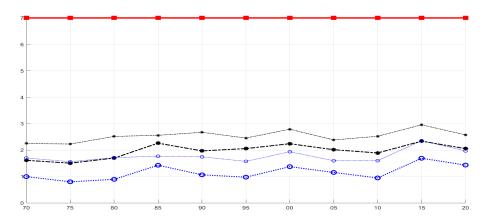
FIGURE 2

Sum of canonical correlations of the group-specific factors

Panel (a) displays the sum of the canonical correlations of the seven specific factors in CZ21 test asset portfolios with: (i) FF3 (blue circles, thick dotted line): (ii) FF5 (black stars, thick dashed line); (iii) the FF3 factors and momentum (blue circles, thin dotted line); (iv) the FF5 factors and momentum (black stars, thin dashed line). Panel (b) displays the sum of the canonical correlations of the seven specific factors in CRSP individual stocks with the same four sets of observable factors. For each year we report results computed on the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020.



CZ21 test assets

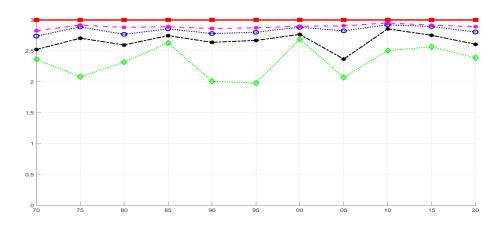


CRSP individual stocks

FIGURE 3

Sum of canonical correlations between 3CF and PCs from the zoo.

The figure displays the sum of the canonical correlations of 3CF factors with the first 15 PCs (magenta circles, dashed line), 10 PCs (blue circles, dotted thin line), 5 PCs (black stars, dashed thin line), and finally the first 3 PCs from the factor zoo (green diamonds, dotted thin line). For each year we report results computed on the block starting in year t-4, for each t=1970,1975,...,2020.



Appendix

Section A.1 covers the proofs of the Propositions appearing in Section II. Section A.2 extends Theorems 1 and 2 in AGGR to the case in which the latent factors are estimated by the variations of PCA described in *Supplementary Material* Section OA.3, namely IPCA for the first group of individual stock returns, and RP_PCA for the second group of portfolio returns, which are not covered in the original work of AGGR, who only considered classical PCA on demeaned data. The estimators of the canonical correlations among the factors from the two groups and the test statistics in Theorems 1 and 2 of AGGR need to be adjusted accordingly. The technical assumptions of the theorems and a sketch of their proofs is provided in the *Supplementary Material*. We will denote the sample mean of a generic sequence z_{τ} , $\tau = 1, ..., T$ as $\bar{z} = \frac{1}{T} \sum_{\tau=1}^{T} z_{\tau}$, the T-dimensional vector of ones as 1 = [1, ..., 1]', and the identity matrix of order T as I_T . Finally, the upper index (c) denotes the upper $(k^c, 1)$ block of a vector, and the upper index (c, c) denotes the upper-left (k^c, k^c) block of a matrix.

A.1. Proof of Propositions in Section II

A. Proof of Proposition 1

The proof of the proposition follows from writing the time-varying beta model in terms of the assumed beta dynamics appearing in equation (4) using the decomposition in Assumption 1.

B. Proof of Proposition 2

Under Assumptions 2 and 3 the portfolios excess returns are such that:

(A.1)
$$y_{j,\tau}^p = [B_j^0 + B_j Z_{\tau-1}^*]' f_{\tau}^c + e_{j,\tau}^p,$$

where $B_j^0 := \mathscr{C}W_j^0$ and $B_j := \mathscr{C}W_j$. By rearranging terms in this scaled factor model, using $(B_j Z_{\tau-1}^*)' f_\tau^c = (f_\tau^c)' B_j Z_{\tau-1}^* = \text{vec}((f_\tau^c)' B_j Z_{\tau-1}^*) = (Z_{\tau-1}^* \otimes f_\tau^c)' \text{ vec}(B_j)$, and by plugging this equation into (7) we get

(A.2)
$$y_{i,\tau}^p = \lambda_i' f_{2,\tau} + e_{i,\tau}^p$$

where $f_{2,\tau}=(f^{c\prime}_{\tau},f^{s\prime}_{2,\tau})'$ with $f^s_{2,\tau}=Z^*_{\tau-1}\otimes f^c_{\tau}$. By the projection argument we get that the factor spaces spanned by f^c_{τ} and $f^s_{2,\tau}$ do not intersect. Q.E.D.

A.2. Group Factor Model: Extension for Time-varying

Loadings and Factors with Generic Mean

A more general version of the group-factor model with constant loadings for both groups considered in AGGR is one where we allow for time varying loadings for the returns of the assets in group 1, while still allowing for constant loadings in group 2. In particular, let us consider the special version of model (6) for group 1 that we estimate in our empirical analysis for individual stocks (our group 1), namely

(A.3)
$$y_{1,i,\tau} = (CZ_{i,\tau-1})' f_{1,\tau} + \varepsilon_{1,i,\tau} = \lambda'_{1,i,\tau} h_{1,\tau} + \varepsilon_{1,i,\tau}$$

where $h_{1,\tau}=f_{1,\tau}=(f_{\tau}^c,f_{1,\tau}^{s\prime})'$, and with $\lambda_{1,i,\tau}=(\lambda_{1,i,\tau}^{c\prime},\lambda_{1,i,\tau}^{s\prime})'=CZ_{i,\tau-1}$. The L-dimensional vector $Z_{i,\tau-1}=[Z_{i,\tau-1,1},...,Z_{i,\tau-1,L}]'$ collects the stock-specific characteristics observable at date t-1, and $C=[C^{c\prime},\ C^{s\prime}]'$ is a full-column $k_1\times L$ matrix, with $k_1\leq L\leq N_1$, and usually $k_1< L\ll N_1$. Moreover, C^c (resp. C^s) is also a full-column $k^c\times L$ (resp. $k_1^s\times L$) matrix, implying $\lambda_{1,i,\tau}^c=C^cZ_{i,\tau-1}$ (resp. $\lambda_{1,i,\tau}^s=C^sZ_{i,\tau-1}$). To keep the appendix self-contained, we also re-write the model (8) that we estimate in our empirical analysis for portfolios (our group 2), as

(A.4)
$$y_{2,i,\tau} = \lambda'_{2,i} h_{2,\tau} + \varepsilon_{2,i,\tau}$$

where $h_{2,\tau}=(f_{\tau}^c,f_{2,\tau}^{s\prime})'$, and with $\lambda_{2,i}=(\lambda_{2,i}^{c\prime},\lambda_{2,i}^{s\prime})'$. Then, the group-factor model appearing with constant loadings in AGGR can be generalized as

(A.5)
$$\begin{bmatrix} y_{1,\tau} \\ y_{2,\tau} \end{bmatrix} = \begin{bmatrix} \Lambda_{1,\tau}^c & \Lambda_{1,\tau}^s & 0 \\ \Lambda_2^c & 0 & \Lambda_2^s \end{bmatrix} \begin{bmatrix} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{bmatrix} + \begin{bmatrix} \varepsilon_{1,\tau} \\ \varepsilon_{2,\tau} \end{bmatrix},$$

for a generic sample of dates $\tau=1,...,T$, where $\Lambda_{1,\tau}^c=[\lambda_{1,1,\tau}^c,...,\lambda_{1,N_1,\tau}^c]'=[Z_{1,\tau-1},...,Z_{N_1,\tau-1}]'C^{c\prime}$ and $\Lambda_{1,\tau}^s=[\lambda_{1,1,\tau}^s,...,\lambda_{1,N_1,\tau}^s]'=[Z_{1,\tau-1},...,Z_{N_1,\tau-1}]'C^{s\prime}$ are the time-varying matrices of factor loadings for group 1. Moreover, $\Lambda_2^c=[\lambda_{2,1}^c,...,\lambda_{2,N_2}^c]'$ and $\Lambda_2^s=[\lambda_{2,1}^s,...,\lambda_{2,N_2}^s]'$ are the matrices of constant factor loadings for group 2, and $\varepsilon_{j,\tau}=[\varepsilon_{j,1\tau},...,\varepsilon_{j,N_j\tau}]'$ and $\varepsilon_{j,\tau}=[\varepsilon_{j,1,\tau},...,\varepsilon_{j,N_j,\tau}]'$ the error terms, with j=1,2.

As in AGGR we assume, without loss of generality, that the group-specific factors $f_{1,\tau}^s$ and $f_{2,\tau}^s$ are unconditionally orthogonal to the common factor f_{τ}^c . Since the unobservable factors can be standardized, we have:

$$\mathbb{E}\left[\begin{array}{c} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{array}\right] = \left[\begin{array}{c} \mu^c \\ \mu_1^s \\ \mu_2^s \end{array}\right], \quad \text{and} \quad \Sigma_F := \mathbb{V}\left[\begin{array}{c} f_{\tau}^c \\ f_{1,\tau}^s \\ f_{2,\tau}^s \end{array}\right] = \left[\begin{array}{c} I_{k^c} & 0 & 0 \\ 0 & I_{k_1^s} & \Upsilon \\ 0 & \Upsilon' & I_{k_2^s} \end{array}\right],$$

where the expected values of the factors are finite, and matrix Σ_F is positive-definite. We allow for a non-zero covariance Υ between group-specific factors, but differently from AGGR, we allow the factors to have expected value different from zero.

Model (A.5) together with Assumption (A.6) is identified by the same arguments based on canonical correlation analysis used by AGGR, which we summarize here as our identification argument is constructive for the estimation. Let $h_{j,\tau} = [f_{\tau}^{c\prime}, f_{j,\tau}^{s\prime}]'$, with j=1,2, and $V_{j\ell} = \operatorname{Cov}(h_{j,\tau},h_{\ell,\tau})$, with $j,\ell=1,2$. The $\underline{k} = \min(k_1,k_2)$ largest eigenvalues of the matrices $R = V_{11}^{-1}V_{12}V_{22}^{-1}V_{21}$ and $R^* = V_{22}^{-1}V_{21}V_{11}^{-1}V_{12}$ are the same, and are equal to the squared canonical correlations ρ_{ℓ}^2 , $\ell=1,...,\underline{k}$, between $h_{1,\tau}$ and $h_{2,\tau}$. The associated eigenvectors $w_{1,\ell}$ (resp. $w_{2,\ell}$), with $\ell=1,...,\underline{k}$, of matrix R (resp. R^*) standardized such that $w_{1,\ell}'V_{11}w_{1,\ell}=1$ (resp. $w_{2,\ell}'V_{22}w_{2,\ell}=1$) are the canonical directions which yield the canonical variables $w_{1,\ell}'h_{1,\tau}$ (resp. $w_{2,\ell}'h_{2,\tau}$). The next Proposition A.1 deals

with determining k^c , the number of common factors, using canonical correlations between the vectors $h_{1,\tau}$ and $h_{2,\tau}$, which are unobserved and can be estimated by performing IPCA as in Kelly et al. (2019) in group 1, and PCA or RP_PCA as in Lettau and Pelger (2020b) in group 2. It corresponds to Proposition 1 in AGGR where the zero mean assumption of the factors is replaced with our new Assumption (A.6).

PROPOSITION A.1 Under Assumption (A.6), the following hold: (i) If $k^c > 0$, the largest k^c canonical correlations between $h_{1,\tau}$ and $h_{2,\tau}$ are equal to 1, and the remaining $\underline{k} - k^c$ canonical correlations are strictly less than 1, (ii) Let W_j be the (k_j, k^c) matrix whose columns are the canonical directions for $h_{j,\tau}$ associated with the k^c canonical correlations equal to 1, for j=1,2. Then, $f_{\tau}^c = W_j' h_{j,\tau}$ (up to an orthogonal matrix), (iii) If $k^c = 0$, all canonical correlations between $h_{1,\tau}$ and $h_{2,\tau}$ are strictly less than 1.

(iv) Let W_1^s (resp. W_2^s) be the (k_1, k_1^s) (resp. (k_2, k_2^s)) matrix whose columns are the eigenvectors of matrix R (resp. R^*) associated with the smallest k_1^s (resp. k_2^s) eigenvalues. Then $f_{j,\tau}^s = W_j^{s\prime} h_{j,\tau}$ (up to an orthogonal matrix) for j=1,2.

Proposition A.1 shows that the number of common factors k^c , the common factor space spanned by f^c_{τ} , and the spaces spanned by group-specific factors, can be identified from the canonical correlations and canonical variables of $h_{1,\tau}$ and $h_{2,\tau}$. Therefore, the factor space dimensions k^c , k^s_j , and factors f^c_{τ} and $f^s_{j,\tau}$, j=1,2, are identifiable (up to a rotation) from information that can be inferred by disjoint IPCA and PCA (or RP_PCA) described in Section OA.3, on the two groups.

A. Estimation of Factors and Loadings

When the true number of factors $k_j>0$ in each subgroup j=1,2 and $k^c>0$ are known, Proposition A.1 suggests the following estimation procedure. Let $\hat{h}_{1,\tau}^{ipc}$ be the IPCA estimator of the k_1 factor values in the vector $h_{1,\tau}$, and \hat{C} be the IPCA estimator of matrix C mapping stock-specific characteristics to factor loadings, and $\hat{\Lambda}_{1,\tau}=\hat{C}'Z_{\tau-1}$ be the IPCA estimator of the $N_1\times k_1$ matrix of factor loadings $\Lambda_{1,\tau}$ at date τ , for the generic dates $\tau=1,...,T$, of model (A.3), as defined in Kelly et al. (2019), and summarized in Section OA.3.2 in the Supplementary Material. Let also $\hat{h}_{2,\tau}$ be the PCA, or RP_PCA, estimate of the k_2 factors $h_{2,\tau}$, and $\hat{\Lambda}_2$ be the PCA, or RP_PCA, estimator of the loading matrix Λ_2 , as summarized in Section OA.3.1.

Let $\hat{V}_{j\ell}$ denote the empirical covariance matrix between $\hat{h}_{j,\tau}$ and $\hat{h}_{\ell,\tau}$, i.e. $\hat{V}_{j\ell} = \sum_{\tau=1}^T \hat{h}_{j,\tau} \hat{h}'_{\ell,\tau} / T - \left(\sum_{\tau=1}^T \hat{h}_{j,\tau} / T\right) \left(\sum_{\tau=1}^T \hat{h}_{\ell,\tau} / T\right)'$, for $j,\ell=1,2$, and let:

$$\hat{R} := \hat{V}_{11}^{-1} \hat{V}_{12} \hat{V}_{22}^{-1} \hat{V}_{21}, \text{ and } \hat{R}^* := \hat{V}_{22}^{-1} \hat{V}_{21} \hat{V}_{11}^{-1} \hat{V}_{12},$$

be the estimators of matrices R and R^* , respectively. Differently from AGGR, the estimators of the variance-covariance matrices $\hat{V}_{j\ell}$ take into account that the estimated factors might have non-zero mean. Matrices \hat{R} and \hat{R}^* have the same non-zero eigenvalues. The k^c largest eigenvalues of \hat{R} (resp. \hat{R}^*), denoted by $\hat{\rho}_{\ell}^2$, $\ell=1,...,k^c$, are the first k^c squared sample canonical correlations between $\hat{h}_{1,\tau}$ and $\hat{h}_{2,\tau}$. The associated k^c canonical directions, collected in the (k_1,k^c) matrix \hat{W}_1 (resp. (k_2,k^c) matrix \hat{W}_2), are the eigenvectors associated with the k^c largest eigenvalues of matrix \hat{R} (resp. \hat{R}^*), normalized to have length 1 with respect to \hat{V}_{11} (resp. \hat{V}_{22}). It also holds that:

(A.8)
$$\hat{W}_1'\hat{V}_{11}\hat{W}_1 = I_{k^c}$$
, and $\hat{W}_2'\hat{V}_{22}\hat{W}_2 = I_{k^c}$.

DEFINITION 1 Two estimators of the common factors vector are $\hat{f}_{\tau}^c = \hat{W}_1' \hat{h}_{1,\tau}$ and $\hat{f}_{\tau}^{c*} = \hat{W}_2' \hat{h}_{2,\tau}$.

Definition 1 and equation (A.8) imply that the estimated common factors have identity sample variance-covariance matrix: $\hat{V}(\hat{f}_{\tau}^c) := \sum_{\tau=1}^T \hat{f}_{\tau}^c \hat{f}_{\tau}^{c\prime} / T - \left(\sum_{\tau=1}^T \hat{f}_{\tau}^c / T\right) \left(\sum_{\tau=1}^T \hat{f}_{\tau}^{c\prime} / T\right) = I_{k^c}$, and analogously $\hat{V}(\hat{f}_{\tau}^{c*}) = \sum_{\tau=1}^T \hat{f}_{\tau}^{c*} \hat{f}_{\tau}^{c*\prime} / T - \left(\sum_{\tau=1}^T \hat{f}_{\tau}^{c*\prime} / T\right) \left(\sum_{\tau=1}^T \hat{f}_{\tau}^{c*\prime} / T\right) = I_{k^c}$, i.e. the estimated common factor values match in-sample the normalization condition of identity variance-covariance matrix in (A.6).

Let matrix \hat{W}_1^s (resp. \hat{W}_2^s) be the (k_1, k_1^s) (resp. (k_2, k_2^s)) matrix collecting k_1^s (resp. k_2^s) eigenvectors associated with the k_1^s (resp. k_2^s) smallest eigenvalues of matrix \hat{R} (resp. \hat{R}^*), normalized to have length 1 with respect to the matrix \hat{V}_{11} (resp. \hat{V}_{22}). It also holds: $\hat{W}_j^s \hat{V}_{jj} \hat{W}_j^s = I_{k_j^s}$, j = 1, 2. The estimators of the group-specific factors can be defined analogously to the estimators of the common factors

DEFINITION 2 Two estimators of the group-specific factors are $\hat{f}_{1,\tau}^s = \hat{W}_1^{s\;\prime} \hat{h}_{1,\tau}$ and $\hat{f}_{2,\tau}^s = \hat{W}_2^{s\;\prime} \hat{h}_{2,\tau}$.

The fact that, by definition, canonical directions satisfy $\hat{W}^{s}_{j}\hat{V}_{jj}\hat{W}^{c}_{j}=0_{(k^{s}_{j},k^{c})}$ for j=1,2, see ch. 17 in Magnus and Neudecker (2007), implies that \hat{f}^{c}_{τ} and $\hat{f}^{s}_{1,\tau}$ (resp. \hat{f}^{c*}_{τ} and $\hat{f}^{s}_{2,\tau}$) are orthogonal in-sample, i.e. their sample covariance is zero.

