

Improving the Measurement of Earnings Dynamics^{*}

Moira Daly[†]

Dmytro Hryshko[‡]

Iouri Manovskii[§]

December 10, 2020

Abstract

The stochastic process for earnings is the key element of incomplete markets models in modern quantitative macroeconomics. We show that a simple modification of the canonical earnings process used in the literature leads to a substantial improvement in the measurement of earnings dynamics in administrative and survey data alike. Empirically, earnings at the start or end of earnings spells are lower and more volatile than the observations in the interior of earnings histories, reflecting mainly the effects of working less than the full year. Ignoring these properties of earnings, as is standard in the literature, leads to a substantial mismeasurement of the variances of permanent and transitory shocks and induces the large and widely documented divergence in the estimates of these variances based on fitting the earnings moments in levels or growth rates. Accounting for these effects enables more accurate analysis using quantitative models with permanent and transitory earnings risk and improves empirical estimates of consumption insurance against permanent earnings shocks.

Keywords: Earnings processes, Incomplete markets, Partial insurance, Estimation

JEL Classifications: D31, D91, E21

^{*}We benefited from comments of participants at numerous conferences and seminars. Support from the National Science Foundation Grants No. SES-0922406, 1357903, 1824520, the Danish Social Science Research Council (FSE) grant No. 300279, and SSHRC IG grant 435-2018-1275 is gratefully acknowledged.

[†]Centre for Economic and Business Research, Copenhagen Business School, Porcelaenshaven 16A, 2000 Frederiksberg, Denmark. E-mail: moda.eco@cbs.dk.

[‡]University of Alberta, Department of Economics, 8-14 HM Tory Building, Edmonton, AB, T6G2H4, Canada. E-mail: dhryshko@ualberta.ca.

[§]University of Pennsylvania, Department of Economics, The Ronald O. Perelman Center for Political Science and Economics, 133 South 36th Street, Philadelphia, PA 19104. E-mail: manovskii@econ.upenn.edu.

1 Introduction

The central element of many models in modern quantitative macroeconomics with heterogeneous agents is either an exogenously specified or an endogenously determined stochastic process for individual earnings. For example, in the models with incomplete insurance markets, the properties of the earnings process serve as key determinants of the evolution of consumption, assets, and other economic choices over the life cycle and across individuals.¹ Following the seminal contribution by Friedman (1957), modern consumption theory recognizes that consumption should respond more to the longer-lasting or permanent, as opposed to transitory, innovations in earnings. This explains the keen interest in the literature in measuring the variances of these components using variants of the permanent/transitory earnings decomposition² written, in its basic form, as:

$$\begin{aligned}y_{it} &= \alpha_i + p_{it} + \tau_{it} \\p_{it} &= \phi_p p_{it-1} + \xi_{it} \\ \tau_{it} &= \theta(L)\epsilon_{it},\end{aligned}\tag{1}$$

where log-earnings y_{it} of individual i at time t consist of the permanent component, p_{it} , and the transitory component, τ_{it} . If ϕ_p is close to 1, the shocks ξ_{it} are highly persistent (and are truly permanent if ϕ_p is 1), and if $\theta(L) = 1$ (where $\theta(L)$ is a moving average polynomial in the lag operator L), the shocks ϵ_{it} are completely transitory.

In addition to determining equilibrium consumption and wealth distributions, the variance and persistence of the shocks ξ_{it} and ϵ_{it} have important implications for policy design. For example, they are crucial for determining the optimal design of the bankruptcy code in Livshits et al. (2007), they govern the impact of the welfare system on household savings in Hubbard et al. (1995), stimulus effects of fiscal policy in Heathcote (2005) and Hagedorn et al. (2019b), the transmission mechanism of monetary policy in Kaplan et al. (2018) and Hagedorn et al. (2019a) as well as the optimal design of the tax system in Banks and Diamond (2010) and Farhi and Werning (2012). Moreover, there is great interest in understanding whether the dramatic increase in earnings dispersion over the last few decades in the U.S. is due to the increase in the variances of persistent or transitory shocks (e.g., Gottschalk and Moffitt, 1994). A better understanding of this increase in earnings dispersion could help determine why consumption inequality did not increase nearly as much (e.g., Krueger and Perri, 2006; Blundell et al., 2008; Heathcote et al., 2010b; Attanasio et al., 2012). Knowing the stochastic nature of earnings is also essential for the design of active labor market policies. For example, Meghir and Pistaferri

¹See, e.g., Deaton (1991), Carroll (1997), Castañeda et al. (2003).

²This decomposition was pioneered by Friedman and Kuznets (1954) and empirically supported by MaCurdy (1982), Abowd and Card (1989), and Meghir and Pistaferri (2004), among others. A prominent alternative, e.g., Guvenen (2009), allows for less persistent shocks but individual-specific trends in earnings.

(2011) suggest that income maintenance policies might be an appropriate response to changes in inequality driven by transitory shocks while training programs are potentially more relevant to counteracting the effects of permanent shocks.

Unfortunately, despite their manifest importance, there is no consensus in the literature on the size of the shocks ϵ_{it} and ξ_{it} . In particular, using the same data, the estimates of the earnings process in Eq. (1) when targeting the moments of log-earnings in levels are dramatically different from the estimates obtained when fitting the moments of log-earnings in differences. Although this discrepancy was first documented using survey-based data, it remained undiminished when the focus of the literature has shifted to relying more on administrative datasets.³ These datasets are typically orders of magnitude larger than survey-based ones; free of sampling issues; do not suffer from the typical issues of attrition; are based on administrative sources, such as tax records, which are considered highly reliable and free of issues of systematic nonresponse or measurement errors that typically plague survey-based data. Yet, despite their numerous attractive properties, these datasets must also have features that lead to the large discrepancy in the estimates based on moments in growth rates and levels.

Such an observation led Heathcote et al. (2010a) to conclude that the widely used model of earnings dynamics in Eq. (1) is misspecified. Unfortunately, in the absence of knowledge of the nature of this potential misspecification one cannot be confident in the conclusions of the models that incorporate this earnings process. Even if this misspecified process is used as a model input because there is no better alternative, whether it is more appropriate to parameterize it using the estimates targeting the data moments in levels or in differences is unclear. Relatedly, in the literature that endogenizes the earnings process,⁴ it is unclear whether the process implied by the model should be compared to the one estimated from the data using the specification in levels or differences, given that estimating the reduced-form process (1) on the model-generated data does not give rise to the observed discrepancy.

In this paper, we uncover an important source of this misspecification. Estimation of the parameters of the earnings process in the literature is based on fitting the entire set of autocovariance moments for levels or differences of log-earnings. However, even when estimation is based on the same set of observations in the data, computation of the autocovariance moments in levels and differences is effectively based on different information. To clarify with an example, consider an individual with a single earnings observation in the sample. This observation will contribute to the estimated variance of earnings in levels, but it will not contribute to any moment in differences. More generally, earnings observations adjacent to a missing one (e.g., observations at the start or at the end of an individual's earnings history) also contribute

³Recent contributions include Blundell et al. (2015), DeBacker et al. (2013), Domeij and Flodén (2010), Guvenen et al. (2014), among others.

⁴E.g., Huggett et al. (2011) and Postel-Vinay and Turon (2010).

differently to the moments in levels and differences. If earnings observations surrounding the missing ones were random draws similar to observations from the rest of earnings histories, this would not matter. However, in the data these earnings observations are much lower and more volatile. We will show formally that this data feature raises the variance of transitory shocks when estimation relies on the moments in levels and raises the variance of permanent shocks recovered by estimation based on the moments in differences.

In the first set of quantitative experiments in the paper we assess the magnitude of these effects using large administrative datasets from Denmark and Germany. The Danish data contain complete earnings histories of each resident of Denmark from 1981 through 2006. The German data are a 2% random sample of social security numbers. For these individuals, we observe the complete earnings history from 1984 through 2008. These samples are sufficiently large to allow analysis at the level of particular age cohorts, making it possible to focus on a parsimonious earnings model in (1), sidestepping the issue of modeling cohort effects. Moreover, the large size of the data enables reliable estimation when replicating the design of samples typically used in the literature. Specifically, we consider a balanced sample spanning 25 (26) years in German (Danish) data, a sample with 9 or more consecutive observations, as in e.g., Browning et al. (2010) and Meghir and Pistaferri (2004), and a sample with 20 or more not necessarily consecutive observations as in, e.g., Guvenen (2009). Our smallest Danish sample is comprised of about 67,000 individuals and 1.7 million observations, while our smallest German sample contains about 10,000 individuals with more than 200,000 observations.

Using the unbalanced samples in both datasets, we find, consistent with the literature, a substantially higher estimated variance of permanent (transitory) shocks targeting the moments of earnings in growth rates (levels). In contrast, we find that the discrepancy is nearly absent in balanced samples drawn from the two datasets. To highlight the special nature of the earnings trajectories in unbalanced samples, in Figure 1 we plot means and variances of (residual) earnings for the German data for individuals whose earnings are first observed in the data only after the start of the sample period in 1984, or for those whose earnings are last observed in the data before the end of the sample window in 2008. Panels (a) and (b) show that earnings are considerably lower on average, and panels (c) and (d) show that they are substantially more volatile than typical earnings observations at the start and end of incomplete earnings spells, respectively.⁵ Appendix Figure A-1 presents a different visualization of these patterns in unbalanced samples and, in addition, highlights their absence in the balanced sample. What then makes balanced and unbalanced samples different?

For the vast majority of individuals in the balanced sample, their first year in the sample does not coincide with the first year of their earnings history. Similarly, their last year in the

⁵Log residual earnings are residuals from cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. We subtracted individual fixed effects from residual earnings before taking means and variances.