The estimators of the loadings can be obtained by noting that $\tilde{W}_j := [\hat{W}_j \ , \ \hat{W}_j^s]'$ is a full-rank $k_j \times k_j$ matrix of (transposed) eigenvectors such that: $\tilde{W}_1 h_{1,\tau} = [\hat{f}_{\tau}^{c\prime}, \ \hat{f}_{1,\tau}^{s\prime}]', \tilde{W}_2 h_{1,\tau} = [\hat{f}_{\tau}^{c*\prime}, \ \hat{f}_{2,\tau}^{s\prime}]', \tilde{W}_j \hat{V}_{jj} \tilde{W}_j' = I_{k_j}$ for j = 1,2, and that following equations holds $y_{1,\tau} = \Lambda_{1,\tau} h_{1,\tau} + \varepsilon_{1,\tau} = \Lambda_{1,\tau} \tilde{W}_1^{-1} \tilde{W}_1 h_{1,\tau} + \varepsilon_{1,\tau}$ and $y_{2,\tau} = \Lambda_2 h_{2,\tau} + \varepsilon_{2,\tau} = \Lambda_2 \tilde{W}_2^{-1} \tilde{W}_2 h_{2,\tau} + \varepsilon_{2,\tau}.^{20}$ Let $(\tilde{W}_j^{-1})^c$ (resp. $(\tilde{W}_j^{-1})^s$) the $k_j \times k_c$ (resp. $k_j \times k_j^s$) matrix collecting the first k_c (resp. last k_j^s) columns of \tilde{W}_j^{-1} , i.e. we partition the inverse of \tilde{W}_j as $\tilde{W}_j^{-1} = [(\tilde{W}_j^{-1})^c, \ (\tilde{W}_j^{-1})^s]$, then the loadings of the common factors and group-specific can be defined as follows. **DEFINITION 3** The estimators of the common factors loadings are $\hat{\Lambda}_{1,\tau}^c = \hat{\Lambda}_{1,\tau} (\tilde{W}_1^{-1})^c = \hat{C}' Z_{\tau-1} (\tilde{W}_1^{-1})^c$ and

DEFINITION 3 The estimators of the common factors loadings are $\hat{\Lambda}_{1,\tau}^c = \hat{\Lambda}_{1,\tau}(\tilde{W}_1^{-1})^c = \hat{C}'Z_{\tau-1}(\tilde{W}_1^{-1})^c$ and $\hat{\Lambda}_2^c = \hat{\Lambda}_2(\tilde{W}_2^{-1})^c$. Moreover, the estimators of the common factors loadings are $\hat{\Lambda}_{1,\tau}^s = \hat{\Lambda}_{1,\tau}(\tilde{W}_1^{-1})^s = \hat{C}'Z_{\tau-1}(\tilde{W}_1^{-1})^s \text{ and } \hat{\Lambda}_2^s = \hat{\Lambda}_2(\tilde{W}_2^{-1})^s.$

B. Inference on the Number of Common Factors

In order to infer the dimension k^c of the common factor space, we consider the case where the number of pervasive factors k_1 and k_2 in each sub-panel is known, hence $\underline{k} = \min(k_1, k_2)$ is also known. As explained in AGGR, all the results remain unchanged when the numbers of pervasive factors k_1 and k_2 are estimated consistently. From Proposition A.1, dimension k^c is the number of unit canonical correlations between $h_{1,\tau}$ and $h_{2,\tau}$. We consider the hypotheses: $H(0) = \{1 > \rho_1 \ge \ldots \ge \rho_{\underline{k}}\}$, $H(1) = \{\rho_1 = 1 > \rho_2 \ge \ldots \ge \rho_{\underline{k}}\}$, ..., $H(k^c) = \{\rho_1 = \ldots = \rho_{k^c} = 1 > \rho_{k^c+1} \ge \ldots \ge \rho_{\underline{k}}\}$, ..., and finally $H(\underline{k}) = \{\rho_1 = \ldots = \rho_{\underline{k}} = 1\}$, where $\rho_1, \ldots, \rho_{\underline{k}}$ are the ordered canonical correlations of $h_{1,\tau}$ and $h_{2,\tau}$. Generically, $H(k^c)$ corresponds to the case of k^c common factors and $k_1 - k^c$ and $k_2 - k^c$ group-specific factors in each group, and H(0) corresponds to the absence of common factors. In order to select the number of common factors, let us consider the following sequence of tests: $H_0 = H(k^c)$ against $H_1 = \bigcup_{0 \le r < k^c} H(r)$, for each $k^c = \underline{k}, \underline{k} - 1, \ldots, 1$. To test H_0 against H_1 , for any given $k^c = \underline{k}, \underline{k} - 1, \ldots, 1$ we consider the test statics $\hat{\xi}(k^c)$ defined in equation (11). The null hypothesis $H_0 = H(k^c)$ is rejected when $\hat{\xi}(k^c) - k^c$ is negative and large. The critical value is obtained from the large sample distribution of the statistic when $\frac{1}{2^0}$ Let \hat{W}_1 and \hat{W}_2 be the "usual" orthonormal eigenvectors of \hat{R} and \hat{R}^* , such that $\hat{W}_1' \hat{W}_1 = I_{k_1}$ and

 $W_1' = I_{k_1}$ and $W_2' = I_{k_2}$, and easily obtainable from most statistical software. Then, \tilde{W}_j can be easily computed from $\tilde{W}_j' := [\hat{W}_j, \hat{W}_j^s]' = \hat{V}_{jj}^{-1/2} \tilde{W}_j$, see e.g. the proof of Theorem 13 of ch. 17 Magnus and Neudecker (2007).

 $N_1, N_2, T \to \infty$, provided below. The number of common factors is estimated by sequentially applying the tests starting from $k^c = \underline{k}$, the maximum number of common factors.

Let us denote $N=\min\{N_1,N_2\}$ and $\mu_N=\sqrt{N_2/N_1}$. Without loss of generality, we set $N=N_2$, which implies $\mu_N\leq 1$. We assume that:

(A.9)
$$\sqrt{T}/N = o(1), \ N/T^2 = o(1) \text{ and } \mu_N \to \mu, \text{ with } \mu \in [0, 1].$$

Notably, the assumption $N/T^2 = o(1)$ is made also by Lettau and Pelger (2020a), and is more restrictive than the assumption $N/T^{5/2} = o(1)$ made by AGGR in their equation (4.1), and simplifies the derivation and the form of the asymptotic distribution of our test statistics. To further simplify the analysis, and similarly to AGGR, we assume that the errors $\varepsilon_{j,i,\tau}$ are a conditionally homoscedastic martingale difference sequence for each individual i, conditional on the sigma field $\mathcal{F}_{\tau}^{ipc} = \{F_{\tau}, F_{\tau-1}, ...; Z_{\tau-1}^{1:N}, Z_{\tau-2}^{1:N}, ...\}$ generated by current and past factor values $F_{\tau} = (f_{\tau}^{c\prime}, f_{1,\tau}^{s\prime}, f_{2,\tau}^{s\prime})'$ and past characteristic values $Z_{\tau-1}^{1:N} = [vec(Z_{1,\tau-1})', ..., vec(Z_{N,\tau-1})']'$, that is,

$$(A.10) E[\varepsilon_{j,i,\tau}|\{\varepsilon_{j,i,\tau-h}\}_{h\geq 1},\mathcal{F}_{\tau}^{ipc}] = 0, E[\varepsilon_{j,i,\tau}^{2}|\{\varepsilon_{j,i,\tau-h}\}_{h\geq 1},\mathcal{F}_{\tau}^{ipc}] = \gamma_{j,ii} > 0 (\text{say}),$$

for all j, i, τ, h , see Assumptions A.7 (b) and (c) in the *Supplementary Material*. We further assume that the errors $\varepsilon_{j,i,\tau}$ are conditionally uncorrelated at all leads and lags, both within and across groups, conditionally on \mathcal{F}_{τ}^{ipc} , that is,

(A.11)
$$Cov(\varepsilon_{j,i,\tau}, \varepsilon_{k,\ell,\tau-h}|\mathcal{F}_{\tau}^{ipc}) = 0, \quad \text{when } h \neq 0, \text{ and for all } j, k, i, \tau.$$

Nevertheless we allow for some contemporaneous correlation between the errors within and across panels, i.e. we allow the conditional variance-covariance matrix $V(\varepsilon_{\tau}|\mathcal{F}_{\tau}^{ipc})$ of the stacked errors from the two panels $\varepsilon_{\tau} = [\varepsilon'_{1,\tau}, \varepsilon'_{2,\tau}]'$, to be sparse and time-invariant, namely

(A.12)
$$Cov(\varepsilon_{j,i,\tau}, \varepsilon_{k,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}) = \gamma_{jk,i\ell}, \quad \text{for all } \tau, i = 1, ..., N_j, k, \ell = 1, ..., N_k,$$

and j,k=1,2. Recall that $\hat{h}^{ipc}_{1,\tau}$ is the IPCA estimator of factors $h_{1,\tau}$, while $\hat{\lambda}_{1,i,\tau}=(\hat{\lambda}^{c\prime}_{1,i,\tau},\hat{\lambda}^{s\prime}_{1,i,\tau})'$ is the IPCA

estimator of loadings $\lambda_{1,i,\tau}$, for all $\tau=1,...,T$ and $i=1,...,N_1$, of model (A.3), as defined in Kelly et al. (2019), and summarized in *Supplementary Material* Section OA.3.2.²¹ Then, Kelly et al. (2020) in their Theorem 4 show that, for $\tau=1,\ldots,T$ the estimator $\hat{h}_{1,\tau}^{ipc}$ is asymptotically equivalent, up to negligible terms, to $\hat{\mathcal{H}}_{1}^{ipc}(h_{1,\tau}+u_{1,\tau}^{ipc}/\sqrt{N_1})$ where

(A.13)
$$u_{1,\tau}^{ipc} = \left(\frac{1}{N_1} \sum_{i=1}^{N_1} \lambda_{1,i,\tau} \lambda'_{1,i,\tau}\right)^{-1} \frac{1}{\sqrt{N_1}} \sum_{i=1}^{N_1} \lambda_{1,i,\tau} \varepsilon_{1,i,\tau},$$

as shown in equation (111) of their Section C.15, and $\hat{\mathcal{H}}_1^{ipc}$ is a nonsingular $k_1 \times k_1$ stochastic factor "rotation" matrix, common to all sample dates $\tau=1,..,T$, which depends on the factor normalization imposed in the IPCA estimation algorithm.²²

Moreover, let $\hat{h}_{2,\tau}$ be the (RP-)PCA estimator of factors $h_{2,\tau}$ and $\hat{\lambda}_{2,i}=(\hat{\lambda}_{2,i}^{c\prime},\hat{\lambda}_{2,i}^{s\prime})'$ be the (RP-)PCA estimator of loadings for all $\tau=1,...,T$ and $i=1,...,N_1$, of model (A.4), as defined in Lettau and Pelger (2020b), and summarized in Section OA.3.1 of the *Supplementary Material*. The estimator $\hat{h}_{2,\tau}$ is asymptotically equivalent, up to negligible terms, to $\hat{\mathcal{H}}_2^{rppc}(h_{2,\tau}+u_{2,\tau}/\sqrt{N_2})$ where

(A.14)
$$u_{2,\tau} = \left(\frac{1}{N_2} \sum_{i=1}^{N_2} \lambda_{2,i} \lambda'_{2,i}\right)^{-1} \frac{1}{\sqrt{N_2}} \sum_{i=1}^{N_2} \lambda_{2,i} \varepsilon_{2,i,\tau},$$

and $\hat{\mathcal{H}}_{2}^{rppc}$ is a nonsingular $k_2 \times k_2$ stochastic factor "rotation" matrix, common to all sample dates $\tau = 1,...,T$, which depends on the factor normalization imposed in the RP_PCA estimation algorithm.²³ The asymptotic expansion of factor estimators $\hat{h}_{1,\tau}$ and $\hat{h}_{2,\tau}$ allows to obtain the infeasible asymptotic distribution of our test statistics $\hat{\xi}(k^c)$, as its asymptotic bias and variance depend on $\tilde{\Sigma}_{u,jk,\tau}^{ipc} = \operatorname{Cov}(u_{j,\tau}, u_{k,\tau} | \mathcal{F}_{\tau}^{ipc})$, that is the true and

²¹The explicit equation for $\lambda_{1,i,\tau}$ is given in (OA.11) in *Supplementary Material* Section OA.3.2.

²²Kelly et al. (2020) discuss extensively in their Sections 3, 4 and 6 the possible choices of the factor normalization and rotation in IPCA.

²³See Section 4 in Lettau and Pelger (2020a) for the definition and the properties of $\hat{\mathcal{H}}_2^{rppc}$.

unknown covariance between $u_{j,\tau}$ and $u_{k,\tau-h}$ conditional on the sigma field \mathcal{F}_{τ}^{ipc} , namely

$$\tilde{\Sigma}_{u,11,\tau}^{ipc} = \left(\frac{1}{N_1} \sum_{i=1}^{N_1} \lambda_{1,i,\tau} \lambda_{1,i,\tau}'\right)^{-1} \frac{1}{N_1} \sum_{i=1}^{N_1} \sum_{\ell=1}^{N_1} \lambda_{1,i,\tau} \lambda_{1,i,\tau}' \operatorname{Cov}(\varepsilon_{1,i,\tau}, \varepsilon_{1,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}) \left(\frac{1}{N_1} \sum_{i=1}^{N_1} \lambda_{1,i,\tau} \lambda_{1,i,\tau}'\right)^{-1}$$

$$(A.16) \qquad \tilde{\Sigma}_{u,22,\tau}^{ipc} = \left(\frac{1}{N_2} \sum_{i=1}^{N_2} \lambda_{2,i} \lambda_{2,i}'\right)^{-1} \frac{1}{N_2} \sum_{i=1}^{N_2} \sum_{\ell=1}^{N_2} \lambda_{2,i} \lambda_{2,i}' \operatorname{Cov}(\varepsilon_{2,i,\tau}, \varepsilon_{2,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}) \left(\frac{1}{N_2} \sum_{i=1}^{N_2} \lambda_{2,i} \lambda_{2,i}'\right)^{-1}$$

$$\tilde{\Sigma}_{u,12,\tau}^{ipc} = \left(\frac{1}{N_{1}} \sum_{i=1}^{N_{1}} \lambda_{1,i,\tau} \lambda_{1,i,\tau}'\right)^{-1} \frac{1}{\sqrt{N_{1}N_{2}}} \sum_{i=1}^{N_{1}} \sum_{\ell=1}^{N_{2}} \lambda_{1,i,\tau} \lambda_{2,i}' \operatorname{Cov}(\varepsilon_{1,i,\tau}, \varepsilon_{2,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}) \left(\frac{1}{N_{2}} \sum_{i=1}^{N_{2}} \lambda_{2,i} \lambda_{2,i}'\right)^{-1}$$

and $\tilde{\Sigma}_{u,21,\tau}^{ipc} = \tilde{\Sigma}_{u,12,\tau}^{ipc\prime}$. We set $\Sigma_{u,jk,\tau}^{ipc} = \underset{N_j,N_k \to \infty}{\text{plim}} \tilde{\Sigma}_{u,jk,\tau}^{ipc}$. The following Theorem A.1 provides the asymptotic distribution of the infeasible test statistic $\hat{\xi}(k^c)$. Note that, conditionally on $Z_{\tau-1}^{1:N}$, the values of the loadings $\lambda_{1,i,\tau} = C'Z_{i,\tau-1}$ are not stochastic.

THEOREM A.1 Under Assumptions A.1 - A.7 (in the OA) (b) and (c), and the null hypothesis $H_0 = H(k^c)$ of k^c common factors, we have:

$$(A.18) \quad \tilde{\xi}_{inf}^{ipc}(k^c) := N\sqrt{T} \cdot \left(\Omega_{U,1}^{ipc}\right)^{-1/2} \cdot \left[\hat{\xi}(k^c) - k^c + \frac{1}{2N}tr\left\{\tilde{\Sigma}_{cc}^{-1}\tilde{\Sigma}_{U}^{ipc}\right\}\right] \xrightarrow{d} N\left(0,1\right),$$

where $\tilde{\Sigma}_{cc} = \frac{1}{T} \sum_{\tau=1}^{T} \breve{f}_{\tau}^{c} \breve{f}_{\tau}^{c\prime}$, and

$$\begin{split} \tilde{\Sigma}_{U}^{ipc} &= & \frac{1}{T} \sum_{\tau=1}^{T} \left(\mu_{N}^{2} \tilde{\Sigma}_{u,11,\tau}^{ipc(cc)} + \tilde{\Sigma}_{u,22,\tau}^{ipc(cc)} - \mu_{N} \tilde{\Sigma}_{u,12,\tau}^{ipc(cc)} - \mu_{N} \tilde{\Sigma}_{u,21,\tau}^{ipc(cc)} \right), \\ \Omega_{U,1}^{ipc} &= & \frac{1}{2} \lim_{T \to \infty} E\left[\frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \Sigma_{U,\tau}^{ipc} \cdot \Sigma_{U,\tau}^{ipc} ' \right\} \right], \\ \Sigma_{U,\tau}^{ipc} &= & \mu^{2} \Sigma_{u,11,\tau}^{ipc(cc)} + \Sigma_{u,22,\tau}^{ipc(cc)} - \mu \Sigma_{u,12,\tau}^{ipc(cc)} - \mu \Sigma_{u,21,\tau}^{ipc(cc)}, \qquad h = ..., -1, 0, 1, ..., \end{split}$$

with
$$\check{f}_{\tau}^{c\prime} := f_{\tau}^c - \bar{f}^c$$
, and $\bar{f}^c = \sum_{\tau=1}^T f_{\tau}^c / T$.

Theorem A.1 corresponds to Theorem 1 in AGGR, where the estimator of the canonical correlations of the estimated

factors (used to compute $\hat{\xi}(k^c)$), and the sample covariance matrix of the true factors $\tilde{\Sigma}_{cc}$ have been modified to take into account that the factors are allowed to have a non-zero mean, and that the factor loadings for group 1 are allowed to be time-varying as in the IPCA model. We note that matrices $\tilde{\Sigma}_{u,jk,\tau}^{ipc}$ and $\Sigma_{u,jk,\tau}^{ipc}$, with j,k=1,2, do depend on time as the loadings are time varying in the IPCA model, differently from AGGR where the loadings were constant over time as they considered a static factor model estimable with classical PCA.

To get the distribution of the feasible counterpart of the infeasible statistic $\tilde{\xi}_{inf}^{ipc}(k^c)$, we need consistent estimators for the unknown scalar tr $\left\{\tilde{\Sigma}_{cc}^{-1}\tilde{\Sigma}_{U}^{ipc}\right\}$ and matrix $\Omega_{U,1}^{ipc}$ in Theorem A.1.

Theorem A.2 shows the asymptotic distribution of the feasible version of the test statistics $\tilde{\xi}_{inf}^{ipc}(k^c)$ obtained by replacing $\tilde{\Sigma}_{cc}$ with its large sample limit I_{k^c} , and both matrices $\tilde{\Sigma}_{U}^{ipc}$ and $\Omega_{U,1}^{ipc}$ by consistent estimators, which are built starting from the following quantities:

$$\hat{\Sigma}_{u,11,\tau}^{ipc} = \left(\frac{1}{N_{1}} \sum_{i=1}^{N_{1}} \hat{\lambda}_{1,i,\tau} \hat{\lambda}'_{1,i,\tau}\right)^{-1} \left[\frac{1}{N_{1}} \sum_{i=1}^{N_{1}} \sum_{\ell=1}^{N_{1}} \hat{\lambda}_{1,i,\tau} \hat{\lambda}'_{1,i,\tau} \cdot \mathbf{1} \{|c\hat{orr}_{T}(\hat{\varepsilon}_{1,i,\tau},\hat{\varepsilon}_{1,\ell,\tau})| > c^{*}\} \cdot \check{\varepsilon}_{1,i,\tau} \cdot \check{\varepsilon}_{1,\ell,\tau}\right] \\
(A.19) \times \left(\frac{1}{N_{1}} \sum_{i=1}^{N_{1}} \hat{\lambda}_{1,i,\tau} \hat{\lambda}'_{1,i,\tau}\right)^{-1}$$

$$\begin{split} \hat{\Sigma}_{u,22,\tau}^{ipc} &= \left(\frac{1}{N_2}\sum_{i=1}^{N_2}\hat{\lambda}_{2,i}\hat{\lambda}_{2,i}'\right)^{-1}\left[\frac{1}{N_2}\sum_{i=1}^{N_2}\sum_{\ell=1}^{N_2}\hat{\lambda}_{2,i}\hat{\lambda}_{2,i}'\cdot\mathbf{1}\{|c\hat{orr}_T(\hat{\varepsilon}_{2,i,\tau},\hat{\varepsilon}_{2,\ell,\tau})|>c^*\}\cdot\check{\varepsilon}_{2,i,\tau}\cdot\check{\varepsilon}_{2,\ell,\tau}\right]\\ &\times\left(\frac{1}{N_2}\sum_{i=1}^{N_2}\hat{\lambda}_{2,i}\hat{\lambda}_{2,i}'\right)^{-1} \end{split}$$

$$\hat{\Sigma}_{u,12,\tau}^{ipc} = \left(\frac{1}{N_1} \sum_{i=1}^{N_1} \hat{\lambda}_{1,i,\tau} \hat{\lambda}'_{1,i,\tau}\right)^{-1} \left[\frac{1}{\sqrt{N_1 N_2}} \sum_{i=1}^{N_2} \sum_{\ell=1}^{N_2} \hat{\lambda}_{1,i,\tau} \hat{\lambda}'_{2,i} \cdot \mathbf{1} \{|c\hat{orr}_T(\hat{\varepsilon}_{1,i,\tau},\hat{\varepsilon}_{2,\ell,\tau})| > c^*\} \cdot \check{\varepsilon}_{1,i,\tau} \cdot \check{\varepsilon}_{2,\ell,\tau}\right] \\
\times \left(\frac{1}{N_2} \sum_{i=1}^{N_2} \hat{\lambda}_{2,i} \hat{\lambda}'_{2,i}\right)^{-1} \tag{A.21}$$

(A.22)
$$\hat{\Sigma}_{u,21,\tau}^{ipc} = \hat{\Sigma}_{u,12,\tau}^{ipc\prime},$$

where $\hat{\varepsilon}_{j,i,\tau} = y_{j,i,\tau} - \hat{\lambda}_{j,i}^{c} f_{\tau}^{c} - \hat{\lambda}_{j,i}^{s} f_{j,\tau}^{c}$, $c\hat{orr}_{T}(\hat{\varepsilon}_{j,i,\tau},\hat{\varepsilon}_{k,\ell,\tau})$ is the sample correlation between $\hat{\varepsilon}_{j,i,\tau}$ and $\hat{\varepsilon}_{k,\ell,\tau}$ computed over the entire sample of T dates, for $j,k=1,2,c^{*}\in(0,1)$ is a trimming constant, $\check{\varepsilon}_{j,i,\tau}:=\hat{\varepsilon}_{j,i,\tau}-\bar{\varepsilon}_{j,i,\tau}$, and $\bar{\hat{\varepsilon}}_{j,i,\tau}=\sum_{\tau=1}^{T}\hat{\varepsilon}_{j,i,\tau}/T$, and $\mathbf{1}\{condition\}$ is the indicator function which is equal to 1 if the condition is satisfied, and 0 otherwise. Importantly, as we allow the errors to be weakly correlated in the cross-section with sparse covariance (and therefore correlation) matrix, but not over time, consistent estimation of $\tilde{\Sigma}_{U}^{ipc}$ and $\Omega_{U,1}^{ipc}$ require thresholding of estimated cross-sectional covariances appearing in the terms $\hat{\Sigma}_{u,jk,\tau}^{ipc}$, with j,k=1,2.

THEOREM A.2 Define $\hat{\Sigma}_{U}^{ipc} := \frac{1}{T} \sum_{\tau=1}^{T} \hat{\Sigma}_{U,\tau}^{ipc}$, and $\hat{\Omega}_{U,1}^{ipc} := \frac{1}{2} \frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \hat{\Sigma}_{U,\tau}^{ipc} \cdot \hat{\Sigma}_{U,\tau}^{ipc} \right\}$, where $\hat{\Sigma}_{U,\tau}^{ipc} := \hat{\mu}_{N}^{2} \hat{\Sigma}_{u,11,\tau}^{ipc(cc)} + \hat{\Sigma}_{u,22,\tau}^{ipc(cc)} - \hat{\mu}_{N} \hat{\Sigma}_{u,12,\tau}^{ipc(cc)} - \hat{\mu}_{N} \hat{\Sigma}_{u,21,\tau}^{ipc(cc)}$, with terms $\hat{\Sigma}_{u,jk,\tau}^{ipc}$, for j, k = 1, 2 defined in equations (A.19) - (A.22), and $\hat{\mu}_{N} := \sqrt{N_{2}/N_{1}}$. Define the test statistic:

$$(\text{A.23}) \qquad \qquad \tilde{\xi}^{ipc}(k^c) := N\sqrt{T} \left(\hat{\Omega}_{U,1}^{ipc}\right)^{-1/2} \left[\hat{\xi}(k^c) - k^c + \frac{1}{2N} tr \left\{\hat{\Sigma}_U^{ipc}\right\}\right],$$

and let Assumptions A.1 - A.7 (in the Supplementary Material) hold. Then: (i) Under the null hypothesis $H_0 = H(k^c)$ we have: $\tilde{\xi}^{ipc}(k^c) \xrightarrow{d} N(0,1)$. (ii) Under the alternative hypothesis $H_1 = \bigcup_{0 \le r < k^c} H(r)$, we have: $\tilde{\xi}^{ipc}(k^c) \xrightarrow{p} -\infty$.

Our Theorem A.2 corresponds to Theorem 2 in AGGR where their original estimators of the residuals' variances has been modified to take into account that, in the more general set up of this paper, i) the observed variables $y_{j,i,\tau}$ are not demeaned and the estimated factor models do not include a constant so the residuals of these models are not necessarily zero mean (in case the estimated factors do not perfectly explain the risk premia of all the assets in both groups)²⁵, and ii) as the errors are allowed to be weakly correlated across individuals, consistent estimation of $\tilde{\Sigma}_U^{ipc}$ and $\Omega_{U,1}^{ipc}$ requires thresholding of the cross-sectional covariances between the residuals $\hat{\varepsilon}_{j,i,\tau}$ and $\hat{\varepsilon}_{k,\ell,\tau}$ when either $j \neq k$, or $i \neq \ell$, or both.

 $^{^{24}}$ As $\hat{corr}_T(\hat{\varepsilon}_{j,i,\tau},\hat{\varepsilon}_{k,\ell,\tau})=1$ when j=k, and $i=\ell$, variances are never trimmed by the indicator function, which therefore thresholds only covariances among different assets, both within group and across the two groups.

²⁵This could not happen in the framework of AGGR who considered classical PCA applied to demeaned data, as typically done in classical PCA. For more details on this subtle issue, with references the different approaches used in the finance literature to estimate PCs, see the discussion in *Supplementary Material* Section OA.3.1.

Three Common Factors

Supplementary Material

Elena Andreou* Patrick Gagliardini[†] Eric Ghysels[‡] Mirco Rubin[§]

This Draft: October 11, 2025

^{*}University of Cyprus and CEPR, e-mail: elena.andreou@ucy.ac.cy.

 $^{^\}dagger Universit\grave{a}$ della Svizzera italiana, Lugano and Swiss Finance Institute, e-mail: patrick.gagliardini@usi.ch.

[‡]University of North Carolina - Chapel Hill and CEPR, e-mail: eghysels@unc.edu.

[§]EDHEC Business School, Nice, e-mail: mirco.rubin@edhec.edu.

OA.1 Data Description - Individual Stocks and Portfolios

We consider three panels of monthly returns in our analysis, namely (i) individual US stock returns from CRSP, (ii) the panel of test asset portfolios from the April 2021 release of the database "Open Source Cross-Sectional Asset Pricing" created by Chen and Zimmermann (2022), CZ21 hearafter, and (iii) the panel of factors from the zoo considered by CZ21.1 For all three panels, we consider two samples: (i) the chronological time sample which includes all data available in each dataset from Jan. 1966 to Dec. 2020, and (ii) the publication time sample which goes from Jan. 1996 to Dec. 2020, where the CZ21 test assets portfolios and factors enter with their publication date in the database. We split the 660 (resp. 300) months in the chronological time (resp. publication time) sample into B = 11 (resp. 5) non-overlapping blocks of 60 months, denoted as b = 1, ..., B. The first block in the chronological time (resp. publication time) sample is from Jan. 1966 to Dec. 1970 (resp. Jan. 1996 to Dec. 2000) and the last block is from Jan. 2016 to Dec. 2020. Within each block, we consider only a balanced sample of individual stocks and test asset portfolios, that is we only include assets with returns available for all the 60 months. We work with 5-year non-overlapping samples to address the concern of survivorship bias if we were to use the full sample of individual stocks. Similar to the arguments in Kim and Korajczyk (2021), one can view the 5-year span as a compromise between a sample large enough for our test procedure to have desirable finite sample properties and the concern of capturing new and disappearing stocks. Figure OA.1 displays the number of individual CRSP stocks, test assets portfolios and factors available in each of the B blocks in the chronological time and publication time samples, respectively. Both samples are described in more detail below.

¹Data for the "Open Source Cross-Sectional Asset Pricing" project are available on: https://www.openassetpricing.com/

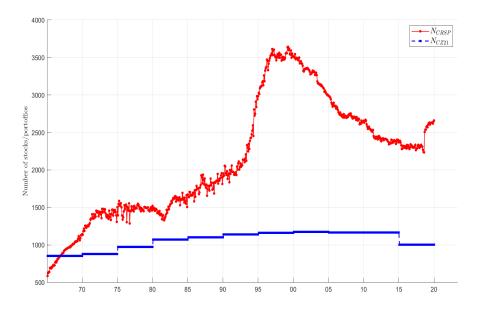
The first chronological time sample panel of test assets consists of individual stocks available from the Center for Research in Security Prices (CRSP) traded on the New York Stock Exchange (NYSE), the American Stock Exchange (AMEX) and the NASDAQ for the period from January 1966 through December 2020. We focus on common stocks (CRSP share codes 10 and 11) and delete all stocks having less than 60 consecutive monthly returns. We end up having an unbalanced panel for the returns of 14948 different stocks. The average cross-sectional size, computed in each month, is about 4270 stocks. In the first (resp. last) block, that is the block 1966-1970 (resp. 2016-2020), we have 1539 (resp. 2668) stocks. The publication time sample is constructed analogously but goes from January 1996 to December 2020. Applying the same filters as above, we end up having an unbalanced panel for the returns of 8131 different stocks. The average cross-sectional size, computed in each month, is about 4170 stocks. In the first (resp. last) block, that is the block 1996-2000 (resp. 2016-2020), we have 3779 (resp. 2668) stocks.

Turning to the test assets portfolios and factors from CZ21, we consider the unbalanced panel of 1215 portfolios formed starting from the 205 firm-level characteristic, or predictors, having predictive ability for firm-level returns according to the four asset pricing meta-studies by McLean and Pontiff (2016), Green, Hand, and Zhang (2017), Hou, Xue, and Zhang (2020), Harvey, Liu, and Zhu (2016). The returns of the 205 factors in our zoo panel are those of long-short portfolios of the upper and bottom quantile portfolios constructed by sorting stocks according to each characteristic. Following CZ21, we consider test asset portfolios and factors associated only with characteristics classified either as "clearly" or "likely" returns predictors in their study.

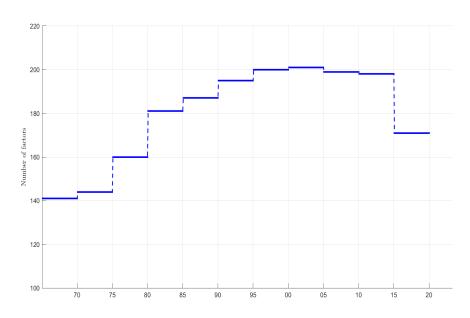
²CZ21 construct factors following the methodology of the papers where they have been introduced, therefore most factors are constructed from long-short portfolios of equal-weighted quintiles. Value-weighting or other quantiles are used in the factor construction only for the few papers that emphasize these constructions.

 $^{^3}$ CZ21 define as "clear predictor" a characteristic which is expected to achieve statistically significant mean raw returns in long-short portfolios (e.g. t-stat > 2.5 in a long-short portfolio, monotonic portfolio sort with 80 bps spread, t-stat > 4 in a regression, t-stat > 3 in 6-month event study). On the other hand, a "likely predictor" is a characteristic

Figure OA.1: Number test assets and factors in the Zoo, full sample: 1966-2020



(a) CRSP and CZ21 test assets



(b) CZ21 Factors in the zoo

Panel (a) displays the number of assets in the balanced panel of CZ21 test assets (blue horizontal segments) included in the 5-years block of monthly data ending in year t, for each t=1970,...,2020, and the number of stocks (red circles) available in our unbalanced panel of individual stocks for every month in our sample ranging from Jan. 1966 to Dec. 2020. Panel (b) displays the number of factors in our zoo, that is the number of factors in the Chen and Zimmermann (2022) dataset (blue horizontal segments and dashed lines). In every 5-year window we include in our analysis and report in this figure quantile portfolios and factors from the zoo (long-short portfolios) with non-missing returns for all the 60 months as available in the CZ21 dataset. On the other hand, for the IPCA estimation we consider the unbalanced panel of individual stocks for which all the returns are available in CRSP and all the 35 stock-specific and time-varying characteristics considered in Kelly, Pruitt, and Su (2019) and originally proposed by Freyberger, Neuhierl, and Weber (2020).

In the chronological time sample we include all the quantile portfolios and factors available in the baseline version of the database of CZ21, leading to an unbalanced panel of 1214 portfolios associated with 205 characteristics.⁴ The average cross-sectional size, computed in each month, is about 1113 portfolios. In the first (resp. last) block, that is the block 1966-2000 (resp. 2016-2020), we have 855 (resp. 1001) test asset portfolios and 141 (resp. 171) factors.

In the publication time sample we include all the quantile portfolios available in the baseline version of the database of CZ21, after excluding all (binary) portfolios associated to binary characteristics. This leads to 1159 portfolios associated to 177 characteristics.⁵ In each 5-year block going from January of year t-4 to December of year t, a factor and the relative test assets portfolios from CZ21 are included for all the dates corresponding to the rolling window only if the paper introducing the factor was published in year t+1, or before.⁶ These choices allow us to have in the first rolling window (resp. the last), that is the window 1996-2000 (resp. 2015-2020), 59 (resp. 171) factors, and 276 (resp. 959) test asset portfolios.

Finally, for both samples we also download from Kenneth French website the 5 Fama and French factors: Market, SMB, HML, Operating Profitability (RMW), and Investment Style (CMA), together with the momentum factor (and based on prior 2-12 months returns), and the 1 month risk free rate which is used to compute excess returns for the panels of test assets.

expected to achieve borderline evidence for the significance of mean raw returns in long-short portfolios (e.g. t-stat = 2.0 in long-short with factor adjustments, t-stat between 2 and 3 in a regression, large t-stat in 3-day event study).

⁴For 28 characteristics only 2 quantile portfolios are available, for 7 characteristics 3 quantile portfolios are available, for 5 characteristics 4 quantile portfolios are available, for 1 characteristic 5 quintile portfolios are available, for 1 characteristic 6 quantile portfolios are available, for 1 characteristic 7 quantile portfolios are available, and finally for 58 characteristics all 10 decile portfolios are available.

⁵More precisely, for 7 characteristics only 3 quantile portfolios are available, for 5 characteristics 4 quantile portfolios are available, for 105 characteristics 5 quintile portfolios are available, for 1 characteristic 6 quantile portfolios are available, for 1 characteristic 7 quantile portfolios are available, and finally for 58 characteristics all 10 decile portfolios are available.

⁶Publication dates are also available in CZ21.

OA.2 Macroeconomic indicators/factors

In Table OA.1 we display the information about the macroeconomic indicators/factors used in Section D

Table OA.1: Macro Indicators/Factors

Full name	Description	Data Source
Term Spread 10y-1y	Difference between the 10-year Treasury bond and the one-year constant maturity Treasury bill rates	FRED
Term Spread 10y-3m	Difference between yields on 10-year Treasury and 3-month T-bill rates	FRED
Term Spread 1y-FEDFund	Term spread is proxied by the difference between yields on 1-year Treasury and Federal Funds Rate.	FRED
Default spread BAA-AAA	The default spread is the difference between Moody's Baa and Aaa corporate bond yields	FRED
Cochrane Piazzesi factor	A single linear combination of forward rates	CRSP
Consumption growth	Consumption growth (non-durables plus services)	FRED
Labor Income growth	Labor income growth	FRED
IP growth	Industrial Production growth	FRED
IP growth innov	AR(1) innovations in Industrial Production growth	FRED
L&N Macro PC (1:3) VAR innov	VAR(1) innovations in the first three PCs (1:3) of 279 macro-finance variables from Ludvigson and Ng (2009)	Ludvigson and Ng (2009)
M&N Macro FRED-MD PCs (1:4)	First four PCs (1:4) estimated from the Macro variables in the FRED-MD panel from McCracken and Ng (2016)	FRED-MD
Inflation	Inflation (t)	FRED
Lagged Inflation	Inflation $(t-1)$	FRED
Expected Inflation	Expected Inflation based on Fama and Gibbons (1984)	FRED
Unexpected Inflation	Unexpected inflation based on Fama and Gibbons (1984)	FRED
Inflation innov	AR(1) innovations in inflation	FRED
EPU	Economic Policy Uncertainty level	FRED
EPU growth	Economic Policy Uncertainty growth rate	FRED
Macro Uncertainty (h=1m, 3m, 12m)	Jurado, Ludvigson, and Ng (2015) Macro Uncertainty indexes for horizons $h=1,3$ and 12 months	Jurado et al. (2015)

All Macro indicators/factors have been downloaded or reconstructed at monthly frequency. The Cochrane Piazzesi factor is based on the computations based on their original replication code and updated Fama-Bliss discount bonds data from CRSP. Similarly, the methodology from Fama and Gibbons (1984) is applied to decompose the growth rate of the CPI index from FRED into Expected and Unexpected inflation.

OA.3 Factor estimation: PCA and its recent extensions

In this section we discuss extensions of PCA to RP-PCA of Lettau and Pelger (2020a and 2020b) and IPCA of Kelly et al. (2019). A subsection is devoted to each case.

OA.3.1 RP-PCA estimation

To simplify the exposition we will use a generic notation here for the discussion of various estimators which can be applied to different panel data settings. Let y_{τ} be N-dimensional vector of returns, and assume that the data generating process of y_{τ} is a linear factor model as the APT of Ross (1976), that is:

$$y_{\tau} = \Lambda h_{\tau} + \varepsilon_{\tau}, \qquad \tau = 1, ..., T,$$
 (OA.1)

where h_{τ} is the (k,1) vector of unobservable factors with expected value $\mu_h \equiv E[h_{\tau}]$, possibly different from zero, $\Lambda = [\lambda_1, ..., \lambda_N]'$ is the (N, k) full column rank matrix of unknown loadings, and the idiosyncratic innovations $\mathbb{E}[\varepsilon_{\tau}] = 0$. These assumptions imply $\mathbb{E}[y_{\tau}] = \Lambda \mu_h$, possibly different from zero. Model (OA.1) can be written as:

$$Y = H\Lambda' + \varepsilon, \tag{OA.2}$$

where $Y = [y_1, ..., y_T]'$ is the (T, N)-dimensional matrix of observed excess returns, and $H = [h_1, ..., h_T]'$ is the (T, k)-dimensional matrix of factor values.