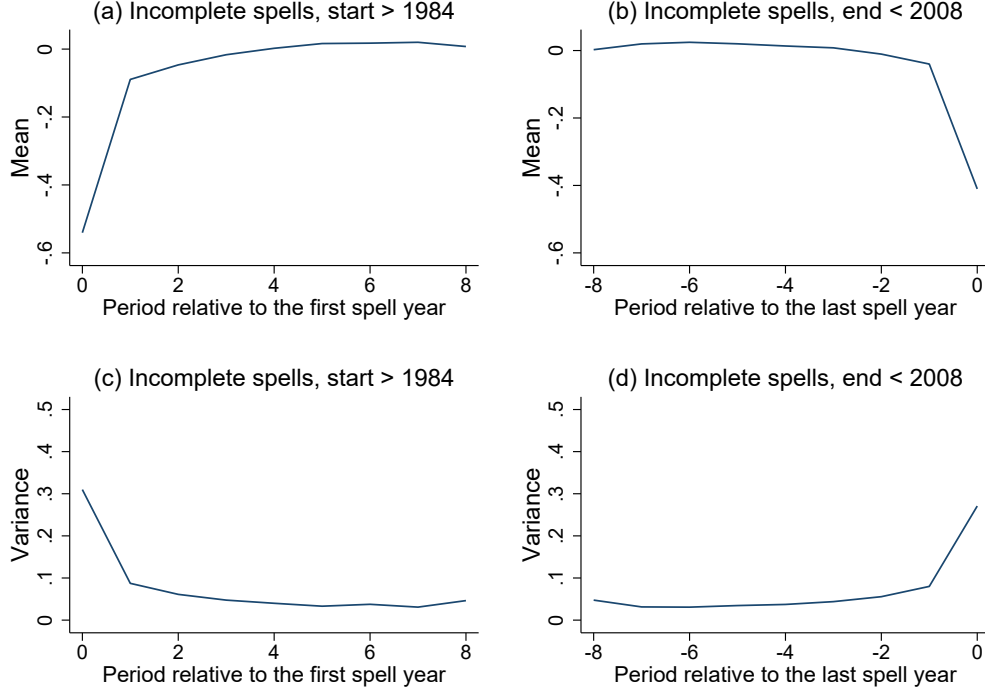


FIGURE 1: PROPERTIES OF RESIDUAL EARNINGS IN INCOMPLETE SPELLS. GERMAN DATA

sample mechanically truncates earnings histories, implying that it is not the last year of the actual earnings spell. Thus, earnings of these individuals in the first and last sample years are not expected to differ systematically from the neighboring observations in the interior of the sample window.

This stands in sharp contrast to the earnings histories of individuals entering and/or exiting the data in the interior of the sample window. Entering the sample in the interior of the sample window implies that these individuals do not have a valid earnings observation, at least in the first year of the sample window. By construction, their first available earnings observation must be preceded by a missing one. Similarly, the last earnings observation of a person exiting the sample early must be followed by a missing one. The reasons why earnings observations could be missing in administrative data will differ across datasets, but most of them do not include earnings of individuals who are self-employed or work for the government, especially in tenured positions (because in many countries these individuals do not contribute to the regular Social Security system from which the data typically come). Individuals who die or move abroad and pay taxes there will be missing from the data. Individuals will not typically have earnings records in the years they devote to education. Long unpaid paternity or maternity leave may also result in missing earnings observations. Individuals will typically not appear in the data while they are out of work and are not receiving unemployment or welfare benefits. Some missing observations are created by researchers themselves through

setting sample restrictions. For example, in an attempt to guard against errors in the data it is common to remove earnings outliers by setting them to missing. Most studies also set to missing observations when individuals do not satisfy some restrictions on, e.g., the minimum amount of time worked during a year. Except for the observations set to missing by researchers when processing the data, in most datasets it is not even possible to know why a particular observation is missing.⁶

Importantly, regardless of why earnings might be missing, earnings that follow or precede these periods ought to be relatively low in expectation and vary widely across individuals for the following reasons. The data in most administrative datasets are annual. Individuals who join the sample in year t but had missing earnings in year $t - 1$ because they were, e.g., unemployed, students, self-employed, or government workers are not expected to become regular employees on January 1 of year t , meaning that their annual earnings (from regular jobs) in the first year of earnings history are going to be relatively low in expectation and vary widely across individuals depending on when they actually started regular employment. Similarly, not all individuals die or move out of the country on December 31, and neither do they necessarily leave regular employment for other activities on that day. Thus, earnings in the last year are also going to be relatively low and more volatile. There are also some wage effects in the first and last years of incomplete earnings spells that we document below, but the transitory effect of incomplete working years is clearly dominant.

Our theoretical argument implies that the nonrandomness of earnings that surround missing observations in the unbalanced samples can induce the discrepancy between the estimates in levels and differences in the data from the unbalanced samples. The quantitative question is how large this effect is. To provide an answer, we proceed in three steps. First, we quantify the contribution of the low mean and high variance of earnings surrounding missing observations in the unbalanced samples to the subset of theoretical autocovariance moments on which the identification argument in levels and differences is based and confirm that they induce the observed discrepancy in the estimates. Second, using unbalanced samples, we remove a few observations adjacent to missing records. We find that estimating the earnings process in levels and in differences on the remaining data yields virtually identical estimates of the variances of permanent and transitory shocks. Third, we simulate artificial data based on these estimates of the earnings process while replicating the structure of the unbalanced samples (by design, observations surrounding missing records are not systematically different from observations in the rest of the earnings histories). We find no discrepancy between the estimates in levels and differences when these artificial data are used. We then draw an additional shock at the start and end of the contiguous earnings spells to replicate the mean and the variance of earnings in those periods in the data. We find that in this case, the estimates of the variance

⁶In Appendix II we exploit the richness of the Danish data to document the relative importance of various reasons for why individuals may not be present in the administrative earnings records.

of permanent and transitory shocks are very different when moments in levels and differences are used but are very close to those in the data from the corresponding unbalanced samples.

The results of these experiments lead us to conclude that the discrepancy in the estimates of the earnings process (1) in growth rates and levels is indeed driven by its misspecification. The nature of the misspecification is surprisingly simple. It is driven by the high variance and the low mean of the observations surrounding missing records. We show that an extended earnings process that includes these data elements can be estimated in the data. Estimating such an extended process results in similar parameters regardless of whether the moments in levels or differences are used.

While the focus of this paper is on the administrative data on individual earnings, the insights apply more broadly. In the appendix, we show that both individual earnings and hourly wage observations surrounding the missing observations in survey data from the Panel Study of Income Dynamics (PSID) exhibit a lower mean and higher variance than typical observations, and that accounting for the properties of these observations helps reconcile the dramatically different estimates of the variances of permanent and transitory shocks when targeting the moments in levels or differences. Hryshko and Manovskii (2018) show that the same findings also hold for total household earnings and net family incomes in the PSID. They also show that not taking these features of the data into account when estimating the stochastic properties of family earnings and income induces a large bias on the estimated degree of consumption insurance achieved through household saving and borrowing as well as on the estimated insurance role of the tax and transfer system.

Although the mechanism described in this paper is powerful in reconciling the estimates of the earnings process in growth rates and levels, it is not the only mechanism that can generate such a discrepancy. For example, Hryshko and Manovskii (2017) show that the discrepancy can also be induced if one restricts the permanent component to a random walk when its true persistence is lower. Importantly, this type of misspecification cannot generate the difference between the theoretical moments that we use to establish identification in levels and differences in this paper because they are identically affected by any such misspecification. These theoretical identifying moments can only differ if the underlying autocovariance moments on which they are based disagree, and we show that this is indeed the consequence of the low mean and high variance of observations at the start and end of earnings spells. We find that this accounts for virtually all of the discrepancy between the estimates in growth rates and levels in the data considered in this paper.

The rest of the paper is organized as follows. In Section 2, we discuss identification of the permanent-transitory decomposition of earnings and derive theoretically the biases in the estimated variances of permanent and transitory shocks when using the moments in levels and differences constructed from an unbalanced panel. In Section 3, we describe the administrative Danish and German data, the estimation procedure and basic results, and document the low

mean and high variance of earnings observations surrounding the missing ones. In Section 4, we show that this property of earnings quantitatively accounts for the difference in the estimates of earnings processes in levels and differences in our administrative data. In the Online Appendix, we confirm that our findings based on administrative data also apply to survey data on earnings and wages from the PSID. Section 5 concludes.

2 Sources of the Differences in Estimates Targeting Earnings Growth Rates and Levels

In this section, we first outline identification of the canonical earnings process that consists of a permanent random walk component and a nonpersistent transitory shock. We then derive the biases in the estimated variances of permanent and transitory shocks in unbalanced panels that feature lower on average and more volatile earnings next to the missing records. We close the section by discussing the biases when permanent shocks have finite persistence or transitory component is modeled as a moving-average process.

2.1 Identifying Moments for the Canonical Earnings Process

In the literature, estimation of the parameters of the earnings process typically relies on the minimum-distance method. In particular, estimation based on the moments in levels targets the entire set of autocovariance moments $E[y_{it}y_{it+j}]$, where $i \in [1, N]$ denotes individuals in the sample, t denotes time, and j denotes all leads and lags of earnings observed in the data. In differences, estimation targets the full set of autocovariance moments $E[\Delta y_{it}\Delta y_{it+j}]$, where Δ is the difference operator between two consecutive observations, $\Delta y_{it+j} \equiv y_{it+j} - y_{it+j-1}$.

Although all available autocovariance moments are used in estimation, identification is usually established using only a subset of autocovariance moments; see, e.g., Meghir and Pistaferri (2004), Blundell et al. (2008), and Heathcote et al. (2014). For example, consider the earnings process that consists of a random walk and an i.i.d. transitory shock, which corresponds to setting $\theta(L)$ and ϕ_p to 1 in Eq. (1). This process was considered by Heathcote et al. (2010a), who, assuming that transitory and permanent shocks are not correlated with each other and with initial conditions, proposed the following moments to identify the variances of permanent and transitory shocks at time t :

Differences:

$$\sigma_{\xi,t}^2 = E[\Delta y_{it}\Delta y_{it-1}] + E[\Delta y_{it}\Delta y_{it}] + E[\Delta y_{it}\Delta y_{it+1}], \quad (\text{D1})$$

$$\sigma_{\epsilon,t}^2 = -E[\Delta y_{it}\Delta y_{it+1}]. \quad (\text{D2})$$

Note that (D1) and (D2) represent linear combinations of autocovariance moments for earnings growth rates. For clarity, we will refer to individual autocovariance moments as simply “moments,” and to a linear combination of autocovariance moments used for identification such as (D1) and (D2) as “identifying moments.”

Expanding (D1) and (D2), we obtain the identifying moments for the variances of permanent and transitory shocks, based on autocovariance moments in levels, at time t :

Levels:

$$\sigma_{\xi,t}^2 = E[y_{it}y_{it+1}] - E[y_{it+1}y_{it-1}] - E[y_{it}y_{it-2}] + E[y_{it-1}y_{it-2}], \quad (\text{L1})$$

$$\sigma_{\epsilon,t}^2 = E[y_{it}y_{it}] - E[y_{it}y_{it+1}] - E[y_{it-1}y_{it}] + E[y_{it-1}y_{it+1}]. \quad (\text{L2})$$

In a sample of individuals whose earnings are nonmissing for the periods $t - 2$ through $t + 1$, the identifying moments (D1)-(D2) and (L1)-(L2) are expected to deliver identical estimates of the variance of permanent and transitory shocks at time t , since they are based on exactly the same earnings information. Moreover, as the moments (L1)-(L2) simply represent an expansion of the moments (D1)-(D2), they will be identically affected by other potential misspecifications of the earnings process. This allows us to isolate and measure the importance of the high variance and low mean of the observations at the start and end of contiguous earnings spells, which, as we show below, contribute differently to the autocovariance moments on which (D1)-(D2) and (L1)-(L2) are based.⁷

For example, the presence of omitted idiosyncratic trends in earnings will not induce a wedge between the estimated variances of permanent shocks using the moments (L1) and (D1) (or transitory shocks using the moments (L2) and (D2)). Specifically, suppose individuals differ in growth rates such that the earnings process is $y_{it} = \alpha_i + \beta_i h_{it} + p_{it} + \epsilon_{it}$, where $\beta_i \sim \text{iid}(0, \sigma_\beta^2)$ and h_{it} counts years of (potential) work experience. In this case (L1) and (D1) will both deliver $3\sigma_\beta^2 + \sigma_{\xi,t}^2$.⁸ It follows that both (L1) and (D1) will recover an upward-biased estimate of the variance of the permanent shock, but the bias will be the same in levels and differences. Relatedly, the typical estimates of σ_ξ^2 using (D1) imply a much steeper profile of

⁷Identifying moments in levels can be constructed using fewer autocovariance moments, such as

$$\sigma_{\xi,t}^2 = E[y_{it}y_{it+1}] - E[y_{it}y_{it-1}], \quad (\text{L1-Short})$$

$$\sigma_{\epsilon,t}^2 = E[y_{it}y_{it}] - E[y_{it}y_{it+1}]. \quad (\text{L2-Short})$$

These moments are analogous to those in Heathcote et al. (2010a) if one relies on the annual data, instead of biennial data used in their paper, for identification of the variances. These identifying moments in levels do not, however, use the same information as the identifying moments (D1)-(D2) in differences. For example, the information on earnings in $t - 2$ is used in (D1) but not in (L1-Short). Moreover, Hryshko and Manovskii (2017) show that a misspecification of the persistence of the permanent component drives a wedge between the estimates based on identifying moments (D1)-(D2) and (L1-Short)-(L2-Short), but not between identifying moments (D1)-(D2) and (L1)-(L2).

⁸This derivation assumes that $\text{corr}(\alpha_i, \beta_i) = 0$ but a similar expression obtains if this assumption is relaxed.

earnings inequality over the life cycle (and time) than that observed in the data. The fit to this profile might be improved if one allows for a negative cross-sectional correlation between initial conditions, α_i , and permanent shocks, ξ_{it} . However, omitting such correlation does not induce a difference in the estimated moments (L1) and (D1). For example, suppose that the correlation is implied by $\xi_{it} = \kappa\alpha_i + \eta_{it}$, where η_{it} is orthogonal to α_i and ϵ_{it} . In this case, (D1) and (L1) will recover identical upward-biased estimate $3\kappa^2\sigma_\alpha^2 + \sigma_{\xi_t}^2$, but the bias will once again be the same in levels and differences.

Importantly, each autocovariance moment is measured as the average across all available observations that contribute to it. This implies that, although the identifying moments (D1)-(D2) and (L1)-(L2) are based on the same earnings data, the autocovariance moments used in estimating (D1)-(D2) and (L1)-(L2) are computed using different sets of observations. Returning to the extreme example used in the Introduction, consider an individual who appears in the sample only once, in period t . This individual will contribute to the autocovariance moment $E[y_{it}y_{it}]$, and thus his only earnings observation will affect the identifying moment (L2), but it will not contribute to any autocovariance moment used to construct the corresponding identifying moment in differences (D2). If earnings of individuals who appear in the sample only once are systematically different, this will induce the difference between identifying moments (L2) and (D2) and lead to different estimates of the variance of transitory shocks using the moments in levels and differences.

Similarly, we will now show that earnings observations at the time individuals enter or exit the sample contribute differently to the autocovariance moments on which the identifying moments (D1)-(D2) and (L1)-(L2) are based. Moreover, our empirical analysis below will reveal that these earnings observations are systematically different (they are typically lower and substantially more volatile). In the rest of this section we formally show that this induces systematic differences in estimated variances of permanent and transitory shocks using the moments in growth rates and levels. In subsequent sections, we quantify the magnitude of the induced difference.

2.2 Potential Bias Induced by Observations Next to the Missing Ones in Various Samples

We will consider three types of samples. Suppose a dataset containing panel data on individual earnings starts in period t_0 and ends in period T . We refer to the sample as balanced if all individuals in the sample have $T - t_0 + 1$ valid earnings observations. While not part of the formal definition, it is convenient to think that earnings spells of individuals in the balanced samples start before t_0 and end after T . In other words, the boundaries of the balanced sample mechanically truncate continuous earnings spells in progress. We refer to samples that include only uninterrupted earnings spells (i.e., no gaps) but with a duration of less than $T - t_0 + 1$ for

at least some individuals as consecutive unbalanced samples. Finally, we refer to unbalanced samples that also include individual earnings spells interrupted by missing observations in any period $t \in (t_0, T)$ as nonconsecutive unbalanced samples.

2.2.1 Biases in Consecutive Unbalanced Samples

The nature of consecutive unbalanced samples is that at least some individuals are observed starting or ending their earnings spells inside the sample window.

As mentioned above and documented below, earnings have a lower mean and are highly volatile in the first and last periods of an incomplete earnings history. Consider modeling this through an additional ν -shock that occurs only in the first and last year of an individual's earnings history, that is

$$y_{it} = \alpha_i + p_{it} + \epsilon_{it} + \nu_{it},$$

where ν_{it} has mean μ_ν (taking a negative value) and variance σ_ν^2 , and is uncorrelated with permanent and transitory shocks and initial conditions. We will now show that ignoring ν_{it} and estimating the process (1) instead leads to an upward bias in the estimated variance of permanent shocks using the moments in differences and the estimated variance of transitory shocks using the moments in levels.

For simplicity, assume there is a set of individuals first entering the sample at time t , in the interior of the sample period $[t_0, T]$, whereas the remaining individuals are continuously observed throughout the sample. Individuals first appearing at time t will contribute to estimation of the autocovariance moments $E[y_{it}y_{it}]$ and $E[y_{it}y_{it+1}]$ in the identifying moment (L2). The estimated moment $E[y_{it}y_{it+1}]$ will be no different for such individuals than for the rest of the sample, and will equal $\sigma_\alpha^2 + \text{var}(p_{it})$. The other moments in (L2), $E[y_{it-1}y_{it}]$ and $E[y_{it-1}y_{it+1}]$, will both equal $\sigma_\alpha^2 + \text{var}(p_{it-1})$. The autocovariance moment $E[y_{it}y_{it}]$ estimated on the full sample, however, will equal $\sigma_\alpha^2 + \text{var}(p_{it}) + \sigma_{\epsilon,t}^2 + s_t(\mu_\nu^2 + \sigma_\nu^2)$, where s_t is the share of individuals, at time t , whose (incomplete) spells start at time t in the total number of individuals at time t with nonmissing earnings. The identifying moment (L2), therefore, will recover an estimate of the variance of transitory shocks equal to $\sigma_{\epsilon,t}^2 + s_t(\mu_\nu^2 + \sigma_\nu^2)$, with an upward bias of $s_t(\mu_\nu^2 + \sigma_\nu^2)$.

The variance of permanent shocks at time $t + 1$, estimated using the identifying moment (D1), will also be biased upward. Individuals first appearing at t will contribute to estimation of the autocovariance moments $E[\Delta y_{it+1}\Delta y_{it+1}]$ and $E[\Delta y_{it+1}\Delta y_{it+2}]$ in the identifying moment (D1). For such individuals, the autocovariance moment $E[\Delta y_{it+1}\Delta y_{it+2}]$ will be no different from the rest of the sample and will equal $-\sigma_{\epsilon,t+1}^2$, while the autocovariance moment $E[\Delta y_{it+1}\Delta y_{it+1}]$ will equal $\sigma_{\xi,t+1}^2 + s_{t,t+1}(\mu_\nu^2 + \sigma_\nu^2) + \sigma_{\epsilon,t}^2 + \sigma_{\epsilon,t+1}^2$, where $s_{t,t+1}$ is the share of individuals who start (incomplete) earnings spells at time t , with nonmissing earnings at times t and $t + 1$, in the number of individuals with nonmissing earnings both at t and $t + 1$. Since

the autocovariance moment $E[\Delta y_{it+1} \Delta y_{it}]$ will be estimated using information for those individuals whose earnings are nonmissing in periods $t - 1$ through $t + 1$ and will equal $-\sigma_{\epsilon_t}^2$, the identifying moment (D1) for time $t + 1$ will recover an estimate of the permanent shock equal to $\sigma_{\xi_{t+1}}^2 + s_{t,t+1}(\mu_\nu^2 + \sigma_\nu^2)$, with an upward bias of $s_{t,t+1}(\mu_\nu^2 + \sigma_\nu^2)$.

Note that if the ν -shock first appears, say, at time $t + 1$, i.e., in the interior of an earnings spell for individuals first entering into the sample at time t , it will simply elevate, by the same magnitude, the estimated variance of transitory shocks in levels and differences at time $t + 1$, with no differential effect on the identifying moments (L2) and (D1).

Summing up, incomplete earnings spells first appearing in the sample at t will bias upward the estimated variance of transitory shocks at time t when targeting the moments in levels and will bias upward the estimated variance of permanent shocks at time $t + 1$ when targeting the moments in differences. They have no effect, at any point in time, on the estimated magnitude of the identifying moments (L1) and (D2).

The same logic extends to the incomplete earnings spells ending at time t , which is different from the last potential sample year T – the presence of such spells will produce upward-biased estimates of permanent variances in differences at t (since these individuals will contribute to the estimation of the moment $E[\Delta y_{it} \Delta y_{it}]$ that is part of the identifying moment D1) and of transitory variances in levels at t .

2.2.2 Biases in Nonconsecutive Unbalanced Samples

We now consider the consequences of missing earnings in the interior points of the earnings history. We assume that individual earnings are realizations of the earnings process (1), with some observations missing in any period $t \in (t_0, T)$ but not in periods $t - 1$ and $t + 1$. We will show below that such periods are often associated in the data with low mean and high variance of earnings in periods $t - 1$ and $t + 1$. We model this by introducing additional transitory ν -shocks with a negative mean μ_ν at the time before and after earnings are missing (ν_{it-1} and ν_{it+1} , respectively) that are assumed to be uncorrelated with fixed effects, permanent and transitory shocks, and uncorrelated with each other:⁹

$$\begin{aligned} y_{it-1} &= \alpha_i + p_{it-1} + \epsilon_{it-1} + \nu_{it-1}, \\ y_{it} &\text{ missing}, \\ y_{it+1} &= \alpha_i + p_{it+1} + \epsilon_{it+1} + \nu_{it+1}. \end{aligned}$$

Assume there is a set of individuals whose earnings are missing in period $t \in (t_0, T)$ while the rest of the individuals have continuously observed earnings throughout the whole sample

⁹For ease of exposition, we assume that the mean and variance of the ν -shock one year before and after earnings are missing are the same, although in the data they slightly differ.

period.