Lettau and Pelger (2020a, 2020b) and Zaffaroni (2025) suggest that estimating model (OA.2) by performing PCA on the demeaned returns $\tilde{y}_{\tau} = y_{\tau} - \bar{y}$, as typically done in the finance and macroeconomics literature, is restrictive as the mean of the factors and the returns should contain

information on the factor structure. Let $\tilde{h}_{\tau}=h_{\tau}-\bar{h}$ be the demeaned factors, and $\bar{y}_i=\frac{1}{T}\sum_{\tau=1}^T y_{i,\tau}$ be the time series mean of the returns of the *i*-th asset, with i=1,...,N. Lettau and Pelger (2020a) address the estimation of the non-demeaned factors in model (OA.1) with their RP-PCA procedure, which consists in solving the following minimization problem:

$$\min_{\lambda_1, \dots, \lambda_N, } \frac{1}{NT} \sum_{i=1}^N \sum_{\tau=1}^T (\tilde{y}_{i,\tau} - \tilde{h}'_{\tau} \lambda_i)^2 + (1 + \gamma_{RP}) \frac{1}{N} \sum_{i=1}^N (\bar{y}_i - \bar{h}' \lambda_i)^2, \qquad (OA.3)$$

$$h_1, \dots, h_T,$$

subject to normalization restrictions. The first double summation in (OA.3) corresponds to the average unexplained (time-series) variation of the data, the second summation correspond to the (cross-sectional) average of the squared "pricing errors" across all N assets, and $\gamma_{RP} \in [-1, +\infty)$ is a constant which can be interpreted as a tuning parameter: as it increases more weight is given to the pricing errors in the factor estimation. Lettau and Pelger (2020a) show that the solution to (OA.3) can be obtained by performing the following two steps:

(i) Estimate the loading matrix Λ is as \sqrt{N} times the (N,k) matrix of the eigenvectors associated to the largest k eigenvalues of matrix

$$M_{RP}(\gamma_{RP}) := \frac{1}{T} \sum_{\tau=1}^{T} y_{\tau} y_{\tau}' + \gamma_{RP} \left(\frac{1}{T} \sum_{\tau=1}^{T} y_{\tau} \right) \left(\frac{1}{T} \sum_{\tau=1}^{T} y_{\tau} \right)'$$
 (OA.4)

The estimated loadings, that we denote as $\hat{\Lambda}_{RP}$, are such that $\hat{\Lambda}'_{RP}\hat{\Lambda}_{RP}/N = I_k$.

 $^{^7}$ Lettau and Pelger (2020a), in their online appendix, show that $\hat{\Lambda}_{RP}$ can be obtained as the conventional PCA estimator of the loadings applied to the "projected" model: $\check{Y} = \check{H}\Lambda + \check{\varepsilon}$ where $\check{Y} := W(\gamma_{RP})Y$, $\check{H} := W(\gamma_{RP})H$, $\check{\varepsilon} := W(\gamma_{RP})\varepsilon$, and $W(\gamma_{RP}) = I_T + (\sqrt{\gamma_{RP}+1}-1)(\mathbb{1}_T\mathbb{1}_T')/T$. That is, the loading matrix $\Lambda_{j,RP}$ can be estimated as \sqrt{N} times the (N,k) matrix of the eigenvectors associated to the largest k eigenvalues of $M_{RP}(\gamma_{RP}) = \check{Y}'\check{Y}$ /T.

(ii) Estimate the latent factors in model (OA.1) at each date τ by a cross-sectional regression of the returns y_{τ} on the estimated loadings $\hat{\Lambda}_{RP}$:

$$\hat{h}_{\tau,RP} := \left(\hat{\Lambda}'_{RP}\hat{\Lambda}_{RP}\right)^{-1}\hat{\Lambda}'_{RP} y_{\tau}. \tag{OA.5}$$

We denote as $\hat{H}_{RP} = [\hat{h}_{1,RP}, ..., \hat{h}_{T,RP}]'$ the (T,k) matrix of estimated factors.

As linear latent factor models are identified up to an invertible transformation, an equivalent estimator \hat{H}_{RP}^* of the factors is obtained by rescaling \hat{H}_{RP} such that the (uncentered) second moment of the estimated factors is $\hat{H}_{RP}^{*\prime}\hat{H}_{RP}^*/T=I_k$, that is $\hat{H}_{RP}^*:=\hat{H}_{RP}\left(\hat{H}_{RP}'\hat{H}_{RP}/T\right)^{-1/2}$. Following Lettau and Pelger (2020a), we refer to \hat{H}_{RP} and \hat{H}_{RP}^* as "RP-PCA estimators". Importantly, the factors estimated by RP-PCA have a mean $\hat{h}=\frac{1}{T}\sum_{\tau=1}^T\hat{h}_{\tau}$ which is not necessarily equal to zero. In fact, equation (OA.5) shows that $\hat{h}_{\tau,RP}$ is a linear combination of the original returns which are not-demeaned.8

Special case of the RP-PCA: $\gamma_{RP}=0$

When $\gamma_{RP}=0$, the matrix $M_{RP}(\gamma_{RP})$ characterizing the RP-PCA estimator (OA.5) coincides with the conventional PCA estimator but with loadings estimated from the uncentered second moment matrix of the returns $M_{RP}(\gamma_{RP}=0)=\frac{1}{T}\sum_{\tau=1}^{T}y_{\tau}y_{\tau}'$. The RP-PCA estimators of the factors and the loadings with $\gamma_{RP}=0$ coincides with those proposed by Zaffaroni (2025).

⁸See Section 3 of Zaffaroni (2025) showing that the estimated factors $\hat{h}_{\tau,RP}$ are portfolio (excess-) returns, and correspond to "the feasible PCA-estimators" of the infeasible "mimicking portfolios" (of the true latent factors) proposed by Huberman, Kandel, and Stambaugh (1987) and Breeden, Gibbons, and Litzenberger (1989). See Lehmann and Modest (2005) for a discussion of factor-mimicking portfolio estimators.

Special case of the RP-PCA: $\gamma_{RP}=-1$

When $\gamma_{RP}=-1$ the RP-PCA estimator of the loadings, denoted by $\hat{\Lambda}_{PCA}$, is computed as \sqrt{N} times the eigenvectors of the sample variance-covariance matrix of the returns $\hat{V}(y_{\tau})=M_{RP}(\gamma_{RP}=-1)=\frac{1}{T}\sum_{\tau=1}^{T}\tilde{y}_{\tau}\tilde{y}_{\tau}'$. We denote the RP-PCA factor estimator in this special case as $\hat{h}_{\tau,PCA}$:

$$\hat{h}_{\tau,PCA} := \left(\hat{\Lambda}'_{PCA}\hat{\Lambda}_{PCA}\right)^{-1}\hat{\Lambda}'_{PCA}y_{\tau}, \tag{OA.6}$$

and name $\hat{\Lambda}_{PCA}$ and $\hat{h}_{\tau,PCA}$ as the "conventional PCA" estimators of the loadings and factors, respectively, as they are used by most of the financial literature. For instance, the factor estimators used in Connor and Korajczyk (1988), Lehmann and Modest (2005), Kozak, Nagel, and Santosh (2018), Kozak, Nagel, and Santosh (2020), Giglio and Xiu (2021), and Pukthuanthong, Roll, and Subrahmanyam (2019), among others all coincide with $\hat{h}_{\tau,PCA}$. Another frequently used estimator, denoted by $\hat{h}_{\tau,PCA}$, is obtained by a cross-sectional regression of the demeaned returns \tilde{y}_{τ} on $\hat{\Lambda}_{PCA}$:

$$\hat{\tilde{h}}_{\tau,PCA} := \left(\hat{\Lambda}'_{PCA}\hat{\Lambda}_{PCA}\right)^{-1}\hat{\Lambda}'_{PCA}\,\tilde{y}_{\tau}. \tag{OA.7}$$

Differently from $\hat{h}_{\tau,PCA}$ and $\hat{h}_{\tau,RP}$, factors $\hat{\tilde{h}}_{\tau,PCA}$ have zero-mean as they are linear combinations of the demeaned data \tilde{y}_{τ} .

Let $\tilde{Y}=[\tilde{y}_1,...,\tilde{y}_T]'$ be (T,N) matrix collecting the demeaned returns. AGGR consider the estimator $\hat{\tilde{H}}^*_{PCA}=[\hat{\tilde{h}}^*_{1,PCA},...,\hat{\tilde{h}}^*_{T,PCA}]'$ of the k factors which is defined as \sqrt{T}

⁹As discussed in Section 2 of Zaffaroni (2025), $\tilde{h}_{\tau,PCA}$ is the estimator of the (demeaned) latent factors $\tilde{h}_{\tau}:=h_{\tau}-\bar{h}$ of model (OA.1) for the demeaned data \tilde{y}_{τ} . This can be easily seen by noting that the model for the demeaned data can be written as: $\tilde{y}_{\tau}=\Lambda_{j}(h_{\tau}-\bar{h})+(\varepsilon_{\tau}-\bar{\varepsilon})$.

times the eigenvectors associated to the k largest eigenvalues of the matrix $\frac{1}{NT}\tilde{Y}\tilde{Y}'$. By construction the estimated factors are zero mean, and their (sample) variance-covariance matrix is $\hat{H}^{*\prime}_{PCA}\hat{H}^*_{PCA}/T=I_k$. Using the arguments in Bai and Ng (2002), it can be shown that $\hat{H}^{*\prime}_{PCA}$ is equal to the PCA estimator in (OA.7) rescaled to have unit variance: $\hat{H}^*_{PCA}=\hat{H}_{PCA}(\hat{H}'_{PCA}\hat{H}_{PCA}/T)^{-1/2}$.

Importantly, $\hat{h}_{\tau,PCA}$ and $\hat{h}_{\tau,PCA}^*$ are consistent estimators of the latent factors only when these are assumed to have zero expected value, as in Assumption A.2 of AGGR. In the next Section A.2 we show that relaxing this assumption does not change the main results of their paper, but requires modifications to their canonical correlations estimator as well as other statistics.

OA.3.2 IPCA estimator

The original IPCA model specification in Kelly et al. (2019) is:

$$y_{\tau} = \Lambda_{\tau}^{ipc} \cdot h_{\tau} + \varepsilon_{\tau} \,, \tag{OA.8}$$

with

$$\Lambda_{\tau}^{ipc} = [\lambda_{1,\tau}^{ipc}, ..., \lambda_{i,\tau}^{ipc}, ..., \lambda_{i,N_{\tau}}^{ipc}]' := Z_{\tau-1}\Gamma_{\Lambda} + \nu_{\tau}, \tag{OA.9}$$

where $Z_{\tau-1}=[z_{1,\tau-1},...,z_{i,\tau-1},...,z_{N_{\tau},\tau-1}]'$ is an $N_{\tau-1}\times L$ matrix containing the lagged values of all the L company-specific characteristics (collected in the L-dimensional vector $z_{i,\tau-1}$ for each stock i) for the $N_{\tau-1}=N_{\tau}$ stocks for which both returns and lagged characteristics are observable at dates τ and $\tau-1$, respectively. The N_{τ} -dimensional vector ε_{τ} collect the idiosyncratic innovations, and is such that $\mathbb{E}[\varepsilon_{\tau}]=0$, the $L\times K$ matrix Γ_{Λ} maps characteristics into loadings of the K factors h_{τ} , and the N_{τ} -dimensional unobservable vector ν_{τ} captures any residual feature

loadings that is orthogonal to the characteristics, and does not affect the estimation of the model. Each row of the system of equations in (OA.8) is

$$y_{i,\tau} = \lambda_{i,\tau}^{ipc\prime} h_{\tau} + \varepsilon_{i,\tau}, \tag{OA.10}$$

where

$$\lambda_{i,\tau}^{ipc} := \Gamma_{\Lambda}' z_{i,\tau-1} + \nu_{i,\tau}. \tag{OA.11}$$

Kelly et al. (2019) propose a recursive procedure that delivers estimates of the matrices Γ_{Λ} and of the $T \times K$ factors $h := [h_1, h_2, ..., h_T]'$. The estimator $\{\hat{h}^{ipc}, \hat{\Gamma}^{ipc}_{\Lambda}\}$ minimizes the sum of squared errors $\sum_{\tau=1}^T (r_{\tau} - Z_{\tau-1} \Gamma^{ipc}_{\Lambda} h^{ipc}_{\tau})'(r_{\tau} - Z_{\tau-1} \Gamma^{ipc}_{\Lambda} h^{ipc}_{\tau})$, under the constraints given by the IPCA identification conditions: factors are orthogonal among each other, have positive mean (i.e. their ex-post risk premium is positive), and $\Gamma'_{\Lambda} \Gamma_{\Lambda} = I_{K}$.

The first order conditions of the constrained optimization problem allow to obtain the factor estimator

$$\hat{h}_{\tau}^{ipc} = \left(\hat{\Gamma}_{\Lambda}^{ipc\prime} Z_{\tau-1}' Z_{\tau-1} \hat{\Gamma}_{\Lambda}^{ipc}\right)^{-1} \cdot \hat{\Gamma}_{\Lambda}^{ipc\prime} Z_{\tau-1}' y_{\tau} , \qquad (OA.12)$$

and the estimator of (the vectorized version) of matrix Γ_{Λ} :

$$\operatorname{vec}(\hat{\Gamma}_{\Lambda}^{ipc\prime}) = \left(\sum_{\tau=1}^{T} Z_{\tau-1}' Z_{\tau-1} \otimes \hat{h}_{\tau} \hat{h}_{\tau}^{ipc\prime}\right)^{-1} \cdot \sum_{\tau=1}^{T} \left[Z_{\tau-1} \otimes \hat{h}_{\tau}^{ipc\prime}\right]' y_{\tau} . \tag{OA.13}$$

The estimates of h_{τ}^{ipc} and Γ_{Λ}^{ipc} are obtain by an alternating least squares (ALS) procedure which iterates between (OA.12) and (OA.13) until convergence, starting from classical PCA estimation

of the $L \times K$ loadings matrix for the panel of returns of the so-called "managed" portfolios defines as $Z'_{\tau-1}r_{\tau} = [z'_{1,\tau}r_{\tau},...,z'_{N_{\tau},\tau}r_{\tau}]$. This approach for the choice of the starting value of the ALS procedure is justified by the observation that, if the stock-specific characteristics were constant over time (which is a reasonable assumption for characteristics such as companies' financial rations observed at monthly frequencies), then $Z_{\tau} = Z$ for all τ , and the panel of returns of the manged portfolios is $Z'r_{\tau} = Z'Z\Gamma_{\Lambda}h_{\tau} + (\varepsilon_{\tau} + Z'\nu_{\tau})$, i.e. it has the same structure of an approximate linear factor model if appropriate assumptions are satisfied for the ν_{τ} and the idiosyncratic innovations ε_{τ} , where the factor loadings loadings $Z'Z\Gamma_{\Lambda}$ are a linear transformation (given by the $L \times L$ constant matrix Z'Z) of matrix Γ_{Λ} .

Finally, we note that, analogously to Kelly et al. (2019) in our specification of the IPCA model we always include the constant characteristic which assumes value of 1 for all stocks and all dates, i.e. the first element in the L-dimensional vector $z_{i,\tau}$ is set to 1 for all stocks i and dates τ . A risk factor that is loading mostly onto this constant characteristic, compared to the others, can be interpreted as a special managed portfolio corresponding to an equally-weighted equity market factor. The fact that we use IPCA, implies that the return i-th generic stock is included in our dataset in month τ only if both its return $r_{i,\tau}$ in month τ and all the 35 stock-specific and timevarying characteristics at time $\tau-1$ are not-missing.

OA.4 Assumptions and proofs

Section OA.4.1 includes all the Assumptions required to prove Proportion A.1, Theorem A.1 and Theorem A.2 in Appendix A. Section OA.4.2 provides the proof of Proposition A.1, while Sections OA.4.3 and OA.4.4 provide the proofs of Theorems A.1 and A.2, respectively.

In this appendix, we denote by $a_{\tau}=[A]_{\tau}$ the column vector corresponding to the τ -th row a'_{τ} of matrix $A=[a_1,...,a_{\tau},...,a_T]'$.

OA.4.1 Assumptions for Proposition A.1 and Theorems A.1 and A.2

We make the following assumptions:

Assumption A.1. We have $N_1, N_2, T \to \infty$ such that the conditions in (A.9) hold, that is: $\sqrt{T}/N = o(1), N/T^2 = o(1),$ and $\mu_N = \sqrt{N_2/N_1} \to \mu$, with $\mu \in [0, 1]$.

Assumption A.2. The unobservable factor process $F_{\tau} = [f_{\tau}^{c\prime}, f_{1,\tau}^{s\prime}, f_{2,\tau}^{s\prime}]'$ has vector of means, and covariance matrix as defined in (A.6), that is:

$$E[F_{ au}] = \left[egin{array}{c} \mu^c \ \mu^s_1 \ \mu^s_2 \end{array}
ight], \qquad ext{and} \qquad \Sigma_F := extbf{Var}(F_{ au}) = \left[egin{array}{ccc} I_{k^c} & 0 & 0 \ 0 & I_{k^s_1} & \Upsilon \ 0 & \Upsilon' & I_{k^s_2} \end{array}
ight],$$

with all the elements of vector $\mathbb{E}[F_{\tau}]$ being finite, and where Σ_F is positive-definite.

Assumption A.3. The loadings matrix $\Lambda_{1,\tau} = [\Lambda_{1,\tau}^c : \Lambda_{1,\tau}^s] = [\lambda_{1,\tau,1}, \ldots, \lambda_{1,\tau,N_j}]'$ is such that $\lim_{N_1 \to \infty} \frac{1}{N_{1,\tau}} \Lambda'_{1,\tau} \Lambda_{1,\tau} = \Sigma_{\lambda,1,\tau}$, where $\Sigma_{\lambda,1,\tau}$ is a positive-definite (k_1, k_1) matrix with distinct eigenvalues and $k_1 = k^c + k_1^s$, for all $\tau = 1, \ldots, T$. Analogously, the loadings matrix $\Lambda_2 = [\Lambda_2^c : \Lambda_2^s] = [\lambda_{2,1}, \ldots, \lambda_{2,N_2}]'$ is such that $\lim_{N_2 \to \infty} \frac{1}{N_2} \Lambda'_2 \Lambda_2 = \Sigma_{\lambda,2}$, where $\Sigma_{\lambda,2}$ is a positive-definite (k_2, k_2) matrix with distinct eigenvalues and $k_2 = k^c + k_2^s$.

Assumption A.4. The error terms $\varepsilon_{j,i,\tau}$ and the factors $h_{j,\tau} = [f_{\tau}^{c\prime}, f_{j,\tau}^{s\prime}]'$ are such that for j = 1, 2 and all $i, \tau \geq 1$: a) $\mathbb{E}[\varepsilon_{j,i,\tau} | \mathcal{F}_{\tau}^{ipc}] = 0$ and $\mathbb{E}[\varepsilon_{j,i,\tau}^2 | \mathcal{F}_{\tau}^{ipc}] \leq M$, a.s., where $\mathcal{F}_{\tau}^{ipc} = \sigma(F_s, s \leq \tau)$,

b) $\mathbb{E}[\varepsilon_{j,i,\tau}^8] \leq M$ and $\mathbb{E}[\|h_{j,\tau}\|^{2r\vee 8}] \leq M$, for a constant $M < \infty$, and r > 2. Moreover, the conditions in (A.10) hold.