In this case, the variance of transitory shocks at times $t - 1$ and $t + 1$ using the moments in levels will be biased upward as the autocovariance moments $E[y_{it-1}y_{it-1}]$ and $E[y_{it+1}y_{it+1}]$ in the identifying moment (L2) are amplified by the variation of the ν -shocks. Similarly, the variance of permanent shocks at times $t - 1$ and $t + 2$ using the moments in differences will be biased upward as the autocovariance moments $E[\Delta y_{it-1}\Delta y_{it-1}]$ and $E[\Delta y_{it+2}\Delta y_{it+2}]$ in the identifying moment (D1) are amplified by the variation of the ν -shocks. Since the ν -shocks are assumed to be uncorrelated, the identifying moment (D2) will not be affected. Since ν -shocks have negative means, there will be a downward bias in the estimated variance of permanent shocks using (L1) due to an amplified moment $E[y_{it-1}y_{it+1}]$.¹⁰

Thus, incomplete earnings spells with missing earnings at t , in the interior of the sample period, will bias upward the estimated variance of transitory shocks at times $t - 1$ and $t + 1$ when targeting the moments in levels, and will bias upward the variance of permanent shocks at times $t - 1$ and $t + 2$ when targeting the moments in differences.

2.3 Extensions

2.3.1 Biases When ξ_{it} Shocks Have Limited Persistence

If ϕ_p in Eq. (1) is less than 1, one must rely on a modified set of identifying moments to recover the permanent and transitory variances. For a given estimate of the persistence ϕ_p , which can be separately identified,¹¹ the set of identifying moments will amount to

Differences:

$$\sigma_{\xi,t}^2 = E[\tilde{\Delta}y_{it}\tilde{\Delta}y_{it+1}] + \phi_p E[\tilde{\Delta}y_{it}\tilde{\Delta}y_{it}] + \phi_p^2 E[\tilde{\Delta}y_{it}\tilde{\Delta}y_{it-1}], \quad (\text{D1-a})$$

$$\sigma_{\epsilon,t}^2 = -\frac{1}{\phi_p} E[\tilde{\Delta}y_{it}\tilde{\Delta}y_{it+1}], \quad (\text{D2-a})$$

where $\tilde{\Delta}y_{it} \equiv y_{it} - \phi_p y_{it-1}$.

Expanding the above moments results in the following set of moments in levels identifying

¹⁰A further downward bias in the variance of permanent shocks at time t using (L1) will occur due to missing observations at $t - 1$ surrounded by nonmissing records at $t - 2$ and t . These downward biases will be absent if one relies on (L1-short) rather than (L1) for estimating the variance of transitory shocks using the moments in levels.

¹¹The persistence ϕ_p can be recovered from the moments $\frac{E[y_{it+k+3}y_{it+k}] - E[y_{it+k+2}y_{it+k}]}{E[y_{it+k+2}y_{it+k}] - E[y_{it+k+1}y_{it+k}]}$ for $k \geq 0$. One can also use the moments in growth rates to identify it; see, e.g., Hryshko (2012). There is also a large literature, reviewed in MaCurdy (2007) and Arellano and Honoré (2001), that does not rely on fitting the autocovariance function of earnings but exploits various orthogonality conditions in a GMM setting to recover the persistence.

the variances at time t :

Levels:

$$\sigma_{\xi,t}^2 = E[y_{it}y_{it+1}] - \phi_p E[y_{it+1}y_{it-1}] - \phi_p^3 E[y_{it}y_{it-2}] + \phi_p^4 E[y_{it-1}y_{it-2}], \quad (\text{L1-a})$$

$$\sigma_{\epsilon,t}^2 = E[y_{it}y_{it}] - \frac{1}{\phi_p} E[y_{it}y_{it+1}] - \phi_p E[y_{it-1}y_{it}] + E[y_{it-1}y_{it+1}]. \quad (\text{L2-a})$$

Although the biases for the variance of transitory shocks in levels will be exactly the same as in the random-walk case, the biases for the variance of permanent shocks recovered using the identifying moments in differences will be scaled by the persistence ϕ_p . Note, however, that the size of the bias is unlikely to decline substantially since ϕ_p is typically estimated at high values in various datasets.

2.3.2 Biases When Transitory Component and/or ν -shocks are Serially Correlated

The transitory component is often estimated to have some persistence. Assume that the transitory component is modeled as $\tau_{it+1} = \epsilon_{it+1} + \theta_\tau \epsilon_{it}$, and that the ν -shock component is modeled as $\chi_{it} = \nu_{it}$, which is nonzero in the beginning and/or end of an incomplete earnings spell, and before/after a missing earnings record, and that $\chi_{it+1} = \theta_\chi \nu_{it}$ – both will be consistent with the autocovariance function for earnings growth rates truncating at the second order, as is often found in the empirical applications.¹² In this case, the moments (L1)–(D2) no longer identify the variances of permanent and transitory shocks. In growth rates, the identifying moment for the variance of permanent shocks should be modified to

$$\sigma_{\xi,t}^2 = E[\Delta y_{it} \Delta y_{it+2}] + E[\Delta y_{it} \Delta y_{it+1}] + E[\Delta y_{it} \Delta y_{it}] + E[\Delta y_{it} \Delta y_{it-1}] + E[\Delta y_{it} \Delta y_{it-2}]. \quad (\text{D1-b})$$

The variance of permanent shocks at time $t + 1$, estimated using (D1-b), will be biased upward by the magnitude $s_{t,t+1}(1 - \theta_\chi)^2(\mu_\nu^2 + \sigma_\nu^2)$ for a sample with consecutive earnings spells where a fraction of individuals enter the sample at time $t > t_0$, for the first time. Note that the bias will remain large for small positive values of θ_χ . If, instead, individuals exit the sample at some time $t < T$, the bias of the permanent variance using the moments in growth rates will be unaffected by serial correlation of the ν -shocks since the earnings paths for such individuals are unobserved past year t ; the bias, in this case, will be the same as in the case of a serially uncorrelated transitory component. The same logic extends to the biases in the nonconsecutive samples. The variance of permanent shocks recovered using the moments in levels will remain unbiased (as can be verified from the identifying moment for permanent shocks in levels obtained by expanding (D1-b)).

¹²This formulation assumes that ν - and ϵ -shocks both die out in two periods, with the difference that the ν -shock process does not renew itself in the next period with a new ν -shock.

Under the assumption of no measurement error in administrative earnings, θ_τ can be identified from the first and second-order autocovariances in earnings growth rates if the transitory component is serially correlated and there are no ν -shocks; see, e.g., Meghir and Pistaferri (2004). One can then identify the variance of transitory shocks dividing (L2) and (D2) by $(1 - \theta_\tau)^2$. If the ν -shock is serially correlated, however, θ_τ will be recovered with a bias using the standard moment. We will label this estimate as $\tilde{\theta}_\tau$. Assuming that the variance of transitory shocks does not change much between adjacent periods, for the data with incomplete consecutive spells that start at t , an estimate of the variance of transitory shocks relying on (L2) will yield $(1 - \tilde{\theta}_\tau)^{-2} [(1 - \theta_\tau)^2 \sigma_{\epsilon_t}^2 + s(1 - \theta_\chi)(\mu_\nu^2 + \sigma_\nu^2)]$, whereas an estimate relying on (D2) will yield $(1 - \tilde{\theta}_\tau)^{-2} [(1 - \theta_\tau)^2 \sigma_{\epsilon_t}^2 - s\theta_\chi(1 - \theta_\chi)(\mu_\nu^2 + \sigma_\nu^2)]$ for $t + 1$.¹³ Clearly, an estimate of the variance of transitory shocks in levels is larger than an estimate using growth rates given θ_χ is nonnegative. This logic extends to other examples of incomplete earnings spells in consecutive and nonconsecutive panels – the estimated variance of transitory shocks using the moments in levels will be higher than the estimated variance of transitory shocks using the moments in growth rates.

2.4 Summary

The analysis above yields three major implications if ν -shocks are present in the data. First, estimating the abbreviated earnings process in (1), one may expect to recover fairly well the variance of transitory shocks using the moments in growth rates if the ν -shock is not serially correlated, and the variance of permanent shocks using the moments in levels. Second, the identifying moments in levels tend to produce upward-biased estimates of the variance of transitory shocks, while the identifying moments in differences produce upward-biased estimates of the variance of permanent shocks. The magnitude of the biases depends positively on the variance of the ν -shocks and on the difference between their mean from the mean of the shocks in the rest of earnings histories. Finally, if one's interest extends beyond identifying the variances of permanent and transitory shocks of the abbreviated earnings process in (1), the remaining parameters of the comprehensive earnings process can also be estimated by introducing the moments identifying the properties of ν -shocks.

3 Data, Estimation Details, and Basic Results

In this section, we first describe administrative Danish and German data used for estimation of earnings processes. We then present estimation of the canonical earnings process in Eq. (1) on balanced and unbalanced samples from the two datasets. Since the results on balanced samples substantively deviate from the results on unbalanced samples, we further present

¹³We assumed that $s_t = s_{t,t+1} = s_{t,t+2} = s$ in the derivation.

descriptive regression results for the level and volatility of residual earnings confirming the special nature of earnings records surrounding the missing ones in the unbalanced samples.

3.1 Data

In this section, we describe the administrative data and construction of the samples that we study. Following the literature, we focus on individuals with a strong attachment to the labor market characterized by sufficiently high earnings and time spent working.¹⁴

3.1.1 Danish data

Several administrative registers provided by Statistics Denmark were used to construct our samples. The Danish Integrated Database for Labor Market Research (IDA), containing earnings from the tax register from 1980–2006, provides panel data on total earnings for more than 99.9 percent of Danish residents between the ages of 15 and 70. Population and education registers are also merged so that demographic variables such as age and educational status could be appended. Our sample window spans the years from 1981 to 2006. The population we study consists of Danish males born in 1951 through 1955. We first remove all individuals who were ever self-employed and drop records in which an individual was making nonpositive labor market earnings. Next, we drop records for those individuals who have worked less than 10 percent of the year as a full-time employee. Annual earnings in a particular year include all earned labor income, taken from tax records, for that calendar year. This variable is considered “high quality” by Statistics Denmark in that it very accurately captures the earnings of individuals. Earnings are expressed in 1981 monetary units (Danish kroner). We calculate the maximum number of consecutive periods in which an individual has nonmissing earnings and use this information to construct two consecutive samples: a sample in which an individual’s maximum spell is at least nine consecutive periods (102,825 individuals), and a balanced sample in which the individual’s maximum spell covers the entire 26 periods (67,008

¹⁴The selection rules we adopt are typical of the literature that utilizes survey data as well as administrative data. For example, Guvenen et al. (2014) use U.S. administrative data on individual wage and salary income and make the following sample selection: “For a statistic computed using data for not necessarily consecutive years t_1, t_2, \dots, t_n , an individual observation is included if the following three conditions are satisfied for all these years: the individual (i) is between the ages of 25 and 60, (ii) has annual wage/salary earnings that exceed a time-varying minimum threshold, and (iii) is not self-employed (i.e., has self-employment earnings less than the same minimum threshold). This minimum, denoted $Y_{min,t}$, is equal to one-half of the legal minimum wage times 520 hours... This condition allows us to focus on workers with a reasonably strong labor market attachment and avoids issues with taking the logarithm of small numbers. It also makes our results more comparable to the income dynamics literature, where this condition is standard.” Similarly, DeBacker et al. (2013) “...exclude earnings (or income) observations below a minimum threshold...” and “...take the relevant threshold to be one-fourth of a full-year, full-time minimum wage.” In line with our selection of consecutive unbalanced samples (with the difference that we use at least nine consecutive earnings observations), Blundell et al. (2015) “...restrict the sample to individuals with at least four subsequent observations with positive market income.”

individuals). For the sample with nine or more consecutive observations, periods outside of the longest spell are discarded. Within the longest spell, an earnings outlier is defined by an increase in earnings of more than 500 percent or a fall of more than 80 percent in adjacent years. Individuals with earnings outliers within their longest spell are dropped. The third sample we consider consists of individuals who have at least 20 not necessarily consecutive periods in which they have nonmissing earnings (90,668 individuals). We also drop individuals from this sample if they have earnings growth outliers. Finally, we drop individuals if their educational status has changed during the spells considered. Table 1 contains basic statistics for selected samples.

3.1.2 German data

We use administrative panel data from the Sample of Integrated Labour Market Biographies (SIAB) that follows individuals selected as a 2% random sample of German Social Security records for the years 1974–2008. A detailed description of the dataset can be found in Dustmann et al. (2009). We use full-time job spells for German males born in 1951–1955, dropping the spells in East Germany. We also drop annual records when an individual was in apprenticeship during any part of the year. Individual real earnings is a sum of earnings from the jobs held within a year expressed in 2005 euros. We set individual education to the maximum schooling attained during the sample years and set the number of days worked to the sum of calendar days on all jobs within a year. As individual earnings are right-censored at the highest level subject to social security contributions, we impute earnings exceeding the limit assuming that daily wages in the upper tail follow a Pareto distribution, the parameters of which differ by year and age group.¹⁵ After 1983, earnings include one-time payments such as bonuses. To make variable definitions consistent throughout, we use only the data since 1984. We also drop individual records on annual earnings if the combined duration of job spells within a year is fewer than 35 calendar days and drop records with very low daily earnings, defined as the earnings records below 14 euros in 2003 prices. As in the Danish data, we construct three samples – balanced, with nine or more consecutive, and with 20 or more not necessarily consecutive earnings observations – and, as with the Danish samples, drop individuals who have earnings growth outliers. The respective samples contain 9,452, 18,130, and 13,635 individuals with 236,300, 379,080, and 330,748 observations, respectively. Table 2

¹⁵We consider the following eight age groups: those younger than 25, six five-year age groups (25–29, 30–34, 35–39, 40–44, 45–49, and 50–54), and those older than 54. We use a “fixed effects” imputation, keeping a uniform draw for each individual affected by the right-censoring limit fixed when creating a Pareto variate in different years. We also experimented with imputation based on the assumption that truncated log-wage distribution is normal, and a simpler imputation when the daily wage is multiplied by the factor 1.2 if it hits the upper censoring limit. These three imputation methods have been used in Dustmann et al. (2009). Our conclusions below are robust with respect to the choice of the imputation method as well as with respect to limiting the sample to individuals whose earnings histories are not affected by the censoring.

provides some descriptive details of the samples.

3.2 Estimation Details

As is standard in the literature, we estimate the earnings process in Eq. (1) using the method of minimum distance, fitting the entire set of autocovariances of log-earnings in levels or first differences to the autocovariance function implied by the model.¹⁶ We allow for an MA(1) transitory component and an unrestricted estimation of the persistence of the permanent component, ϕ_p .¹⁷ Thus, we estimate five parameters in total – the persistence and the variance of permanent shocks, ϕ_p and σ_ξ^2 ; the persistence and the variance of transitory shocks, θ and σ_ϵ^2 ; and the variance of individual fixed effects, σ_α^2 . We assume that individuals start accumulating permanent and transitory shocks at the age of twenty-five so that part of the estimated variance of fixed effects captures the accumulated permanent and transitory components prior to that age. We also assume that permanent and transitory shocks are not correlated with each other and with fixed effects. We remove predictable variation in earnings by estimating cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. Our measure of idiosyncratic earnings, consistent with the literature, is the residual from those regressions. Since our samples are large, we estimate the model using the optimal weighting matrix, which is an inverse of the variance-covariance matrix of the data moments.

3.3 Basic Results: Estimates of the Canonical Earnings Process

3.3.1 Samples with nine or more consecutive observations

Columns (1)–(4) in Table 3 contain estimation results for the samples with nine or more consecutive observations in the German and Danish data.¹⁸ The permanent component is estimated to be close to a random walk using the moments in differences, but slightly less persistent using the moments in levels. Importantly, in both datasets the variance of the permanent shock is about two times larger in the estimation that uses the moments in growth rates, while the estimated variance of the transitory shock is larger using the moments in levels. Thus, our administrative data exhibit the same large discrepancy that is endemic in this literature. The pattern is less pronounced in the Danish data, which is consistent with the mechanism we describe. In the Danish data, 65% of individuals have complete earnings

¹⁶One of the recent exceptions is Browning et al. (2010) who, apart from selected moments in levels and differences, fit a variety of other data moments studied in the literature on earnings dynamics.

¹⁷Allowing, instead, for an AR(1) transitory component has no substantive influence on the results.

¹⁸In differences, the variance of fixed effects is not identified. It can be identified using quasi-differences when $\phi_p < 1$ as in, e.g. Blundell et al. (2015). We do not pursue this approach as our focus is on the standard estimation in differences (and levels).

spells, while in the German data this number is only 52%. Consequently, fewer individuals have irregular earnings observations adjacent to the missing ones in the Danish data.¹⁹

3.3.2 Samples with 20 or more not necessarily consecutive observations

Columns (5)–(8) in Table 3 contain the results for the samples with 20 or more not necessarily consecutive observations. The variances of persistent shocks are somewhat smaller than those in columns (1)–(4), whereas the variances of transitory shocks are similar in magnitude. Importantly, we still observe that estimations using the moments in differences deliver relatively higher estimates of the variance of permanent shocks, while estimations in levels deliver relatively higher estimates of the variance of transitory shocks, once again confirming the widely documented discrepancy.²⁰

3.3.3 Balanced samples

Estimation results based on the balanced samples are reported in columns (9)–(12) of Table 3. Relative to the estimates on the unbalanced samples discussed above, the use of balanced samples results in a more than 50% reduction of the variance of permanent shocks when using the moments in differences. There is a similarly striking reduction of at least 50% in the variance of transitory shocks when using the moments in levels. It appears that the use of balanced samples largely eliminates the discrepancy between the estimates of the earnings process in levels and differences.

3.4 A Closer Look at Unbalanced Samples

In this section, we first show, in the regression setting, that earnings observations around missing ones are lower and more volatile. We then discuss the economic forces behind these data features.

3.4.1 The means and variances of earnings records surrounding missing observations

The results of estimation on balanced and unbalanced samples indicate that the discrepancy between the estimates based on the moments in levels and differences is specific to unbalanced samples. A defining feature of unbalanced samples is that some observations are missing. As discussed in Section 2, if earnings observations surrounding the missing ones are systematically

¹⁹Randomly dropping individuals with incomplete earnings histories in the German data to match their share in the Danish data results in similar discrepancies across the two datasets.

²⁰As was the case with nine or more consecutive observations, the discrepancy is less pronounced in the Danish data because the share of individuals with complete earnings spells is larger, and the share of missing earnings observations is smaller than in the German data.

different, this can induce the difference in the estimates of the earnings process in growth rates or levels. This data feature has the potential to explain why the discrepancy arises when the estimation is based on unbalanced samples only.

Indeed, Figure 1 in the Introduction revealed a clear pattern that the earnings at the start and/or the end of incomplete earnings spells are considerably lower on average and substantially more volatile than typical earnings observations. To explore these patterns more formally, in columns (1)–(4) of Table 4 we report the estimates from the fixed-effects panel regressions of residual earnings on dummies for the first and last years of individual earnings spells inside the overall sample window. Specifically, the dummies “Year observed: first”–“Year observed: third” equal one if an individual’s first earnings record in the sample occurs later than 1984 in the German data (later than 1981 in the Danish data), and zero otherwise, while the dummies “Year observed: second-to-last”–“Year observed: last” equal one if an individual’s last earnings record is prior to 2008 in the German data (2006 in the Danish data), and zero otherwise.²¹

In both samples and both datasets, earnings are about 0.50 to 0.60 log points lower than an individual’s average in the first year of the spell, whereas the last earnings record is below an individual’s average by about 0.30 to 0.40 log points. Earnings are still somewhat lower in the two years following the first earnings record as well as in the two years preceding the last earnings record. Moreover, earnings are, on average, also lower in the years preceding and following a missing earnings record in the nonconsecutive samples. Interestingly, the dummies for the few first and last earnings records within a spell explain 5 to 13 percent of the within-individual variation in residual earnings. This number is quite high taking into account that a variety of observable factors in a typical Mincer-style regression explain less than 30 percent of within-individual variation in earnings.

Performing the same experiment in reverse, we use our samples with 20 or more not necessarily consecutive observations to assess the predictive power of earnings dynamics for the incidence of missing earnings. Specifically, in Table A-5, the dependent variable is a dummy that equals 100 if individual earnings are missing and 0 otherwise. We find that the predictive power of observables – earnings growth rates before and after missing earnings records, together with education dummies and age – on the incidence of missing earnings is quite low, in line with Fitzgerald et al. (1998) who made a similar observation using PSID data. The strong (weak) earnings growth after (before) missing records lacks high explanatory power for a missing record because there are also many declines and subsequent recoveries of earnings

²¹To reinforce the conclusion that patterns in Table 4 are actually driven by starting and ending of the earnings spells, in Table A-3 of Appendix III we repeat the same analysis by focusing on individuals whose earnings spells begin in the first sample year or end in the last sample year. For the vast majority of these individuals, such cutoffs do not represent an actual start or end of their earnings spells; instead, the sample window mechanically truncates earnings spells in progress. Accordingly, the first (last) few dummies equal one if an individual’s first (last) earnings record is in the first (last) sample year, and zero otherwise.

inside uninterrupted earnings spells. Nonetheless, missing observations are associated with positive earnings growth in the periods following a missing record and with negative earnings growth in the periods preceding a missing earnings record, implying that these individual realizations of residual earnings are not random draws from the earnings distribution. As pointed out by Moffitt and Gottschalk (2012), little is known about the effect of attrition on the autocovariance function of earnings and, therefore, on the estimates of the earnings process. Our results indicate that the effect can be large.