Assumption A.5. Define the variables $\xi_{1,\tau}^{ipc} = \frac{1}{\sqrt{N_1}} \sum_{i=1}^{N_1} \lambda_{1,i,\tau} \varepsilon_{1,i,\tau}$ and $\xi_{2,\tau} = \frac{1}{\sqrt{N_2}} \sum_{i=1}^{N_2} \lambda_{1,i} \varepsilon_{2,i,\tau}$ a) For any $\tau \geq 1$ and $h \geq 0$ have:

$$[\xi_{1,\tau}^{ipc\prime}, \xi_{2,\tau}']' \xrightarrow{d} N(0, \Omega_{\tau}^{ipc}), \quad (\mathcal{F}_{\tau}^{ipc}\text{-stably}),$$

as $N_1, N_2 \to \infty$, where the asymptotic variance matrix is:

$$\Omega_{\tau}^{ipc}(h) = \begin{bmatrix} \Omega_{11,\tau}^{ipc} & \Omega_{12,\tau}^{ipc} \\ & & \\ & \Omega_{22,\tau}^{ipc} \end{bmatrix},$$

where $\Omega_{11,\tau}^{ipc} = \lim_{N_1 \to \infty} \frac{1}{N_1} \sum_{i=1}^{N_1} \sum_{\ell=1}^{N_1} \lambda_{1,i,\tau} \lambda'_{1,\ell,\tau} cov(\varepsilon_{1,i,\tau}, \varepsilon_{1,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}),$

$$\Omega^{ipc}_{22, au} = \lim_{N_2 o \infty} rac{1}{N_2} \sum_{i=1}^{N_2} \sum_{\ell=1}^{N_2} \lambda_{2,i} \lambda'_{2,\ell} cov(arepsilon_{2,i, au}, arepsilon_{2,\ell, au} | \mathcal{F}^{ipc}_{ au}),$$
 and

$$\Omega_{12,\tau}^{ipc} = \underset{N_1N_2 \to \infty}{\text{plim}} \frac{1}{\sqrt{N_1N_2}} \sum_{i=1}^{N_1} \sum_{\ell=1}^{N_2} \lambda_{1,i,\tau} \lambda_{2,\ell}' cov(\varepsilon_{1,i,\tau}, \varepsilon_{2,\ell,\tau} | \mathcal{F}_{\tau}^{ipc}), \text{ for all } h.$$

Moreover, for all $N_1, N_2 \ge 1$ and j = 1, 2, we have: b) $\mathbb{E}(\|\xi_{j,\tau}\|^{2r}|\mathcal{F}_{\tau}^{ipc}) \le M$, and r > 2.

Assumption A.6. a) The triangular array process $V_{\tau} \equiv V_{N_1,N_2,\tau} = [h_{j,\tau},\xi'_{j,\tau},j=1,2]'$ is strong mixing of size $-\frac{r}{r-2}$, uniformly in $N_1,N_2 \geq 1$, with $\xi_{1,\tau} \equiv \xi^{ipc}_{1,\tau}$. Moreover,

b)
$$\|\mathbb{E}(\xi_{j,\tau}\xi'_{k,\tau}|\mathcal{F}_{\tau}^{ipc}) - \mathbb{E}(\xi_{j,\tau}\xi'_{k,\tau}|F_{\tau},...,F_{\tau-m};\tilde{Z}_{\tau-1}^{1:N},...,\tilde{Z}_{\tau-m-1}^{1:N})\|_{2} = O(m^{-\psi})$$
, as $m \to \infty$, uniformly in $N_{1}, N_{2} \geq 1$.

Assumption A.7. The errors $\varepsilon_{j,i,\tau}$ are i) uncorrelated at all leads and lags, both within and across

groups, conditionally on $\mathcal{F}_{ au}^{ipc}$, that is,

$$Cov(\varepsilon_{j,i,\tau},\varepsilon_{k,\ell,\tau-h}|\mathcal{F}_{\tau}^{ipc})=0, \quad when \ h\neq 0, \ and for \ all \ j, \ k, \ i, \ \tau.$$
 (OA.14)

The contemporaneous conditional variance-covariance matrix of the stacked error terms $\varepsilon_{\tau} = [\varepsilon'_{1,\tau}, \varepsilon'_{2,\tau}]'$, denoted as $\mathbb{V}(\varepsilon_{\tau} | \mathcal{F}_{\tau}^{ipc})$, where $\varepsilon_{\tau} = [\varepsilon'_{1,\tau}, \varepsilon'_{2,\tau}]'$ is sparse and time-invariant.

OA.4.2 Proof of Proposition A.1

From the covariance matrix Σ_F of the factor vector $(f_{\tau}^c, f_{1,\tau}^s, f_{2,\tau}^s)'$ in equation (A.6), and the definition of matrices R and R^* given in Section A.2, it follows that:

$$R = \begin{pmatrix} I_{k^c} & 0 \\ 0 & \Upsilon \Upsilon' \end{pmatrix}, \qquad R^* = \begin{pmatrix} I_{k^c} & 0 \\ 0 & \Upsilon' \Upsilon \end{pmatrix}. \tag{OA.15}$$

Noting that also in our set-up matrix Σ_F is assumed to be positive definite (see Assumption A.2), then the proof of Proposition A.1 is omitted as it is analogous to the proof of Proposition 1 in AGGR (see Section C.1 in their OA).

OA.4.3 Proof of Theorem A.1

OA.4.3.1 Asymptotic expansion of $\sum_{\ell=1}^{k^c} \hat{\rho}_{\ell}$ for IPCA

The next Lemma OA.4.1, which corresponds to Lemma B.5 in AGGR modified to take into account that latent factors in the first group, namely $h_{1,\tau}$ are estimated by IPCA, while the latent factors on the second group are estimated by (RP-) PCA. The Lemma provides the asymptotic expansion for

the sum of the k^c largest canonical correlations $\sum_{\ell=1}^{k^c} \hat{\rho}_{\ell}$, which are (square root of) the largest k^c eigenvalues of matrix \hat{R} in equation (A.7), computed with $\hat{h}_{1,\tau}$ being the IPCA estimator of $h_{1,\tau}$, and $\hat{h}_{2,\tau}$ being the (RP-) PCA estimator of $h_{2,\tau}$.

LEMMA OA.4.1. Under the Assumptions in Section OA.4.1, and other technical assumptions on the higher order cross-moments of the errors $\varepsilon_{j,i,\tau}$ and factors in F_{τ} we have:

$$\begin{split} \sum_{\ell=1}^{k^c} \hat{\rho}_{\ell} &= k^c - \frac{1}{2N} tr \left\{ \tilde{\Sigma}_{cc}^{-1} \frac{1}{T} \sum_{\tau=1}^{T} E[(\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)}) (\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)})' | \mathcal{F}_{\tau}^{ipc}] \right\} \\ &- \frac{1}{2N\sqrt{T}} tr \left\{ \frac{1}{\sqrt{T}} \sum_{\tau=1}^{T} \left[(\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)}) (\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)})' - E[(\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)}) (\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)})' | \mathcal{F}_{\tau}^{ipc}] \right] \right\} \\ &+ o_p \left(\epsilon_{N,T} \right), \end{split}$$

where $\epsilon_{N,T} := \frac{1}{N\sqrt{T}}$. The terms in the curly brackets are $O_p(1)$.

Recalling the definitions in (A.15) -(A.17), we get

$$\frac{1}{T} \sum_{\tau=1}^{T} E[(\mu_N u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)})(\mu_N u_{1\tau}^{ipc(c)} - u_{2t}^{(c)})' | \mathcal{F}_{\tau}^{ipc}] \quad = \quad \frac{1}{T} \sum_{\tau=1}^{T} \left(\mu_N^2 \tilde{\Sigma}_{u,11,\tau}^{ipc(cc)} + \tilde{\Sigma}_{u,22,\tau}^{ipc(cc)} - \mu_N \tilde{\Sigma}_{u,12,\tau}^{ipc(cc)} - \mu_N \tilde{\Sigma}_{u,21,\tau}^{ipc(cc)} \right),$$

where the last term is equal to $\tilde{\Sigma}_U^{ipc}$ defined in Theorem A.1. Moreover, let us define the process

$$U_{\tau} := \mu_{N} u_{1\tau}^{ipc(c)} - u_{2\tau}^{(c)}$$

$$= \mu_{N} \left[\left(\frac{1}{N_{1}} \sum_{i=1}^{N_{1}} \lambda_{1,i,\tau} \lambda'_{1,i,\tau} \right)^{-1} \frac{1}{\sqrt{N_{1}}} \sum_{i=1}^{N_{1}} \lambda_{1,i,\tau} \varepsilon_{1,i,\tau} \right]^{(c)} - \left[\left(\frac{1}{N_{2}} \sum_{i=1}^{N_{2}} \lambda_{2,i} \lambda'_{2,i} \right)^{-1} \frac{1}{\sqrt{N_{2}}} \sum_{i=1}^{N_{2}} \lambda_{2,i} \varepsilon_{2,i,\tau} \right]^{(c)}$$

$$= \mu_{N} \left[\left(\frac{1}{N_{2}} \sum_{i=1}^{N_{1}} \lambda_{1,i,\tau} \lambda'_{1,i,\tau} \right)^{-1} \xi_{1,\tau}^{ipc} \right]^{(c)} - \left[\left(\frac{1}{N_{2}} \sum_{i=1}^{N_{2}} \lambda_{2,i} \lambda'_{2,i} \right)^{-1} \xi_{2,\tau} \right]^{(c)} .$$
(OA.16)

Process U_{τ} depends on N_1 , N_2 , but we do not make this dependence explicit for expository purpose. By using these definitions, from Lemma OA.4.1 we get:

$$\sum_{\ell=1}^{k^{c}} \hat{\rho}_{\ell} - k^{c} + \frac{1}{2N} tr \left\{ \tilde{\Sigma}_{cc}^{-1} \tilde{\Sigma}_{U}^{ipc} \right\} = -\frac{1}{2N\sqrt{T}} \left(\frac{1}{\sqrt{T}} \sum_{\tau=1}^{T} \left[U_{\tau}' U_{\tau} - E(U_{\tau}' U_{\tau} | \mathcal{F}_{\tau}^{ipc}) \right] \right) + o_{p} \left(\epsilon_{N,T} \right). \quad (OA.18)$$

Under our set of assumptions the term $\frac{1}{\sqrt{T}}\sum_{\tau=1}^{T}\left[U_{\tau}'U_{\tau}-E(U_{\tau}'U_{\tau}|\mathcal{F}_{\tau}^{ipc})\right]$ is $O_{p}(1)$, as in the next subsection we show that it is asymptotically Gaussian distributed. The remainder term $o_{p}\left(\epsilon_{N,T}\right)$ in the r.h.s. of equation (OA.18) is negligible with respect to the first term in the r.h.s. The result in equation (OA.18) is analogous to the one in equation (B.15) in AGGR, with the notable differences being the new definitions of the term in U_{τ} provided in our equation (OA.16), and of matrix $\tilde{\Sigma}_{U}$.

OA.4.3.2 Asymptotic distribution of the test statistic under the null hypothesis $H(k^c)$

From the asymptotic expansion (OA.18) we obtain the asymptotic distribution of $\hat{\xi}(k^c) = \sum_{\ell=1}^{k^c} \hat{\rho}_{\ell}$ under the null hypothesis $H(k^c)$ of k^c common factors. First, we apply a CLT for weakly dependent triangular array data to prove the asymptotic normality of $\frac{1}{\sqrt{T}} \sum_{\tau=1}^{T} \mathcal{Z}_{N,\tau}$ as $N, T \to \infty$, where

$$\mathcal{Z}_{N,\tau} := U_{\tau}' U_{\tau} - E(U_{\tau}' U_{\tau} | \mathcal{F}_{\tau}^{ipc})$$

depends on N_1, N_2 via process U_{τ} defined in (OA.16).

i) CLT for Near-Epoch Dependent (NED) processes

Consider the assumptions A.5 and A.6, and let process $V_{N_1,N_2,\tau} \equiv V_{\tau}$ be as defined in Assumption A.6, and let $\mathcal{V}_{\tau-m}^{\tau+m} = \sigma(V_s,\tau-m \leq s \leq \tau+m)$ for any positive integer m, with $\mathcal{V}_{\tau} \equiv \mathcal{V}_{-\infty}^{\tau}$.

LEMMA OA.4.2. Under the Assumptions in Section OA.4.1 and other technical assumptions on the higher order cross-moments of the errors $\varepsilon_{j,i,\tau}$ and factors in F_{τ} we have:

(i)
$$\mathcal{Z}_{N,\tau}$$
 is measurable w.r.t. \mathcal{V}_{τ} , and $\mathbb{E}[\mathcal{Z}_{N,\tau}] = 0$ for all $\tau \geq 1$ and $N_1, N_2 \geq 1$,

(ii)
$$\sup_{\tau \geq 1, N_1, N_2 \geq 1} \mathbb{E}\left[\|\mathcal{Z}_{N,\tau}\|^r\right] < \infty$$
, for a constant $r > 2$,

- (iii) Process $(\mathcal{Z}_{N,\tau})$ is L^2 Near Epoch Dependent $(L^2\text{-NED})$ of size -1 on process (V_τ) , and (V_τ) is strong mixing of size -r/(r-2), uniformly in $N_1,N_2\geq 1$.
- (iv) Matrix $\Omega_U := \lim_{T,N \to \infty} \mathbb{V}\left(\frac{1}{\sqrt{T}} \sum_{\tau=1}^T \mathcal{Z}_{N,\tau}\right)$ is positive definite and equal to

$$\Omega_{U} = \lim_{T,N\to\infty} \frac{1}{T} \sum_{\tau=1}^{T} \sum_{s=1}^{T} Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,s}\right)$$
(OA.19)

$$= \lim_{T,N\to\infty} \left(\frac{1}{T} \sum_{\tau=1}^{T} \sum_{h=-(T-t)}^{\tau-1} Cov\left(\mathcal{Z}_{N,\tau}, \mathcal{Z}_{N,\tau-h}\right) \right)$$
(OA.20)

$$= \lim_{T,N\to\infty} \left(\frac{1}{T} \sum_{\tau=1}^{T} Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,\tau}\right) \right). \tag{OA.21}$$

Then, by an application of the univariate CLT in Corollary 24.7 in Davidson (1994) and the Cramér-Wold device, we have that:

$$\frac{1}{\sqrt{T}} \sum_{\tau=1}^{T} \mathcal{Z}_{N,\tau} \stackrel{d}{\longrightarrow} N\left(0, \Omega_{U}\right), \tag{OA.22}$$

as $T, N \to \infty$.

Let us now compute the limit (auto)covariance matrix $Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,\tau}\right)$ explicitly. By the Law of Iterated Expectation and $\mathbb{E}[\mathcal{Z}_{N,\tau}|\mathcal{F}_{\tau}^{ipc}]=0$, we have:

$$Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,\tau}\right) = \lim_{N \to \infty} \mathbb{E}\left[Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,\tau}|\mathcal{F}_{\tau}^{ipc}\right)\right]. \tag{OA.23}$$

Moreover, from Assumptions A.3 and A.5 a), vector $(U'_{\tau}, U'_{\tau})'$ is asymptotically Gaussian for

any h, τ as $N \to \infty$:

$$\begin{pmatrix} U_{\tau} \\ U_{\tau} \end{pmatrix} \xrightarrow{d} \begin{pmatrix} U_{\tau}^{\infty} \\ U_{\tau}^{\infty} \end{pmatrix} \sim N \begin{pmatrix} 0, \begin{bmatrix} \Sigma_{U,\tau}^{ipc} & \Sigma_{U,\tau}^{ipc} \\ \Sigma_{U,\tau}^{ipc} & \Sigma_{U,\tau}^{ipc} \end{bmatrix} \end{pmatrix}, \quad (\mathcal{F}_{\tau}^{ipc}\text{-stably}). \tag{OA.24}$$

We use the Lebesgue Lemma to interchange the limes for $N \to \infty$ and the outer expectation in the r.h.s. of (OA.23), and the fact that convergence in distribution plus uniform integrability imply convergence of the expectation for a sequence of random variables (see Theorem 25.12 in Billingsley (1995)) to show the next lemma.

LEMMA OA.4.3. *Under Assumptions A.3 and A.5 b), we have:*

$$Cov\left(\mathcal{Z}_{N,\tau},\mathcal{Z}_{N,\tau}\right) = \mathbb{E}\left[Cov(U_{\tau}^{\infty}{}'U_{\tau}^{\infty},U_{\tau}^{\infty}{}'U_{\tau}^{\infty}|\mathcal{F}_{\tau}^{ipc})\right].$$

Lemma OA.4.3 allows to deploy the joint asymptotic Gaussian distribution of $(U_{\tau}^{\infty}{}', U_{\tau}^{\infty}{}')'$ to compute the limit autocovariance matrices $Cov(\mathcal{Z}_{N,\tau}, \mathcal{Z}_{N,\tau})$: we do so by using Theorem 12 p. 284 in Magnus and Neudecker (2007) and Theorem 10.21 in Schott (2005). We get

$$Cov(U_{\tau}^{\infty} 'U_{\tau}^{\infty}, U_{\tau}^{\infty} 'U_{\tau}^{\infty} | \mathcal{F}_{\tau}^{ipc}) = 2tr \left\{ \Sigma_{U,\tau}^{ipc} \Sigma_{U,\tau}^{ipc} \right\}$$
 (OA.25)

Therefore from (OA.19) and Lemma OA.4.3 we get:

$$\Omega_{U} = 2 \cdot \left(\lim_{T \to \infty} \frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \mathbb{E} \left[\Sigma_{U,\tau}^{ipc} \Sigma_{U,\tau}^{ipc\prime} \right] \right\} \right) = 4\Omega_{U,1}^{ipc}$$
 (OA.26)

ii) Asymptotic Gaussian distribution of the test statistic

Let us define the constant $D_{N,T}=\frac{1}{2N\sqrt{T}}$. From equations (OA.18) and (OA.26), and by using: $[D_{N,T}^2\Omega_U^{ipc}]^{1/2}=\frac{1}{N\sqrt{T}}\Omega_{U,1}^{1/2}$, and $N\sqrt{T}[\Omega_{U,1}^{ipc}]^{-1/2}=O\left(N\sqrt{T}\right)=O(\epsilon_{N,T}^{-1})$, under the hypothesis of k^c common factors in each group the statistics $\hat{\xi}(k^c)=\sum_{\ell=1}^{k^c}\hat{\rho}_\ell$ is such that:

$$N\sqrt{T}(\Omega_{U,1}^{ipc})^{-1/2} \left[\hat{\xi}(k^c) - k^c + \frac{1}{2N} tr \left\{ \tilde{\Sigma}_{cc}^{-1} \tilde{\Sigma}_{U}^{ipc} \right\} \right] = -(D_{N,T}^2 \Omega_{U}^{ipc})^{-1/2} D_{N,T} \frac{1}{\sqrt{T}} \sum_{\tau=1}^{T} \mathcal{Z}_{N,T} + o_p(1).$$

From equation (OA.22), the r.h.s. converges in distribution to a standard normal distribution, which yields Theorem A.1.

OA.4.4 Proof of Theorem A.2

To establish the asymptotic distribution of the feasible statistic in Theorem A.2 we need to control the effect of replacing the re-centering and scaling terms by means of their estimates.

Proof of Theorem A.2 Part (i)

Let us first consider the asymptotic distribution of $\tilde{\xi}_{inf}^{ipc}(k^c)$ under the null hypothesis of k^c common factors.