In columns (5)–(8) of Table 4, we proceed to explore the *volatility* of idiosyncratic earnings at the start and end of earnings spells. The size of squared residual earnings is mechanically higher in the few first and last earnings records since, as we have just seen, residual earnings are more negative, on average, in those periods. To remove the influence of this mechanical effect, we take the (individually demeaned) residuals from the regressions of columns (1)–(4) and then square them. In the German data, the overall mean of squared residual earnings is about 0.15 in both unbalanced samples, while in the Danish data, the corresponding mean in both samples is 0.11. The results imply that earnings are significantly more volatile in the first and last years of individual spells. For example, in the German data, the mean of squared residual earnings in the first year is about 153% ($100 \times 0.23/0.15$) larger than the typical size measured by the mean of squared residual earnings in the sample. In the German consecutive sample, about 23% of individuals have their first earnings record after 1984, the first calendar year of the sample, and about 31% of individuals have their last record before 2008, the last year of the sample. The same numbers for Danish data are 18% and 22%, respectively. This is a nontrivial number of individuals with pronounced differences in the level and volatility of residual earnings in the few first and last periods of earnings spells. In the nonconsecutive samples, earnings in the periods preceding and following interior missing earnings records are also highly volatile. In the German data, for instance, the volatility of earnings observations one year before an interior missing record is about 100% ($100 \times 0.15/0.15$) larger than the volatility of typical earnings observations. However, the interior missing observations are much less prevalent than missing observations at the start and end of earnings spells. Tables 1 and 2 indicate that less than 1.5% of observations are missing on the interior of the nonconsecutive samples in our Danish and German data.²²

²²Browning and Ejrnæs (2013) argue that the first-stage regression (removing predictable variation in earnings due to observables such as time effects, age, education, race, etc.) should be ideally integrated with the second-stage analysis of modeling the earnings dynamics. As a special case, they show that removing time effects in a first-stage regression could potentially distort the estimated distribution of an autoregressive parameter in a model allowing for parameters to vary across individuals. Although we follow the common practice of separating the first-stage regression from the second stage and focus on the canonical earnings process whose parameters are shared by all individuals, we have obtained very similar results when performing the analysis of Table 4 on raw earnings using the regressions that do not control for education and predictable time variation; see Table A-4.

3.4.2 Forces behind low and volatile earnings surrounding missing records

Finally, we consider some of the forces leading to low and volatile earnings at the start and end of the earnings spells. One obvious explanation is based on the fact that the data on earnings are typically recorded at an annual frequency. An individual who is, say, entering the sample for the first time is (statistically) expected to enter in the middle of the year but may enter at any point throughout the year. Thus, earnings in that year are expected to be lower and have a larger variance than interior earnings observations from contiguous earnings histories. We can assess this conjecture using our German data, which contain information on the number of days worked on all jobs and the average daily wage from all jobs held during a year. We use these data to decompose earnings cuts in the years around missing earnings records due to a reduction in days worked and wages. As can be seen from Table 5, most of the reduction in earnings in the first or last year of the earnings spell is due to the reduction in days worked.²³ The reduction in wages in those years is nontrivial as well but is relatively less important in inducing earnings fluctuations.^{24,25}

In Table A-6, we report that years at the start and at the end of earnings spells are associated with a somewhat elevated probability of occupation (three-digit), industry (two-digit), and employer change. Individuals continue experiencing elevated mobility rates in the second and third years of incomplete earnings spells, as well as in a few years prior to the end of incomplete earnings spells. Since many individuals take some time off in-between jobs, the elevated mobility rates help explain why the effect on work hours can last beyond the first and last year of the spell, although the effect in those years is much smaller.²⁶

²³The extent of the reduction in days worked is directly affected by the sample inclusion restrictions. For example, our sample includes only observations when individuals worked at least 10% of the year. In Appendix IV, we present the results of Tables 3–4 for alternative cutoffs, ranging from 0% (when we keep all nonzero earnings) to 100% (when we restrict the sample to those working the entire year). As expected, in the sample that includes all nonzero earnings, observations surrounding the missing records have even lower mean and higher variance because they now include observations with the particularly low number of days worked. In contrast, in the sample restricted to observations where individuals work the entire year, the earnings records surrounding missing observations are not very different in mean and variance to the records in the interior of the spells (because the variation in the number of days worked is eliminated). As a consequence, the difference in the estimated variances of permanent and transitory shocks using the moments in levels and differences shrinks as the restriction on the minimum number of days worked becomes tighter. In the extreme, when no variation in the number of days worked is allowed, the results on unbalanced samples become similar to those obtained for the balanced samples.

²⁴Consistent with the logic of our argument, we find a much smaller discrepancy in the estimated variances of permanent and transitory shocks when using the moments in levels and differences in *wages*, as much of the variability in earnings in our administrative datasets at the start and end of contiguous spells is due to the variability in hours. This is consistent with the observation of Krueger et al. (2010) that the discrepancy is larger for the estimates of earnings processes than wage processes in a broad cross-section of countries.

²⁵Somewhat relatedly, Hoffmann and Malacrino (2019) study the business-cycle variation in earnings and find that the tails of the distribution of the annual earnings growth are driven by the variation in weeks worked.

²⁶To define mobility, we consider how many individuals report more than one occupation, industry, or employer in a given year of the spell. The mobility measures are relatively low in the first and last years of incomplete spells because they are typically based only on the part of the year when individuals are present in our sample.

4 Quantitative Evaluation of the Importance of Observations Around Missing Ones for the Biases

In this section, we first directly verify the contribution of the low mean and high variance of earnings surrounding missing observations in the unbalanced samples to the variances of permanent and transitory shocks estimated using the theoretical moments (L1)–(L2) and (D1)–(D2). We also highlight the importance of those observations for the estimated earnings process by simply dropping them prior to estimation – an experiment that recovers the same variances of transitory and permanent shocks in levels and differences. We then proceed by augmenting the model in Eq. (1) with extra transitory shocks whose means and variances are estimated by matching the means and variances of earnings records surrounding the missing ones in the data. Since we rely on the full set of the data autocovariance moments in our estimation, as is standard in the literature, our minimum-distance estimator is overidentified and does not offer transparent mapping between the estimated parameters and data moments. We, therefore, further show in simulations replicating our unbalanced samples that our results based on the overidentified minimum-distance estimation hold for various plausible parameterizations of the extended earnings process. We close this section by offering some thoughts on the implications of our results for quantitative modeling that relies on the external estimation of the income process parameters.

4.1 Direct Evaluation of the Biases Using the Permanent-Transitory Decomposition Moments

In this section, we directly verify that irregular observations surrounding the missing ones induce most of the difference between permanent and transitory shock variances implied by identifying moments (L1)–(L2) and (D1)–(D2). This result provides evidence that the identifying moments on which our theoretical argument is based are indeed the relevant ones and largely responsible for the results of our full estimation targeting all available autocovariance moments.

As an example of computing these implied variances, we calculate an estimate of the permanent variance at time t using the identifying moment in levels (L1) as

$$\sigma_{\xi,l,t}^2 = \frac{\sum_i y_{i,t} y_{i,t+1}}{\sum_i I_{t,t+1}^i} + \frac{\sum_i y_{i,t-2} y_{i,t-1}}{\sum_i I_{t-2,t-1}^i} - \frac{\sum_i y_{i,t+1} y_{i,t-1}}{\sum_i I_{t-1,t+1}^i} - \frac{\sum_i y_{i,t} y_{i,t-2}}{\sum_i I_{t-2,t}^i}, \quad (2)$$

where the subscript l indicates that we are estimating the variance using information on log-earnings in levels, and $I_{t,t'}^i$ is an indicator function taking the value of one if individual earnings observations are nonmissing in both years t and t' , and taking the value of zero otherwise. Note that individual i will not contribute to the estimated variance of the permanent shock

at time t only if all of the earnings cross-products for that individual – $y_{it}y_{it+1}$, $y_{it-2}y_{it-1}$, $y_{it+1}y_{it-1}$, and $y_{it}y_{it-2}$ – are missing.

Let I_{it}^m be an indicator function that equals one if an individual's earnings is missing at one of the periods $t-2, \dots, t+2$ defined as $\mathbb{1}\left(\sum_{j=-2}^2 (1 - I_{it+j}) > 0\right)$, where $\mathbb{1}(\cdot)$ equals one if the expression in brackets is true and zero otherwise. We calculate the variance of permanent shocks due to irregular observations surrounding the missing earnings records, $\sigma_{\xi,l,o,t}^2$, as

$$\sigma_{\xi,l,o,t}^2 = \frac{\sum_i y_{i,t}y_{i,t+1}I_{it}^m}{\sum_i I_{t,t+1}^m} + \frac{\sum_i y_{i,t-2}y_{i,t-1}I_{it}^m}{\sum_i I_{t-2,t-1}^m} - \frac{\sum_i y_{i,t+1}y_{i,t-1}I_{it}^m}{\sum_i I_{t-1,t+1}^m} - \frac{\sum_i y_{i,t}y_{i,t-2}I_{it}^m}{\sum_i I_{t-2,t}^m}. \quad (3)$$

An estimate of the permanent variance in levels, net of the effects of irregular observations surrounding the missing ones, $\sigma_{\xi,l,n,t}^2$, can then be calculated as

$$\begin{aligned} \sigma_{\xi,l,n,t}^2 = & \frac{\sum_i y_{i,t}y_{i,t+1}(1 - I_{it}^m)}{\sum_i I_{t,t+1}^m(1 - I_{it}^m)} + \frac{\sum_i y_{i,t-2}y_{i,t-1}(1 - I_{it}^m)}{\sum_i I_{t-2,t-1}^m(1 - I_{it}^m)} \\ & - \frac{\sum_i y_{i,t+1}y_{i,t-1}(1 - I_{it}^m)}{\sum_i I_{t-1,t+1}^m(1 - I_{it}^m)} - \frac{\sum_i y_{i,t}y_{i,t-2}(1 - I_{it}^m)}{\sum_i I_{t-2,t}^m(1 - I_{it}^m)}. \end{aligned} \quad (4)$$

We can similarly define the variances of permanent and transitory shocks in levels and differences for the consecutive unbalanced panels – e.g., the permanent variance utilizing all sample information ($\sigma_{\xi,l,t}^2$ for levels and $\sigma_{\xi,d,t}^2$ for differences), the permanent variance due to irregular observations in the first and last few periods of an individual's earnings spell ($\sigma_{\xi,l,o,t}^2$ and $\sigma_{\xi,d,o,t}^2$), and the permanent variance net of their effects ($\sigma_{\xi,l,n,t}^2$ and $\sigma_{\xi,d,n,t}^2$).