Theorem A.2 i) follows, if we prove:

$$tr\left\{\hat{\Sigma}_{U}^{ipc}\right\} = tr\left\{\tilde{\Sigma}_{cc}^{-1}\tilde{\Sigma}_{U}^{ipc}\right\} + o_{p}\left(\frac{1}{\sqrt{T}}\right),$$
 (OA.27)

$$\frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \hat{\Sigma}_{U,\tau}^{ipc} \cdot \hat{\Sigma}_{U,\tau}^{ipc}' \right\} = E \left[\frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \Sigma_{U,\tau}^{ipc} \cdot \Sigma_{U,\tau}^{ipc'} \right\} \right] + o_p(1). \tag{OA.28}$$

Indeed, the statistic $\tilde{\xi}^{ipc}_{inf}(k^c)$ can be rewritten as:

$$\tilde{\xi}^{ipc}(k^c) = \left[\hat{\Omega}_{U,1}^{ipc} / \Omega_{U,1}^{ipc} \right]^{-1/2} \left\{ N \sqrt{T} (\Omega_{U,1}^{ipc})^{-1/2} \left[\hat{\xi}(k^c) - k^c + \frac{1}{2N} tr \left\{ \tilde{\Sigma}_{cc}^{-1} \tilde{\Sigma}_{U}^{ipc} \right\} \right] + O_p \left(\sqrt{T} \left[tr \left\{ \hat{\Sigma}_{U}^{ipc} \right\} - tr \left\{ \tilde{\Sigma}_{cc}^{-1} \tilde{\Sigma}_{U}^{ipc} \right\} \right] \right) \right\},$$

where the ratio $\hat{\Omega}_{U,1}^{ipc}/\Omega_{U,1}^{ipc}$ converges in probability to 1 from (OA.28), the term within the curly brackets in the first line in the r.h.s. converges in distribution to a standard normal distribution from (A.18), and the term on the second line on the r.h.s. is $o_p(1)$ from (OA.27).

Let us now prove equations (OA.27) and (OA.28) by deriving the asymptotic expansions of $\hat{\Sigma}_{U}^{ipc}$, $\tilde{\Sigma}_{cc}^{-1}$ and $\hat{\Omega}_{U,1}^{ipc}$. To derive the asymptotic expansion of $\hat{\Sigma}_{U}^{ipc}$, we use its definition

$$\hat{\Sigma}_{U}^{ipc} \ := \ \frac{1}{T} \sum_{\tau=1}^{T} \hat{\Sigma}_{U,\tau}^{ipc}, \qquad \text{and} \qquad \hat{\Omega}_{U,1}^{ipc} \ := \ \frac{1}{2} \frac{1}{T} \sum_{\tau=1}^{T} tr \left\{ \hat{\Sigma}_{U,\tau}^{ipc} \cdot \hat{\Sigma}_{U,\tau}^{ipc} \right\},$$

with

$$\hat{\Sigma}_{U,\tau}^{ipc} \ := \ \hat{\mu}_N^2 \hat{\Sigma}_{u,11,\tau}^{ipc(cc)} + \hat{\Sigma}_{u,22,\tau}^{ipc(cc)} - \hat{\mu}_N \hat{\Sigma}_{u,12,\tau}^{ipc(cc)} - \hat{\mu}_N \hat{\Sigma}_{u,21,\tau}^{ipc(cc)},$$

where the matrices $\hat{\Sigma}_{u,ij,\tau}$ i,j=1,2, defined in equation (A.19) - (A.22), involve the estimated loadings and residuals.

LEMMA OA.4.4. Under the Assumptions in Section OA.4.1, and other technical assumptions on the higher order cross-moments of the errors $\varepsilon_{j,i,\tau}$ and factors in F_{τ} i) the following asymptotic

expansion hold:

$$\frac{\hat{\Lambda}'_{1,\tau}\hat{\Lambda}_{1,\tau}}{N_1} = \hat{\mathcal{U}}'_1 \tilde{\Sigma}_{\Lambda,1,\tau} \hat{\mathcal{U}}_1 + o_p \left(\frac{1}{\sqrt{T}}\right), \tag{OA.29}$$

$$\frac{\hat{\Lambda}_2'\hat{\Lambda}_2}{N_2} = \hat{\mathcal{U}}_2' \left[\tilde{\Sigma}_{\Lambda,2} + \frac{1}{\sqrt{T}} \left(L_{\Lambda,2} + L_{\Lambda,2}' \right) \right] \hat{\mathcal{U}}_2 + o_p \left(\frac{1}{\sqrt{T}} \right), \tag{OA.30}$$

where

$$\hat{\mathcal{U}}_j = \left[egin{array}{cc} \hat{\mathcal{H}}_c & 0 \\ 0 & \hat{\mathcal{H}}_{s,j} \end{array}
ight], \qquad j=1,2$$

and $\hat{\mathcal{H}}_c$, $\hat{\mathcal{H}}_{s,j}$ are non-singular matrices w.p.a. 1.,

$$\tilde{\Sigma}_{\Lambda,1,\tau} = \frac{1}{N_1} \Lambda'_{1,\tau} \Lambda_{1,\tau}, \qquad \tilde{\Sigma}_{\Lambda,2} = \frac{1}{N_2} \Lambda'_2 \Lambda_2$$

with $\Lambda_{1,\tau}=[\Lambda_{1,\tau}^c\ \vdots\ \Lambda_{1,\tau}^s]$,, and $\Lambda_2=[\Lambda_2^c\ \vdots\ \Lambda_2^s]$, $L_{\Lambda,2}=\tilde{\Sigma}_{\Lambda,2}Q_2$ with

$$Q_2 = \begin{bmatrix} 0 & 0 \\ \sqrt{T}\tilde{\Sigma}_{2,c}\tilde{\Sigma}_{cc}^{-1} & 0 \end{bmatrix}.$$

ii) Under the same assumptions, we also have

$$\frac{1}{T} \sum_{\tau=1}^{T} \hat{\Sigma}_{u,11,\tau}^{ipc(cc)} = \frac{1}{T} \sum_{\tau=1}^{T} \left[\left(\frac{\hat{\Lambda}'_{1,\tau} \hat{\Lambda}_{1,\tau}}{N_{1}} \right)^{-1} \frac{\hat{\Lambda}'_{1,\tau} \left[\mathcal{I}^{\mathcal{T}}_{\varepsilon,11} (c^{*}) \odot \check{\varepsilon}_{1,\tau} \check{\varepsilon}'_{1,\tau} \right] \hat{\Lambda}_{1,\tau}}{N_{1}} \left(\frac{\hat{\Lambda}'_{1,\tau} \hat{\Lambda}_{1,\tau}}{N_{1}} \right)^{-1} \right]^{(cc)}$$

$$= \frac{1}{T} \sum_{\tau=1}^{T} \tilde{\Sigma}_{u,11,\tau}^{ipc(cc)} + o_{p} \left(\frac{1}{\sqrt{T}} \right)$$

$$\frac{1}{T} \sum_{\tau=1}^{T} \hat{\Sigma}_{u,22,\tau}^{ipc(cc)} = \frac{1}{T} \sum_{\tau=1}^{T} \left[\left(\frac{\hat{\Lambda}_{2}' \hat{\Lambda}_{2}}{N_{2}} \right)^{-1} \frac{\hat{\Lambda}_{2}' \left[\mathcal{I}_{\varepsilon,22}^{T}(c^{*}) \odot \check{\varepsilon}_{2,\tau} \check{\varepsilon}_{2,\tau}' \right] \hat{\Lambda}_{2}}{N_{2}} \left(\frac{\hat{\Lambda}_{2}' \hat{\Lambda}_{2}}{N_{2}} \right)^{-1} \right]^{(cc)}$$

$$= \frac{1}{T} \sum_{\tau=1}^{T} \tilde{\Sigma}_{u,22,\tau}^{ipc(cc)} + o_{p} \left(\frac{1}{\sqrt{T}} \right)$$

and

$$\frac{1}{T} \sum_{\tau=1}^{T} \hat{\Sigma}_{u,12,\tau}^{ipc(cc)} = \frac{1}{T} \sum_{\tau=1}^{T} \left[\left(\frac{\hat{\Lambda}'_{1,\tau} \hat{\Lambda}_{1,\tau}}{N_1} \right)^{-1} \frac{\hat{\Lambda}'_{1,\tau} \left[\mathcal{I}^{\mathcal{T}}_{\varepsilon,12} (c^*) \odot \breve{\varepsilon}_{1,\tau} \breve{\varepsilon}'_{2,\tau} \right] \hat{\Lambda}_2}{\sqrt{N_1 N_2}} \left(\frac{\hat{\Lambda}'_{2} \hat{\Lambda}_{2}}{N_2} \right)^{-1} \right]^{(cc)}$$

$$= \frac{1}{T} \sum_{\tau=1}^{T} \tilde{\Sigma}_{u,12,\tau}^{ipc(cc)} + o_p \left(\frac{1}{\sqrt{T}} \right),$$

where \odot denotes the Hadamard product (i.e. the element-wise product) among two matrices of the same dimensions, $c^* \in (0,1)$ is a trimming constant, $\check{\varepsilon}_{j,i,\tau} := [\check{\varepsilon}_{j,1,\tau},...,\check{\varepsilon}_{j,N_j,\tau}]'$, with $\check{\varepsilon}_{j,i,\tau} := \hat{\varepsilon}_{j,i,\tau} - \bar{\hat{\varepsilon}}_{j,i,\tau}$, and $\bar{\hat{\varepsilon}}_{j,i,\tau} = \sum_{\tau=1}^T \hat{\varepsilon}_{j,i,\tau}/T$, for j=1,2. Moreover, $\mathcal{I}_{\varepsilon,jk}^T(c^*)$ is a $N_j \times N_k$ matrix defined such that its generic element in position (i,ℓ) is $[\mathcal{I}_{\varepsilon,jk}^T(c^*)]_{i,\ell} = \mathbf{1}\{c\hat{orr}_T(\hat{\varepsilon}_{j,i,\tau},\hat{\varepsilon}_{k,\ell,\tau}) > c^*\}$ with $i=1,...,N_j$, and $\ell=1,...,N_k$, and j,k=1,2.

Therefore, from Lemma OA.4.4 and the definition of $\hat{\Sigma}_U^{ipc}$ from Theorem A.1 :

$$\hat{\Sigma}_{U}^{ipc} := \frac{1}{T} \sum_{\tau=1}^{T} \left(\hat{\mu}_{N}^{2} \hat{\Sigma}_{u,11,\tau}^{ipc(cc)} + \hat{\Sigma}_{u,22,\tau}^{ipc(cc)} - \hat{\mu}_{N} \hat{\Sigma}_{u,12,\tau}^{ipc(cc)} - \hat{\mu}_{N} \hat{\Sigma}_{u,21,\tau}^{ipc(cc)} \right)$$

we get the asymptotic expansion:

$$\hat{\Sigma}_{U}^{ipc} = \hat{\mathcal{H}}_{c}^{-1} \tilde{\Sigma}_{U}^{ipc} \left(\hat{\mathcal{H}}_{c}^{\prime} \right)^{-1} + o_{p} \left(\frac{1}{\sqrt{T}} \right). \tag{OA.31}$$

Moreover, adapting the arguments in the proof of Theorem 2 part (i) in AGGR it can be shown

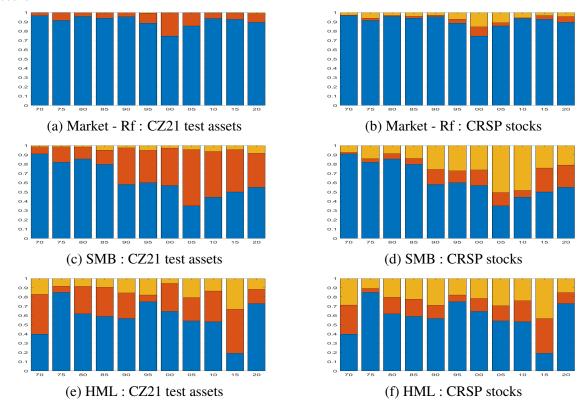
that $\tilde{\Sigma}_{cc}^{-1} = \left(\hat{\mathcal{H}}_c^{-1}\right)'\hat{\mathcal{H}}_c^{-1} + o_p\left(\frac{1}{\sqrt{T}}\right)$. The last equation, together with the asymptotic expansion (OA.31) and the commutative property of the trace operator, imply equation (OA.27). Similarly, the asymptotic expansion (OA.31) and the convergence $\tilde{\Sigma}_{u,jk,\tau}^{ipc} \to \Sigma_{u,jk,\tau}$ for all j,k=1,2 and τ , imply equation (OA.28).

Proof of Theorem A.2 Part (ii)

The proof of Part (ii) of Theorem A.2 uses the same arguments as the Proof of part (ii) of Theorem 2 in AGGR, adapted to new assumptions in Section OA.4.1 and therefore is omitted.

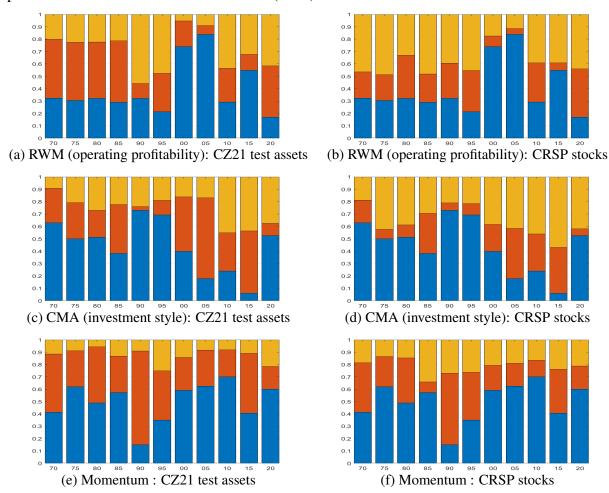
OA.5 Supplementary empirical results

Figure OA.2: Variability of the Fama and French 3 factors explained by common and specific factors



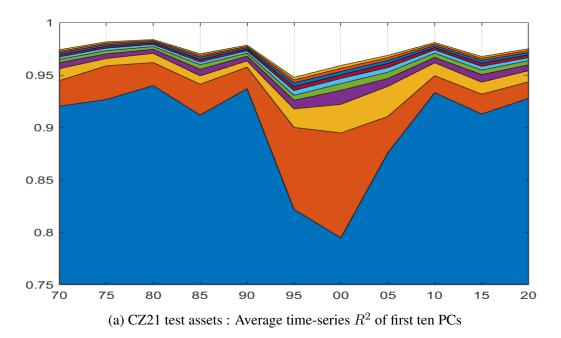
For each of the FF3 factors the figure displays the fraction of variance (R^2) explained by the three common factors (blue bars which are the same in left and right column panels), seven CZ21's group-specific factors (orange bars, left column panels), seven CRSP group-specific factors (orange bars, right column panels), and unexplained by common and group-specific factors (yellow bars). For each year we report results based on the block starting in year t-4 and ending in year t, for each t=1970, ..., 2020.

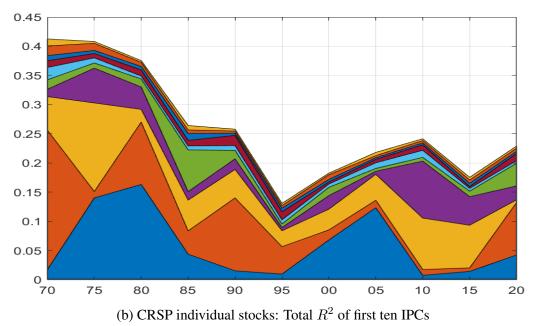
Figure OA.3: Variability of the FF factors RMW, CMA and Momentum explained by common and specific factors in Chen and Zimmermann (2022) test assets and CRSP individual stocks.



For the two Fama and French factors related to Operating Profitability and Investment style, and Momentum, the figure displays the fraction of variance (R^2) explained by the three common factors between CRSP and CZ21 test assets (blue bars which are the same in both panels), CZ21's group-specific factors (orange bars, left panels), CRSP group-specific factors (orange bars, right panels), and unexplained by common and group-specific factors (yellow bars). For each year we report results based on the block starting in year t-4 and ending in year t, for each t=1970, ..., 2020.

Figure OA.4: Explanatory power of first 10 PCs for Chen and Zimmermann (2022) test assets and IPCs from CRSP individual stocks, full sample: 1966-2020





Panel (a) displays the average fraction of variance (time-series regression R^2) of the individuals on the balanced panel of CZ21 test assets explained by the first 10 RP-PCs extracted from the same panel. The bottom blue area in each panel represents the average R^2 of the first RP-PC, the second (from the bottom) orange area represents the average R^2 of the second RP-PC, and so on. Lettau and Pelger's RP-PCs are computed (fixing $\gamma_{RP}=-1$) on balanced panel of portfolios. In every year each panel of assets is constructed for the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020. In every 5-years rolling window we only include assets with non-missing returns for all the 60 months. Panel (b) displays the fraction of (in-sample) the Total R^2 , computed as described in Section OA.6 of the unbalanced panel of returns of individual stocks explained by the first 10 IPCs extracted by estimating the IPCA model from the same unbalanced panel starting in year t-4 and ending in year t, for each t=1970,...,2020. The bottom blue area in each panel represents the Total R^2 of the first IPC, the second (from the bottom) orange area represents the Total R^2 of the second IPC, and so on.