We present the estimates of those variances, averaged across all sample years, for both datasets in Table 6. For the German data, in the consecutive sample, the estimates of the variance of permanent shocks in levels and differences using all sample information are 0.013 and 0.024, respectively.²⁷ Net of the effect of observations surrounding the missing ones, the estimated variances are $\hat{\sigma}_{\xi,l,n}^2 = 0.010$ in levels and $\hat{\sigma}_{\xi,d,n}^2 = 0.010$ in differences. The unadjusted variances of transitory shocks in levels and differences are estimated at 0.020 and 0.008, respectively, while the variances net of outliers in levels and differences are both estimated at 0.007. The results for the Danish data are qualitatively similar. Clearly, the discrepancy between the estimates of permanent and transitory shock variances in levels and differences is virtually eliminated when netting out the effects of irregular observations surrounding the missing ones.

Similarly, in the German nonconsecutive sample, the variances of permanent shocks are $\sigma_{\xi,l}^2 = 0.0096$, $\sigma_{\xi,l,n}^2 = 0.0097$, $\sigma_{\xi,d}^2 = 0.018$, $\sigma_{\xi,d,n}^2 = 0.0097$, while the variances of transitory

²⁷The estimates deviate from the values in Table 3 because we do not impose the exact permanent-transitory decomposition on the data in the minimum-distance estimation of Table 3. The difference between the estimated variance of permanent shocks in levels and differences is not as drastic as in Table 3 because the estimated persistence of the permanent shocks in levels is lower than in differences in the minimum-distance estimation.

shocks are $\sigma_{\epsilon,l}^2 = 0.018$, $\sigma_{\epsilon,l,n}^2 = 0.007$, $\sigma_{\epsilon,d}^2 = 0.007$, $\sigma_{\epsilon,d,n}^2 = 0.007$. Netting out the influence of missing observations and the influence of the first and last records in the earnings spells eliminates most of the discrepancy between the variances of permanent and transitory shocks in differences and levels.

4.2 Restricting Unbalanced Samples

One approach to eliminating the impact of low mean and high variance of observations at the start and end of earnings spells on the estimates of the permanent/transitory decomposition is to simply drop those observations. Accordingly, in Tables 7 and 8, columns (3) and (4) of Panel A, we repeat our analysis of Table 3 using the German and Danish samples with nine or more consecutive observations after dropping the first three observations for individuals whose earnings spells start after the first sample year and the last three observations for individuals whose earnings spells end before the last sample year.²⁸ In the same columns in Panel B we, in addition, drop three observations before and after a missing earnings record in the nonconsecutive samples.

For the sample with nine or more consecutive observations, doing so barely affects the persistence of permanent shocks, although their variance estimated using the moments in differences is reduced by about 70%. The variance of transitory shocks estimated using the moments in levels is reduced by about 60%. We observe a similar pattern in the nonconsecutive sample in Panel B. As a result, the estimated earnings process is virtually identical in estimations utilizing the moments for growth rates and levels.

A comparison with Table 3 also indicates, consistent with the analysis in Section 2, that in both datasets the variance of the permanent component is more robustly estimated using the moments in levels, whereas the variance of the transitory component, net of the transitory variation in earnings due to ν -shocks, is more robustly estimated using the moments in differences.

4.2.1 Validation Using Simulated Data

We now simulate artificial earnings panels, consistent with the properties of the consecutive and nonconsecutive samples in our German data. In this experiment, we know the true parameters of the earnings process, which allows us to assess the performance of the proposed

²⁸Geweke and Keane (2000) drop the first and Baker and Solon (2003) drop the first and last earnings records in their analysis because individuals are likely to work only part of the year the first and last time they are observed, which is consistent with our results. Dropping the first and last records, however, may not be enough to eliminate the biases in the estimated variances of permanent and transitory shocks in levels and differences as the records in a few subsequent or preceding years may still be somewhat different from the rest of the earnings observations. In this case, dropping, say, the first observation only leads to the second irregular observation being next to the missing one and induces the associated biases.

empirical methods and procedures in recovering the true values of these parameters. Specifically, we verify that: (1) the observed low mean and high variance of observations surrounding the missing ones induces a large discrepancy in the estimated variances of shocks when targeting the moments in levels and differences, (2) the variance of permanent (transitory) shocks is more robustly estimated using the moments in levels (differences), and (3) dropping observations surrounding the missing ones reconciles the estimates when using the moments in differences and levels and yields accurate estimates of the permanent and transitory shocks.

We replicate our German unbalanced samples in terms of the number of person-year observations and the age distribution and assume that incomes in the spells starting (ending) in the years other than the first (last) year of the sample are, in addition, affected by the ν -shocks. For the consecutive sample, we assume that the persistence of the permanent component is 0.980, the variance of permanent shocks is 0.008, the persistence of the transitory component is 0.170, the variance of transitory shocks is 0.010, and the variance of fixed effects is 0.025. These values are similar to the estimates of the transitory component using the moments in growth rates and of the permanent component using the moments in levels in Table 3, columns (1)–(2). We assume that the shocks and fixed effects are drawn from Student t -distributions with four degrees of freedom, since our samples have high excess kurtosis.²⁹ We take the means and variances of the ν -shocks in the first three and last three periods from columns (1) and (5) of Table 4. Following the same strategy for the nonconsecutive sample, we assume that the persistence of the permanent component is 0.999, the variance of permanent shocks is 0.005, the persistence of the transitory component is 0.20, the variance of transitory shocks is 0.01, and the variance of fixed effects is 0.025. This is in line with the estimated permanent component in column (5) and transitory component in column (6) of Table 3. We take the means and variances of ν -shocks in the first three and last three periods, and three years before and three years after missing earnings records, from columns (3) and (7) of Table 4, and assume that the shocks follow the moving-average structure of order one with the persistence equal to 0.20.

The results for estimations fitting the entire set of autocovariances of the (simulated) data in levels or growth rates are in Table 9. The results are averages across 100 simulations. Utilizing the full sample results in an overestimation of the variance of the permanent (transitory) shock in differences (levels) as we observed in the actual data in Table 3. The permanent component is recovered fairly well utilizing the moments in levels, while the transitory component is closer to the truth utilizing the moments in differences.³⁰ Dropping the three observations adjacent to missing records aligns the results of estimations in levels and

²⁹Assuming normal shocks instead has no impact on our findings. The choice of degrees of freedom for a Student t -distribution of the shocks is consistent with the data; see footnote (33).

³⁰We present additional simulations in Appendix VI, in which we vary the size of the true variance and persistence of permanent and transitory shocks and find that this pattern remains robust.

differences and correctly recovers the parameters of the underlying earnings process.

4.3 Modeling Earnings Records Surrounding the Missing Ones

In columns (5)–(8) of Tables 7 and 8, instead of dropping observations surrounding the missing ones, we estimate an extended earnings process that explicitly models them.³¹ Specifically, we estimate the following model:

$$\begin{aligned}
y_{it} &= \alpha_i + p_{it} + \tau_{it} + \chi_{it}, \quad t = t_0, \dots, T \\
p_{it} &= \phi_p p_{it-1} + \xi_{it} \\
\tau_{it} &= \epsilon_{it} + \theta \epsilon_{it-1} \\
\chi_{it+j} &= \begin{cases} \nu_{it} & \text{if } y_{it-1} \text{ or } y_{it+1} \text{ is missing and } t-1 \geq t_0, t+1 \leq T, j=0 \\ \theta \nu_{it} & j=1 \\ 0 & \text{otherwise,} \end{cases}
\end{aligned} \tag{5}$$

Columns (5) and (6) of Tables 7 and 8 contain the estimates of this earnings process, where we model the means and variances of the immediate observations around missing records by matching the regression coefficients from Table 4. We assume that ν_{it} 's are drawn from distributions with means and variances that depend on whether an individual has missing observations in the interior of a nonconsecutive earnings spell, at the beginning of a consecutive earnings spell, or at the end of a consecutive earnings spell (the corresponding means and variances have superscripts m , f , and l , respectively, in Appendix Table A-11, which contains full estimation results). We further assume that the persistence of the shock ν_{it} , θ , is the same as the persistence of the ϵ_{it} shock because the estimated persistence of the transitory component barely changes when we drop observations surrounding the missing ones, which can be verified by comparing the results in columns (3)–(4) with the results in columns (1)–(2).³² As before, the other moments used for estimating the model are the autocovariance moments in either levels or differences. We rely on the simulated minimum distance method, assuming that observations are missing at random, all of the innovations are i.i.d. normal, and utilizing the optimal weighting matrix estimated by block-bootstrap.³³ As expected,

³¹An alternative exposition of this earnings process, cast in terms of a selection model, is presented in Online Appendix VIII.

³²We allow the persistence for symmetry with the modeling of the transitory shock, but imposing that it is zero (or another plausible value) does not affect the results.

³³Because it is well known that earnings shocks are non-Gaussian, we have also tried estimations which assume that the shocks are drawn from a Student t-distribution with the degrees of freedom estimated from the data by matching kurtosis of the growth in earnings observed in the data. We found that the point estimates in Table A-11 were virtually the same (with the estimated degrees of freedom of the Student t-distribution equal to about 4, implying a leptokurtic distribution of the shocks). This is not surprising since the discrepancy in the estimated variances is the feature of the second moments of the data – and not the higher-order moments – as is highlighted in Eq. (D1)–(L2).

estimating the extended earnings process results in a substantial reduction of the estimated variance of permanent shocks in differences and transitory shocks in levels.^{34,35}

4.4 Implications

The experiments described in this section suggest that no other mechanisms but the presence of irregular earnings observations around missing earnings records are responsible for the discrepancy in the estimated variances of the shocks in levels and differences in our administrative earnings data. The practical implications of these findings depend on the objective of the analysis. If one is interested in the properties of permanent and transitory components as well as in the detailed analysis of earnings at the start and end of employment spells, one can estimate the extended process in Eq. (5) where the mean and the variance of the shocks at the beginning and the end of contiguous earnings histories are readily identified from the mean and the variance of earnings in those periods. Many macro models are too stylized, however, to incorporate explicit treatment of these observations.³⁶ They use as an input only the permanent/transitory components of the earnings process, as in Eq. (1), so that it becomes crucial to obtain the correct estimates of the stochastic properties of these components in the data. These components can be estimated using the extended earnings process in Eq. (5), although simpler alternatives are also available. First, we have shown theoretically and verified empirically that the variance of the transitory shock is estimated with no bias when estimation is based on the moments for earnings growth rates if the ν -shock is not serially correlated and the variance of the permanent shock is recovered well when estimation targets the moments in levels. One could, therefore, use the estimated permanent component from targeting the moments in levels and the estimated transitory component from targeting the moments in

³⁴It is clear from the precision of our estimates that allowing for the means and variances of ν -shocks is not redundant in fitting the data moments; e.g., the quasi-likelihood ratio test's p-values for excluding $\sigma_{\nu_t^f}^2$, $\sigma_{\nu_t^l}^2$, $\mu_{\nu_t^f}$, and $\mu_{\nu_t^l}$ in the estimation of column (1) Table A-11 are all well below 1%.