OA.5.1 RP-PCs estimation of factors for the CZ21 test assets

Table OA.2: In- and Out-of-sample Total R^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP}

	lı	n-Samp	le		Out-of-Sample				
1	3	4	5	6	1	3	4	5	6
$\gamma_{RP} =$	+1								
	23.5					25.7			
	93.4					91.3			
		24.2	24.8	25.4			26.8	27.7	28.5
		94.3	94.7	95.2			92.1	92.4	92.7
24.5	31.9	34.2	36.6	38.7	18.7	18.5	17.8	17.8	17.0
89.2	94.5	95.1	95.7	96.0	87.3	91.8	92.6	92.9	93.1
$\gamma_{RP} =$	+5								
	23.4					25.7			
	93.4					91.3			
		24.1	24.6	25.3			26.8	27.7	28.4
		94.2	94.7	95.1			92.1	92.4	92.7
24.5	31.8	34.1	36.2	38.4	18.7	18.8	17.7	17.8	17.0
89.1	94.3	95.0	95.6	95.9	87.3	91.7	92.5	92.9	93.1
$\gamma_{RP} =$	+10								
	23.5					25.7			
	93.4					91.3			
		24.1	24.6	25.2			26.8	27.7	28.5
		94.2	94.7	95.1			92.1	92.4	92.7
24.4	31.6	34.0	36.1	38.4	18.8	18.9	17.7	17.8	17.0
89.0	94.1	95.0	95.6	95.9	87.3	91.7	92.5	92.9	93.1
rison									
21.5	29.4	32.2	33.1	35.8	16.3	15.9	13.1	11.9	9.4
73.9	89.8	91.5	90.7	92.3	70.8	83.6	84.7	83.6	84.8
5.4	17.8	22.2	24.2	25.1	5.3	18.7	22.1	24.8	27.1
21.4	65.6	80.3	86.5	88.8	9.0	52.9	60.7	68.9	73.4
	$\gamma_{RP} =$ 24.5 89.2 $\gamma_{RP} =$ 24.5 89.1 $\gamma_{RP} =$ 24.4 89.0 arison 21.5 73.9 5.4	$\gamma_{RP} = +1$ 23.5 93.4 24.5 31.9 89.2 94.5 $\gamma_{RP} = +5$ 23.4 93.4 24.5 31.8 89.1 94.3 $\gamma_{RP} = +10$ 23.5 93.4 24.4 31.6 89.0 94.1 arison 21.5 29.4 73.9 89.8 5.4 17.8	$\gamma_{RP} = +1$ 23.5 93.4 24.2 94.3 24.5 31.9 34.2 89.2 94.5 95.1	$\gamma_{RP} = +1$ 23.5 93.4 24.2 24.8 94.3 94.7 24.5 31.9 34.2 36.6 89.2 94.5 95.1 95.7 $\gamma_{RP} = +5$ 23.4 93.4 24.1 24.6 94.2 94.7 24.5 31.8 34.1 36.2 89.1 94.3 95.0 95.6 $\gamma_{RP} = +10$ 23.5 93.4 24.1 24.6 94.2 94.7 24.4 31.6 34.0 36.1 89.0 94.1 95.0 95.6 rrison 21.5 29.4 32.2 33.1 73.9 89.8 91.5 90.7 5.4 17.8 22.2 24.2	$\gamma_{RP} = +1$ 23.5 93.4 24.2 24.8 25.4 94.3 94.7 95.2 24.5 31.9 34.2 36.6 38.7 89.2 94.5 95.1 95.7 96.0 $\gamma_{RP} = +5$ 23.4 93.4 24.1 24.6 25.3 94.2 94.7 95.1 24.5 31.8 34.1 36.2 38.4 89.1 94.3 95.0 95.6 95.9 $\gamma_{RP} = +10$ 23.5 93.4 24.1 24.6 25.2 94.2 94.7 95.1 24.4 31.6 34.0 36.1 38.4 89.0 94.1 95.0 95.6 95.9 Prison 21.5 29.4 32.2 33.1 35.8 73.9 89.8 91.5 90.7 92.3 5.4 17.8 22.2 24.2 25.1	$\gamma_{RP} = +1$ 23.5 93.4 24.2 24.8 25.4 94.3 94.7 95.2 24.5 31.9 34.2 36.6 38.7 89.2 94.5 95.1 95.7 96.0 87.3 $\gamma_{RP} = +5$ 23.4 93.4 24.1 24.6 25.3 94.2 94.7 95.1 24.5 31.8 34.1 36.2 38.4 89.1 94.3 95.0 95.6 95.9 87.3 $\gamma_{RP} = +10$ 23.5 93.4 24.1 24.6 25.2 24.7	$\gamma_{RP} = +1$ 23.5 93.4 24.2 24.8 24.7 94.3 94.7 95.2 24.5 31.9 34.2 36.6 38.7 89.2 94.5 95.1 95.7 96.0 87.3 91.8	$\gamma_{RP} = +1$ 23.5 93.4 24.2 24.8 94.3 94.7 95.2 24.5 31.9 34.2 36.6 38.7 89.2 94.5 95.1 95.7 96.0 37.3 91.8 92.6 37.4 93.4 24.1 24.6 25.7 91.3 24.1 24.6 25.7 91.3 24.1 24.6 25.7 91.3 24.1 24.6 25.7 91.3 24.1 24.6 25.7 91.3 24.1 24.6 25.7 91.3 26.8 94.2 94.7 95.1 24.5 31.8 34.1 36.2 38.4 89.1 94.3 95.0 95.6 95.9 37.3 91.7 92.5 38.4 91.3 25.7 92.1 24.5 31.8 34.1 36.2 38.4 89.1 94.3 95.0 95.6 95.9 25.7 91.3 26.8 92.1 26.8 94.2 94.7 95.1 26.8 91.3 92.5 26.8 94.2 94.7 95.1 24.1 24.6 25.2 94.2 94.7 95.1 26.8 91.3 91.7 92.5 26.8 94.2 94.7 95.1 92.1 24.4 31.6 34.0 36.1 38.4 89.0 94.1 95.0 95.6 95.9 87.3 91.7 92.5 26.8 92.1 25.7 92.1 26.8 91.3 92.5 70.8 83.6 84.7 73.9 89.8 91.5 90.7 92.3 70.8 83.6 84.7 5.4 17.8 22.2 24.2 25.1 5.3 16.3 15.9 13.1 73.9 89.8 91.5 90.7 92.3 70.8 83.6 84.7 5.4 17.8 22.2 24.2 25.1 5.3 18.7 22.1	$\gamma_{RP} = +1$ $\begin{array}{cccccccccccccccccccccccccccccccccccc$

Panels A, B and C of the table report Total R^2 s in percent for models involving estimates of the latent factors in the CZ21 test assets by means of RP-PCA with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. Models with observable factors and IPCs estimated on individual stocks form CRSP are reported in Panel D to ease the comparison. In each of the first three panels, we report results for a latent factor model with 3 factors common between individual stocks and CZ21 portfolios (lines 1-2 in each panel), the same 3 common factors together with 1, 2, or 3 CRSP-specific factors (line 3 in each panel), again the same 3 common factors together with 1, 2, or 3 CZ21-specific factors (line 4 in each panel), and a latent factor model where the factors are K RP-PCs extracted from the CZ21 portfolios only (lines 5-6 in each panel). Observable factor model specifications are CAPM, FF3, FF3 + Momentum, FF5, and FF5 + Momentum in the K=1,3,4,5,6 columns, respectively. The models (IPCA on individual stocks and PCs on CZ21 portfolios, and group-factor model based on the previous two models) are estimated on non-overlapping windows starting in year t-4 and ending in year t, for each t=1970,...,2020. The R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

Table OA.3: In- and Out-of-sample Pricing \mathbb{R}^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP}

Pricing R^2		Iı	n-Samp	le						
N. of factors, K	1	3	4	5	6	1	3	4	5	6
Panel A: RP-PCA on CZ21 with	$\gamma_{RP} =$	+1								
r: CRSP, f : Comm.		44.6					35.1			
r: CZ21, f: Comm.		91.9					92.7			
r: CRSP, f : Comm. + CRSP spec.			49.3	50.0	48.4			36.7	38.3	37.7
r: CZ21, f : Comm. + CZ21 spec.			95.0	95.4	96.3			93.6	94.4	95.0
r: CRSP, f : RP-PCA on CZ21	43.4	60.0	65.5	66.2	70.5	8.9	6.6	5.1	6.4	4.5
r: CZ21, f: RP-PCA on CZ21	88.7	95.4	96.6	97.2	97.6	90.7	93.7	94.7	95.5	95.6
Panel B: RP-PCA on CZ21 with	$\gamma_{RP} =$	+5								
r: CRSP, f : Comm.		44.2					35.2			
r: CZ21, f: Comm.		91.7					92.7			
r: CRSP, f : Comm. + CRSP spec.			47.1	48.9	46.9			36.7	38.9	37.3
r: CZ21, f: Comm. + CZ21 spec.			94.7	95.3	96.0			93.6	94.4	94.7
r: CRSP, f: RP-PCA on CZ21	43.1	65.3	67.3	67.5	75.0	8.8	7.4	4.9	6.5	4.5
r: CZ21, f: RP-PCA on CZ21	88.2	96.2	97.0	97.3	97.5	90.8	93.5	95.1	95.7	95.9
Panel C: RP-PCA on CZ21 with	$\gamma_{RP} =$	+10								
r: CRSP, f : Comm.		44.3					35.1			
r: CZ21, f: Comm.		91.7					92.7			
r: CRSP, f : Comm. + CRSP spec.			47.0	48.4	47.5			36.7	38.8	37.6
r: CZ21, f : Comm. + CZ21 spec.			94.7	95.2	95.8			93.6	94.4	94.7
r: CRSP, f: RP-PCA on CZ21	42.9	66.9	68.2	69.2	76.1	8.8	7.8	4.9	6.5	4.5
r: CZ21, f : RP-PCA on CZ21	87.7	96.6	97.0	97.1	97.4	90.8	93.5	95.2	95.8	95.9
Panel D: other factors for compa	rison									
r: CRSP, f : FF + mom	32.5	51.0	56.3	54.3	59.1	7.5	8.1	6.5	4.2	1.2
r: CZ21, f: FF + mom	75.2	89.0	88.4	90.0	89.7	78.8	85.5	84.4	86.7	86.4
r: CRSP, f : IPCA on CRSP	17.1	33.9	40.8	45.8	45.8	6.9	33.3	35.6	39.9	42.4
r: CZ21, f : IPCA on CRSP	66.1	94.0	91.6	94.7	95.5	32.0	81.1	85.1	91.9	93.5

Panels A, B and C of the table report Pricing R^2 s in percent for models involving estimates of the latent factors in the CZ21 test assets by means of RP-PCA with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. Models with observable factors and IPCs estimated on individual stocks form CRSP are reported in Panel D to ease the comparison. In each of the first three panels, we report results for a latent factor model with 3 factors common between individual stocks and CZ21 portfolios (lines 1-2 in each panel), the same 3 common factors together with 1, 2, or 3 CRSP-specific factors (line 3 in each panel), again the same 3 common factors together with 1, 2, or 3 CZ21-specific factors (line 4 in each panel), and a latent factor model where the factors are K RP-PCs extracted from the CZ21 portfolios only (lines 5-6 in each panel). Observable factor model specifications are CAPM, FF3, FF3 + Momentum, FF5, and FF5 + Momentum in the K=1,3,4,5,6 columns, respectively. The models (IPCA on individual stocks and PCs on CZ21 portfolios, and group-factor model based on the previous two models) are estimated on non-overlapping windows starting in year t-4 and ending in year t, for each t=1970,...,2020. The R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

Table OA.4: In- and Out-of-sample Predictive R^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP}

Predictive R^2		Iı	ı-Samp	le			Out	-of-Sar	nple	
N. of factors, K	1	3	4	5	6	1	3	4	5	6
Panel A: RP-PCA on CZ21 with	$\gamma_{RP} =$	+1								
r: CRSP, f : Comm.		0.75					0.79			
r: CZ21, f: Comm.		4.06					4.73			
r: CRSP, f : Comm. + CRSP spec.			0.80	0.89	0.90			0.80	0.82	0.83
r: CZ21, f : Comm. + CZ21 spec.			4.22	4.15	4.15			4.75	4.76	4.87
r: CRSP, f : RP-PCA on CZ21	0.77	1.10	1.21	1.26	1.32	0.46	0.34	0.17	0.10	0.11
r: CZ21, f: RP-PCA on CZ21	3.99	4.31	4.33	4.37	4.37	4.47	4.66	4.65	4.69	4.71
Panel B: RP-PCA on CZ21 with	$\gamma_{RP} =$	+5								
r: CRSP, f : Comm.		0.76					0.78			
r: CZ21, f: Comm.		4.03					4.71			
r: CRSP, f : Comm. + CRSP spec.			0.80	0.84	0.85			0.80	0.78	0.84
r: CZ21, f: Comm. + CZ21 spec.			4.18	4.11	4.09			4.73	4.75	4.83
r: CRSP, f: RP-PCA on CZ21	0.77	1.21	1.25	1.29	1.42	0.45	0.24	<0	<0	0.05
r: CZ21, f : RP-PCA on CZ21	3.96	4.20	4.31	4.32	4.33	4.46	4.48	4.62	4.66	4.69
Panel C: RP-PCA on CZ21 with	$\gamma_{RP} =$	+10								
r: CRSP, f : Comm.		0.77					0.78			
r: CZ21, f: Comm.		4.02					4.71			
r: CRSP, f : Comm. + CRSP spec.			0.81	0.83	0.83			0.79	0.78	0.84
r: CZ21, f: Comm. + CZ21 spec.			4.17	4.10	4.08			4.73	4.75	4.83
r: CRSP, f: RP-PCA on CZ21	0.77	1.23	1.28	1.33	1.45	0.44	0.11	<0	<0	0.04
r: CZ21, f: RP-PCA on CZ21	3.94	4.19	4.30	4.30	4.31	4.44	4.43	4.61	4.65	4.69
Panel D: other factors for compa	rison									
r: CRSP, f : FF + mom	0.51	0.89	1.03	0.96	1.07	0.35	0.40	0.36	0.17	0.17
r: CZ21, f: FF + mom	2.60	4.00	4.03	4.07	4.12	2.63	3.92	4.09	3.91	4.11
r: CRSP, f: IPCA on CRSP	0.39	0.85	0.95	1.05	1.02	0.36	0.79	0.94	0.97	0.94
r: CZ21, f : IPCA on CRSP	1.90	2.54	3.79	4.15	4.22	1.94	3.75	4.30	4.36	4.14

Panels A, B and C of the table report Predictive R^2 s in percent for models involving estimates of the latent factors in the CZ21 test assets by means of RP-PCA with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. Models with observable factors and IPCs estimated on individual stocks form CRSP are reported in Panel D to ease the comparison. In each of the first three panels, we report results for a latent factor model with 3 factors common between individual stocks and CZ21 portfolios (lines 1-2 in each panel), the same 3 common factors together with 1, 2, or 3 CRSP-specific factors (line 3 in each panel), again the same 3 common factors together with 1, 2, or 3 CZ21-specific factors (line 4 in each panel), and a latent factor model where the factors are K RP-PCs extracted from the CZ21 portfolios only (lines 5-6 in each panel). Observable factor model specifications are CAPM, FF3, FF3 + Momentum, FF5, and FF5 + Momentum in the K=1, 3, 4, 5, 6 columns, respectively. The models (IPCA on individual stocks and PCs on CZ21 portfolios, and group-factor model based on the previous two models) are estimated on non-overlapping windows starting in year t-4 and ending in year t, for each t=1970, ..., 2020. The R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

Table OA.5: Significance of the CZ21-specific factors in Fama-MacBeth regressions, RP-PCs on CZ21 test assets with $\gamma_{RP}=+1,+5$, and +10

Sample	66-70	71-75	76-80	81-85	86-90	90-95	96-00	01-05	06-10	11-15	16-20		
Panel A :	Panel A : $\gamma_{RP}=+1$												
Wald stat. p-val.	7.727 0.052	3.707 0.295	1.504 0.681	15.053 0.002	3.483 0.323	8.611 0.035	0.117 0.990	8.329 0.040	4.806 0.187	0.769 0.857	1.777 0.620		
Panel B :	Panel B : $\gamma_{RP} = +5$												
Wald stat. p-val.	8.498 0.037	6.986 0.072	1.202 0.753	15.449 0.001	1.735 0.629	3.596 0.309	0.049 0.997	5.376 0.146	4.587 0.205	0.733 0.865	1.801 0.615		
Panel C :	$\gamma_{RP} = -$	-10											
Wald stat. p-val.	8.005 0.046	6.585 0.086	1.231 0.746	19.415 0.000	2.417 0.491	3.130 0.372	0.069 0.995	4.718 0.194	3.680 0.298	0.799 0.850	1.028 0.795		

The table reports the values of the Wald test statistics and associated p-values for the joint significance of the risk premia of the three CZ21-specific factors estimated by cross-sectional regressions of the risk premia of the CZ21 portfolios on the betas of 3CF, and the three CZ21-specific factors considered in Table OA.2, i.e. those obtained by first performing RP-PCA with different values of the RP-PCA tuning parameter set to +1 (Panel A), +5 (Panel B) and +10 (Panel C). Cross-sectional regressions are estimated by Weighted Least Squares, where each cross-sectional observation is scaled by the the square root of the pricing error from the regression itself. The variance covariance matrix of estimated factors' risk premia is computed using the Fama-MacBeth procedure, with 3-lags (optimally selected) Newey-West weighting. P-values are obtained using the asymptotic distribution of the Wald test for the joint significance of the three coefficients, namely a χ^2 with 3 degrees of freedom. The model (RP-PCs on CZ21 portfolios, and group-factor model) and cross-sectional WLS regressions (of risk premia of CZ21 quantile portfolios on factors' betas) are estimated on non-overlapping windows of 60 months starting in year t-4 and ending in year t, for each t=1970, ..., 2020.

Table OA.6: Out-of-sample Sharpe ratios of factor portfolios, RP-PCs on CZ21 test assets with $\gamma_{RP}=+1,+5$ and +10

N. of factors, K	1	3	4	5	6								
Panel A: RP-PCA	on CZ2	1 with -	$\gamma_{RP} = +1$										
Comm.		0.66											
Comm. CZ21 spec.			0.89	1.20	1.26								
Comm. CRSP spec.			0.14	0.25	0.87								
RP-PCA on CZ21	0.49	0.62	1.03	1.32	1.02								
Panel B : RP-PCA on CZ21 with $\gamma_{RP}=+5$													
Comm.		0.66											
Comm. CZ21 spec.			0.91	1.18	1.22								
Comm. CRSP spec.			0.15	0.24	0.91								
RP-PCA on CZ21	0.49	0.51	1.17	1.44	1.15								
Panel C: RP-PCA	on CZ2	1 with \sim	$\gamma_{RP} = +10$										
Comm.		0.65											
Comm. CZ21 spec.			0.92	1.18	1.21								
Comm. CRSP spec.			0.16	0.25	0.91								
RP-PCA on CZ21	0.49	0.46	1.22	1.47	1.19								
Panel D: other fact	ors for	compar	rison										
FF + mom	0.39	0.30	0.62	0.63	0.81								
IPCA on CRSP	0.31	0.77	1.09	1.10	1.44								

The table reports out-of-sample annualized Sharpe ratios for the mean-variance efficient portfolio of factors in each model involving RP-PCs estimated with different values of the RP-PCA tuning parameter γ_{RP} , i.e. +1 (Panel A), +5 (Panel B) and +10 (Panel C). Models with observable factors and IPCs estimated on individual stocks form CRSP are reported in Panel D to ease the comparison. See caption Table OA.2 for further details.

Table OA.7: In- and Out-of-sample Total R^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP} . Common and group specific factors versus zoo

Total \mathbb{R}^2		Ir	ı-Samp	le			Out	t-of-Saı	nple	
N. of factors, K	1	3	4	5	6	1	3	4	5	6
Panel A : $\gamma_{RP} = +1$										
r: CRSP, f : Comm.		23.5					25.7			
r: CZ21, f: Comm.		93.4					91.3			
r: CRSP, f : RP-PCA on Zoo	14.9	25.6	29.3	32.2	34.5	6.0	0.3	1.4	0.8	1.1
r: CZ21, f: RP-PCA on Zoo	40.1	68.5	74.9	78.3	79.8	25.8	24.3	31.1	32.6	35.6
Panel B : $\gamma_{RP} = +5$										
r: CRSP, f : Comm.		23.4					25.7			
r: CZ21, f: Comm.		93.4					91.3			
r: CRSP, f : RP-PCA on Zoo	12.1	24.3	28.9	32.0	34.3	3.4	<0	<0	<0	<0
r: CZ21, f: RP-PCA on Zoo	27.4	67.2	75.4	78.3	79.9	12.7	<0	5.1	6.6	4.4
Panel C : $\gamma_{RP} = +10$										
r: CRSP, f : Comm.		23.5					25.7			
r: CZ21, f: Comm.		93.4					91.3			
r: CRSP, f: RP-PCA on Zoo	8.7	24.0	28.9	32.0	34.3	2.7	<0	<0	<0	<0
r: CZ21, f: RP-PCA on Zoo	21.3	67.3	75.5	78.3	79.9	8.5	<0	<0	<0	<0

Panels A, B and C of the table report Total R^2 s in percent for models involving estimates of the 3CF only (first two lines in each panel), and latent factors from the factors in the Zoo by means of K RP-PCs (last two lines in each panel) with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. The models are estimated on the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020. R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

Table OA.8: In- and Out-of-sample Pricing \mathbb{R}^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP} . Common and group specific factors versus zoo

Pricing R^2		Ir	ı-Samp	le		Out-of-Sample						
N. of factors, K	1	3	4	5	6	1	3	4	5	6		
Panel A : $\gamma_{RP} = +1$												
r: CRSP, f: Comm.		44.6					35.1					
r: CZ21, f: Comm.		91.9					92.7					
r: CRSP, f :RP- PCA on Zoo	<0	14.6	43.7	67.5	68.1	5.8	<0	<0	<0	<0		
r: CZ21, f :RP- PCA on Zoo	<0	21.5	71.1	83.5	82.9	<0	34.0	26.6	26.0	33.3		
Panel B : $\gamma_{RP} = +5$												
r: CRSP, f : Comm.		44.2					35.2					
r: CZ21, f: Comm.		91.7					92.7					
r: CRSP, f :RP- PCA on Zoo	<0	42.7	63.8	74.9	74.8	4.9	<0	<0	<0	<0		
r: CZ21, f :RP- PCA on Zoo	2.3	59.2	85.1	88.4	87.5	<0	40.5	31.8	29.2	33.5		
Panel C : $\gamma_{RP} = +10$												
r: CRSP, f : Comm.		44.3					35.1					
r: CZ21, f: Comm.		91.7					92.7					
r: CRSP, f:RP- PCA on Zoo	<0	56.1	71.3	76.2	76.1	5.0	<0	<0	<0	<0		
r: CZ21, f :RP- PCA on Zoo	31.4	68.2	88.2	89.2	88.3	0.2	50.9	43.5	41.6	44.6		

Panels A, B and C of the table report Pricing R^2 s in percent for models involving estimates of the 3CF only (first two lines in each panel), and latent factors from the factors in the Zoo by means of K RP-PCs (last two lines in each panel) with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. The models are estimated on the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020. R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

Table OA.9: In- and Out-of-sample Predictive R^2 of factor models, RP-PCs on CZ21 test assets with different values of γ_{RP} . Common and group specific factors versus zoo

Predictive \mathbb{R}^2		I	n-Samj	ple		Out-of-Sample					
N. of factors, K	1	3	4	5	6	1	3	4	5	6	
Panel A : $\gamma_{RP} = +1$											
r: CRSP, f : Comm.		0.75					0.79				
r: CZ21, f: Comm.		4.06					4.73				
r: CRSP, f :RP- PCA on Zoo	<0	0.22	0.87	1.26	1.28	<0	<0	<0	<0	<0	
r: CZ21, f:RP- PCA on Zoo	<0	1.18	3.37	3.93	3.94	0.02	0.50	1.70	2.31	2.35	
Panel B : $\gamma_{RP} = +5$											
r: CRSP, f : Comm.		0.76					0.78				
r: CZ21, f: Comm.		4.03					4.71				
r: CRSP, f :RP- PCA on Zoo	<0	0.75	1.24	1.41	1.41	<0	<0	<0	<0	<0	
r: CZ21, f:RP- PCA on Zoo	<0	2.78	3.97	4.15	4.14	<0	0.25	0.47	1.07	0.99	
Panel C : $\gamma_{RP} = +10$											
r: CRSP, f : Comm.		0.77					0.78				
r: CZ21, f: Comm.		4.02					4.71				
r: CRSP, f :RP- PCA on Zoo	<0	1.01	1.36	1.44	1.44	<0	<0	<0	<0	<0	
r: CZ21, f :RP- PCA on Zoo	<0	3.17	4.13	4.19	4.17	<0	0.26	0.52	1.14	1.04	

Panels A, B and C of the table report Predictive R^2 s in percent for models involving estimates of the 3CF only (first two lines in each panel), and latent factors from the factors in the Zoo by means of K RP-PCs (last two lines in each panel) with tuning parameter γ_{RP} set to +1, +5 and +10, respectively. The models are estimated on the rolling window starting in year t-4 and ending in year t, for each t=1970,...,2020. R^2 's in-sample (left table) and out-of-sample (right table) are computed either for the excess returns of individual stocks (r: CRSP) or CZ21 portfolios (r: CZ21) as described in Section B.

OA.6 Performance evaluation measures

We describe the various performance evaluation measures both in-sample and out-of-sample starting with the former.

Let $y_{j,i,\tau}$ be the excess return in month τ belonging to block b of the i-th asset in group j, with j=1 corresponding to individual stocks, and j=2 to CZ21 test asset portfolios. Each model m for the returns of the CZ21 test assets $y_{2,i,\tau}$ can be expressed as

$$y_{2,i,\tau} = \beta_{2,i,b}^{m'} f_{2,\tau}^m + \varepsilon_{2,i,\tau}^m , \quad \text{with} \quad \tau \in b ,$$
 (OA.32)

where $f_{j,\tau}^m = [f_{\tau}^{m,c\prime}, f_{\tau}^{j,m,s\prime}]'$ for j=1,2, $\beta_{2,i,b}^m = [\lambda_{2,i,b}^{m,c\prime}, \lambda_{2,i,b}^{m,s\prime}]'$, while $f_{\tau}^{m,c}$ and $\lambda_{2,i,b}^{m,c\prime}$ (resp. $f_{2,\tau}^{m,s}$ and $\lambda_{2,i,b}^{m,s\prime}$) are the common (resp. group-specific) factors and betas/loadings. We also use the same model with constant loadings (within each block b) also for the individual stocks when the factors are observable. Similarly, each IPCA model m for the individual stocks, with the exclusion of the models with observable factors, can be written as:

$$y_{1,i,\tau} = (z'_{i,\tau-1}\Gamma^m_{\beta,i,b}) f^m_{1,\tau} + \varepsilon^m_{1,i,\tau}, \quad \text{with} \quad \tau \in b.$$
 (OA.33)

where $z_{i,\tau-1}$ is the vector collecting the values of the 36 characteristics (35 time-varying and a constant) of stock i in month $\tau-1$, and $\Gamma^m_{\beta,i,b}$ is a $36\times(k^c+k^s_1)$ matrix of constant coefficients. Therefore, $\Gamma^{m\prime}_{\beta,i,b}z_{i,\tau-1}$ is the $(k^c+k^s_1)$ - dimensional vector collecting the time-varying (with each block b) individual stocks' factor loadings implied by the IPCA model.

OA.6.1 In-sample performance evaluation

Let $N_{j,b}$ be the total number of assets for which the full sample of returns is available in group j and block b, with j=1,2 and b=1,...,B. For each model m and for each group of assets j, we compute the following six performance measures across the entire sample, that is across all B blocks:

1. Total R^2 of Kelly et al. (2019), which for our model with betas changing across blocks can be expressed as:

$$Tot. R_j^2(m) = 1 - \frac{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \sum_{\tau \in b} \left(y_{j,i,\tau} - \hat{\beta}_{j,i,b,\tau}^m f_{\tau}^m \right)^2}{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \sum_{\tau \in b} y_{j,i,\tau}^2}.$$

It represents the fraction of return variance for all the assets present in group j explained by both the dynamic behavior of the loadings across different blocks, as well as by the contemporaneous factor realizations, aggregated over all assets and all time periods, that is across all B blocks. The $Total\ R^2$ summarizes how well the systematic factor in a given model specification describes the realized riskiness in the panel of individual stocks. In the case of observable factors, i.e. models in (i) in Section A, the coefficients $\beta^m_{j,i,b,\tau}$ are constant across all the dates τ in block b, and are estimated by an OLS regression without intercept of excess returns on factors, compatibly with model (OA.32). By construction, the $\beta^m_{j,i,b,\tau}$ and factors for all other models (ii) - (vi) in Section A are also estimated by PCA, or variation of it, compatible with a linear model without intercept. For all models not involving IPCA, the loadings $\beta^m_{j,i,b,\tau}$ are constant across all the dates τ in block b. ¹⁰

Note that in Kelly et al. (2019) the $\hat{\beta}_{i,\tau}$ can effectively be computed only using information up to time τ , as they are a function of the stock-specific characteristics z_{it} observed at time t. Kelly, Palhares, and Pruitt (2020) notes

2. Predictive \mathbb{R}^2 from Kelly et al. (2019), which for our model with betas changing across blocks can be expressed as:

Pred.
$$R_j^2(m) = 1 - \frac{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \sum_{\tau \in b} \left(y_{j,i,\tau} - \hat{\beta}_{j,i,b,\tau}^m \bar{f}_b^m \right)^2}{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \sum_{\tau \in b} y_{j,i,\tau}^2}$$
.

where $\bar{f}_b^m = \frac{1}{T_b} \sum_{\tau \in b} f_\tau^m$ is the sample average of the factors' realizations within all the T_b dates in block b only, that is the same block in which the $\beta_{j,i,\tau}^m$ are estimated compatibly with a model without intercept. *Predictive* R^2 represents the fraction of realized return variation explained by the model's description of conditional expected returns, and summarizes the model's ability to describe risk compensation only through exposure to systematic risk. Our measure of the *Predictive* R^2 is slightly different from the in-sample *Predictive* R^2 of Kelly et al. (2019) as ours allows for factor risk premia which vary across different blocks, while theirs imposes constant risk premia across dates.¹¹

3. Pricing error R^2 of Kelly et al. (2020), which is defined as:

$$Pr.Err. R_j^2(m) = 1 - \frac{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \left[\frac{1}{T_b} \sum_{\tau \in b} (y_{j,i,\tau} - \hat{\beta}_{j,i,b,\tau}^m)^T f_{\tau}^m) \right]^2}{\sum_{b=1}^B \sum_{i=1}^{N_{j,b}} \left(\frac{1}{T_b} \sum_{\tau \in b} y_{j,i,\tau} \right)^2},$$

that is the fraction of the squared unconditional mean excess returns that is described by factors and betas. In contrast to the previous two R^2 measures, this focuses on whether the model's fitted values do a good job of explaining assets' average returns. This metric is close in flavor to formal statistical tests (like the GRS test) of whether or not a cross section of test

that: "the factor realization estimate \hat{f}_{t+1} is constructed from a combination of time t+1 realized returns, time τ instruments $z_{i,\tau}$, and $\hat{\Gamma}_{\beta}$ estimated from data through time τ , according to equation."

¹¹Due to the rotational indeterminacy of factors which are re-estimated across different blocks, we cannot impose a constant factor average across different blocks.

assets' pricing errors are zero.

OA.6.2 Out-of-sample performance evaluation

We implement the out-of-sample version of the $Total\ R^2$, $Pricing\ R^2$ and $Predicitve\ R^2$ where betas and factor loadings, needed to reconstruct the latent factors out-of-sample for date τ in block b are computed using information from the previous block b-1. Analogously to Lettau and Pelger (2020b) we also compute the annualized Sharpe Ratio of the "Maximum Sharpe-ratio portfolio" that can be obtained by an optimal (in a mean-variance sense) linear combination of the factors, which are ultimately portfolios of individual stocks. Our out-of-sample performance measures are defined as:

1. OOS Total \mathbb{R}^2 , which for our model with betas changing across blocks can be expressed as:

OOS Tot.
$$R_j^2(m) = 1 - \frac{\sum_{b=2}^B \sum_{i=1}^{N_{j,b-1}} \sum_{\tau \in b} \left(y_{j,i,\tau} - \hat{\beta}_{j,i,b-1,\tau}^{m'} f_{\tau|b-1}^m \right)^2}{\sum_{b=2}^B \sum_{i=1}^{N_{j,b-1}} \sum_{\tau \in b} y_{j,i,\tau}^2}$$
.

The beta coefficients $\hat{\beta}^m_{j,i,b-1,\tau}$ are estimated using information available in block b-1 only (this also explains why we use only the $N_{j,b-1}$ stocks or portfolios available at the end of the block b-1), while returns $y_{j,i,\tau}$ are observed for the same stocks at dates τ in block b. For models with observable factors, $f^m_{\tau|b-1}$ is simply the observed value of the factor at date τ in block b, as all our observable factors are returns of portfolios of individual stocks observed at date τ with weights computed at date $\tau-1$. Instead, when a model includes latent factors, we compute their values at date τ by running cross-sectional regressions of the returns $y_{j,i,\tau}$ for all assets available both in the previous block b-1, and in the current one b, on the factor loadings estimated in the previous block b-1 only.

For instance, model (v) in Section A implies that in block b the DGP for the returns of the CZ21 portfolios is:

$$y_{2,\tau} = \Lambda_{2,b-1}^m f_{\tau|b-1}^m + \varepsilon_{2,\tau}^m$$
, with $\tau \in b$. (OA.34)

Let $\hat{\Lambda}_{2,b-1}^m$ be the (RP-)PC estimator of matrix $\Lambda_{2,b-1}$ obtained using the returns of all assets in group j for dates $\tau \in b-1$. Then, compatibly with model (OA.34), factors $f_{\tau|b-1}^m$ are computed as $f_{\tau|b-1}^m = (\hat{\Lambda}_{2,b-1}^m \hat{\Lambda}_{2,b-1}^m)^{-1} \hat{\Lambda}_{2,b-1}^m y_{2,\tau}$, for all dates $\tau \in b$. Analogously, for the models involving IPCs the loadings (and therefore the factors) are computed using the values of the characteristics Z_{τ} in block b (appropriately lagged), but with the values of matrix C linking the loadings to characteristics estimated in block b-1.

2. OOS Pricing R^2 , which can be expressed as:

OOS Pr.Err.
$$R_j^2(m) = 1 - \frac{\sum_{b=2}^B \sum_{i=1}^{N_{j,b-1}} \left[\frac{1}{T_b} \sum_{\tau \in b} (y_{j,i,\tau} - \hat{\beta}_{j,i,b-1,\tau}^m f_{\tau|b-1}^m) \right]^2}{\sum_{b=1}^B \sum_{i=1}^{N_{j,b-1}} \left(\frac{1}{T_b} \sum_{\tau \in b} y_{j,i,\tau} \right)^2},$$

where all quantities are computed as described for the $OOS\ Total\ R^2$.

3. OOS Predictive R^2 , which can be expressed as:

OOS Pred.
$$R_j^2(m) = 1 - \frac{\sum_{b=2}^B \sum_{i=1}^{N_{j,b-1}} \sum_{\tau \in b} \left(y_{j,i,\tau} - \hat{\beta}_{j,i,b,\tau}^m \bar{f}_{\tau|b-1,\tau-1}^m \right)^2}{\sum_{b=2}^B \sum_{i=1}^{N_{j,b-1}} \sum_{\tau \in b} y_{j,i,\tau}^2}$$
.

where $\bar{f}_{\tau|b-1,\tau-1}^m$ is the sample average of the factor realizations computed over the 60 months ending at date $\tau-1$, and where the factor is reconstructed (if necessary) for each date as

described for the computation of the OOS Total \mathbb{R}^2 , that is regressing returns in each month τ on loadings estimated in the previous block b-1.

4. Maximum Sharpe-ratio, Max. SR, that is the realized Sharpe Ratio of a portfolio of "factors" (returns) of each model $f^m_{\tau|b-1}$ combined at each date τ in block b with weights $w^m_{f,b} = (\hat{\Sigma}^m_{f,b-1})^{-1}\hat{\mu}^m_{f,b-1}$, where $\hat{\mu}^m_{f,b-1}$ and $\hat{\Sigma}^m_{f,b-1}$ are the sample mean and covariance, respectively, of all factors in model m computed using their observations in block b-1. Therefore, both factors and their weights in the Maximum Sharpe Ratio portfolio in block b are computed using the factor loadings estimated in block b-1.

References

BAI, J. AND S. NG (2002): "Determining the Number of Factors in Approximate Factor Models," *Econometrica*, 70, 191–221.

BILLINGSLEY, P. (1995): Probability and Measure, Wiley.

Breeden, D. T., M. R. Gibbons, and R. H. Litzenberger (1989): "Empirical Test of the Consumption-Oriented CAPM," *Journal of Finance*, 2, 231–262.

CHEN, A. Y. AND T. ZIMMERMANN (2022): "Open Source Cross-Sectional Asset Pricing," Critical Finance Review, 11, 207–264.

CONNOR, G. AND R. KORAJCZYK (1988): "Risk and Return in an Equilibrium APT: Application of a New Test Methodology," *Journal of Financial Economics*, 21, 255–289.

DAVIDSON, J. (1994): Stochastic Limit Theory, Oxford University Press.

- FAMA, E. AND M. GIBBONS (1984): "A comparison of inflation forecasts," *Journa, of Monetary Economics*, 13, 327–348.
- FREYBERGER, J., A. NEUHIERL, AND M. WEBER (2020): "Dissecting Characteristics Nonparametrically," *Review of Financial Studies*, 33, 2326–2377.
- GIGLIO, S. AND D. XIU (2021): "Asset Pricing with Omitted Factors," *Journal of Political Economy*, 129, 1947–1990.
- GREEN, J., J. R. M. HAND, AND X. F. ZHANG (2017): "The Characteristics that Provide Independent Information about Average US Monthly Stock Returns," *Review of Financial Studies*, 30, 4389–4436.
- HARVEY, C. R., Y. LIU, AND H. ZHU (2016): "... and the Cross-Section of Expected Returns," *Review of Financial Studies*, 29, 5–68.
- HOU, K., C. XUE, AND L. ZHANG (2020): "Replicating Anomalies," *Review of Financial Studies*, 35, 2019–2133.
- HUBERMAN, G., S. KANDEL, AND R. F. STAMBAUGH (1987): "Mimicking Portfolios and Exact Arbitrage Pricing," *Journal of Finance*, 42, 1–9.
- JURADO, K., S. C. LUDVIGSON, AND S. NG (2015): "Measuring Uncertainty," *American Economic Review*, 105, 1177–1216.
- KELLY, B. T., D. PALHARES, AND S. PRUITT (2020): "Modeling Corporate Bond Returns," Working Paper.

- KELLY, B. T., S. PRUITT, AND Y. Su (2019): "Characteristics are Covariances: A Unified Model of Risk and Return," *Journal of Financial Economics*, 134, 501–524.
- KIM, S. AND R. KORAJCZYK (2021): "Large Sample Estimators of the Stochastic Discount Factor," *Working paper*.
- KOZAK, S., S. NAGEL, AND S. SANTOSH (2018): "Interpreting Factor Models," *Journal of Finance*, 73, 1183–1223.
- ——— (2020): "Shrinking the Cross-Section," *Journal of Financial Economics*, 135, 271–292.
- LEHMANN, B. N. AND D. M. MODEST (2005): "Diversification and the Optimal Construction of Basis Portfolios," *Management Science*, 51, 581–598.
- LETTAU, M. AND M. PELGER (2020a): "Estimating Latent Asset-Pricing Factors," *Journal of Econometrics*, 218, 1–31.
- ——— (2020b): "Factors That Fit the Time Series and Cross-Section of Stock Returns," *Review of Financial Studies*, 33, 2274–2325.
- LUDVIGSON, S. AND S. NG (2009): "Macro factors in bond risk premia," *Review of Financial Studies*, 22, 5027–5067.
- MAGNUS, J. R. AND H. NEUDECKER (2007): Matrix Differential Calculus with Applications in Statistics and Econometrics, John Wiley and Sons: Chichester/New York.
- MCCRACKEN, M. AND S. NG (2016): "FRED-MD: A Monthly Database for Macroeconomic Research," *Journal of Business and Economic Statistics*, 36, 574–589.

- MCLEAN, R. D. AND J. PONTIFF (2016): "Does Academic Research Destroy Stock Return Predictability?" *Journal of Finance*, 71, 5–32.
- PUKTHUANTHONG, K., R. ROLL, AND A. SUBRAHMANYAM (2019): "A Protocol for Factor Identification," *Review of Financial Studies*, 32, 1573–1607.
- Ross, S. A. (1976): "The Arbitrage Theory of Capital Asset Pricing," *Journal of Economic Theory*, 13, 341–360.
- SCHOTT, J. R. (2005): Matrix Analysis for Statistics, Wiley, New York, 2 ed.
- ZAFFARONI, P. (2025): "Factor Models for Conditional Asset Pricing," *Journal of Political Economy* (forthcoming).