³⁵In the previous version of the paper, we allowed for modeling of three observations adjacent to the missing ones, with little difference to our results on the variances of permanent and transitory earnings shocks and their persistence. This is not surprising. When estimating the extended earnings process, it is only essential to account for the observations immediately adjacent to the missing ones. The other irregular – but much less different – observations will be subsumed in the estimated variance of transitory shocks (as the data include earnings before and after those earnings records that allows to detect mean-reversion of those shocks if the properties of the observation immediately adjacent to the missing one are controlled for).

³⁶The approach in the macro literature can be partially justified by noting that these properties of observations adjacent to the missing ones are largely induced by the measurement timing and frequency in the data that is not actually relevant in the model. To see this, consider an example. An individual graduates from college and starts working on July 1 of some year, works at the same job making \$100 a day for 40 years, and retires on June 31. If a researcher could observe daily earnings, she would conclude that this individual faces no risk to earnings. Even if the data were annual, but the year was defined to run from July 1 to June 31, a researcher will reach the same conclusion. But if a researcher can only observe annual income based on a calendar year running from January 1 to December 31, she will conclude that there is high risk in the first and last year when the individual is observed with positive earnings. Clearly, this “risk” is just an artifact of available data and not the risk that needs to be modeled.

growth rates. An alternative approach is to estimate the earnings process in Eq. (1) on the data that do not include the observations surrounding the missing ones. As we have shown, this solution recovers the true parameters of this abbreviated process quite well.

In Appendix VIII, we assess the performance of these approaches using the standard calibrated incomplete markets model. Specifically, we first calibrate a model using the canonical earnings process with permanent and transitory shocks extended to include missing observations and low mean/high variance of observations surrounding the missing ones. The parameters of this process come from our estimates based on the German administrative data. We then simulate a dataset from this model and treat it as if true data available to a researcher. We ask this hypothetical researcher first to estimate the canonical earnings process (that includes only permanent and transitory shocks) using the standard approach that does not consider the presence of low-mean and high-variance observations surrounding the missing ones. Next, we ask her to use the approach proposed in this paper. After that, the researcher is asked to use the two estimated earnings processes as inputs into a standard incomplete markets model that does not feature missing observations. Finally, we ask which of the two calibrated models serves as a good guide to the objects of interest in the original model that generated the data. We find that if the researcher followed the strategy proposed in this paper by measuring the variances of permanent and transitory shocks while accounting for the properties of observations next to the missing ones, she would have reached largely correct conclusions. In contrast, had she followed the standard practice in the literature and ignored the features of earnings observations surrounding the missing ones when estimating the earnings process, some of her conclusions would be grossly erroneous. This exercise suggests that obtaining correct variances of permanent and transitory shocks in the data is of first-order importance while the loss from not modeling missing observations and the properties of observations surrounding them in incomplete markets models is much smaller.

We, of course, recognize that ignoring missing observations in quantitative models is a shortcut. A precise assessment of the amount of risk that individuals face does depend on what happens to earnings when they are missing in the data. However, it seems very difficult to know this because of the great variety of reasons for why earnings could be missing in administrative data, as we discuss in Appendix II. Depending on the reason, true unobserved earnings may plausibly be the same or higher/lower than the observed ones, and one would have to separately model these pathways into and out of missing observations (along the lines of, e.g., Altonji et al. 2013). This would, in theory, allow one to obtain more precise estimates for the variances of permanent and transitory shocks to earnings and a more accurate assessment in a quantitative model. Even if we knew why a particular observation is missing, such an exercise would be extremely challenging.

We do not expect the literature to overcome these challenges in the foreseeable future. Until that happens, however, we offer a constructive approach that follows the literature in utilizing

in estimation all computable individual-level moments that do not involve missing values and not modeling missing observations in quantitative incomplete markets models. Researchers taking this currently unavoidable shortcut need to be careful. As we show, ignoring the unique properties of earnings observations around the missing observations would lead to severely biased estimates of the variances of permanent and transitory shocks in growth rates and levels and erroneous inference in quantitative models that use these shock variances as an input.

5 Conclusion

Properties of the earnings process play an important role in various areas of macro and labor economics. Different specifications of this process have been explored in the literature. The most widely used specification is based on decomposing earnings into the sum of persistent and transitory components and often assumes that the persistent component follows a random walk. The parameters of such a process can be identified using the moments based on earnings growth rates (first-difference in log earnings) or the moments based on log earnings levels. Historically, the former approach is more common in labor economics, while the latter is more common in the macroeconomics literature. Unfortunately, these two approaches lead to dramatically different estimates of the variances of permanent and transitory components. In particular, using the same set of observations in the data, the variance of persistent shocks is typically estimated to be much higher when the moments in growth rates are targeted, while the variance of transitory shocks is found to be much higher when the estimation is based on fitting the moments in levels. This discrepancy has important implications for substantive economic analysis. For example, the earnings process drives the heterogeneity in Bewley-type models with incomplete markets, and the variances of earnings components determine not only economic choices, such as consumption and savings, but also the optimal design of policies, such as taxes and transfers. Moreover, the standard approach to estimating the amount of insurance that individuals have against permanent and transitory shocks in the data relies on the estimated variances of permanent and transitory components. The uncertainty about the size of these variances translates into uncertainty regarding the right amount of insurance that should be generated by the widely used incomplete markets models and the associated uncertainty about the results of welfare analyses using those models.

In this paper, we uncovered the data features that can quantitatively account for the large difference in the estimates based on earnings growth rates and levels, at least in the administrative data from Denmark and Germany. In particular, we found that earnings are lower on average and more volatile at the start and end of continuous earnings spells. We have shown theoretically that these irregular earnings observations, which are either preceded

or followed by a missing observation, induce an upward bias in the estimates of the variance of permanent shocks based on the moments in differences and the variance of transitory shocks when estimation is based on the moments in levels. Thus, even when working with very large administrative datasets with highly reliable information, one must remain vigilant because natural features of the datasets, like the low mean and high variance of earnings at the start and end of earnings spells, can induce extremely large biases in the estimated earnings processes.

While the primary focus of this paper is on estimating earnings processes on large administrative datasets that are becoming central in the literature, the mechanism we describe also applies to survey-based PSID data on hourly wages, individual earnings, combined earnings of husbands and wives, and net family income. This is important, in part, because the PSID data are the primary source of information on the extent of consumption insurance achieved by U.S. households and on the sources of insurance, such as household savings and borrowing or the public transfer and tax system. Hryshko and Manovskii (2018) find that correcting the measurement of family income dynamics for the effects of the low mean and high variance of the observations surrounding the missing ones as proposed in this paper leads to a very different assessment of the degree and sources of consumption insurance relative to the key empirical benchmark for incomplete markets models provided by Blundell et al. (2008).

Our punchline is that there clearly must be a tight correspondence between the components of the earnings process used in the model and the ones estimated in the data. If the model utilizes a version of the permanent/transitory process as in Eq. (1), the measurement of these permanent/transitory components in the data must not be biased due to the special nature of observations adjacent to the missing ones. If, on the other hand, one wishes to incorporate the features of these observations into the model, one needs to estimate the extended process in Eq. (5) in the data and use it as an input into the model. The paper provides a way to implement both of these approaches.

References

- ABOWD, J. AND D. CARD (1989): “On the Covariance Structure of Earnings and Hours Changes,” *Econometrica*, 57, 411–445.
- ALTONJI, J. G., J. ANTHONY A. SMITH, AND I. VIDANGOS (2013): “Modeling Earnings Dynamics,” *Econometrica*, 81, 1395–1454.
- ARELLANO, M. AND B. HONORÉ (2001): “Panel Data Models: Some Recent Developments,” in *Handbook of Econometrics*, Vol. 5, ed. by J. J. Heckman and E. E. Leamer, Amsterdam: Elsevier, chap. 53, 3229–3296.
- ATTANASIO, O., E. HURST, AND L. PISTAFERRI (2012): “The Evolution of Income, Consumption, and Leisure Inequality in the US, 1980–2010,” NBER Working Paper # 17982.
- BAKER, M. AND G. SOLON (2003): “Earnings Dynamics and Inequality among Canadian Men, 1976–1992: Evidence from Longitudinal Income Tax Records,” *Journal of Labor Economics*, 21, 289–321.
- BANKS, J. AND P. DIAMOND (2010): “The Base for Direct Taxation,” in *Dimensions of Tax Design: The Mirrlees Review*, ed. by The Institute for Fiscal Studies (IFS), Oxford: Oxford University Press for The Institute for Fiscal Studies, chap. 6, 548–648.
- BLUNDELL, R., M. GRABER, AND M. MOGSTAD (2015): “Labor Income Dynamics and the Insurance from Taxes, Transfers, and the Family,” *Journal of Public Economics*, 127, 58–73.
- BLUNDELL, R., L. PISTAFERRI, AND I. PRESTON (2008): “Consumption Inequality and Partial Insurance,” *American Economic Review*, 98, 1887–1921.
- BROWNING, M. AND M. EJRNÆS (2013): “Heterogeneity in the Dynamics of Labor Earnings,” *Annual Review of Economics*, 5, 219–245.
- BROWNING, M., M. EJRNÆS, AND J. ALVAREZ (2010): “Modelling Income Processes with Lots of Heterogeneity,” *Review of Economic Studies*, 77, 1353–1381.
- CARROLL, C. D. (1997): “Buffer-Stock Saving and the Life Cycle/Permanent Income Hypothesis,” *Quarterly Journal of Economics*, 112, 1–55.
- CASTAÑEDA, A., J. DÍAZ-GIMÉNEZ, AND J.-V. RÍOS-RULL (2003): “Accounting for the U.S. Earnings and Wealth Inequality,” *Journal of Political Economy*, 111, 818–857.
- DEATON, A. (1991): “Saving and Liquidity Constraints,” *Econometrica*, 59, 1121–1248.
- DEBACKER, J., B. HEIM, V. PANOUSI, S. RAMNATH, AND I. VIDANGOS (2013): “Rising Inequality: Transitory or Persistent? New Evidence from a Panel of U.S. Tax Returns,” *Brookings Papers on Economic Activity*, 67–142.

- DEMYANYK, Y., D. HRYSHKO, M. J. LUENGO-PRADO, AND B. E. SØRENSEN (2017): “Moving to a Job: The Role of Home Equity, Debt, and Access to Credit,” *American Economic Journal: Macroeconomics*, 9, 149–181.
- DOMELIJ, D. AND M. FLODÉN (2010): “Inequality Trends in Sweden 1978–2004,” *Review of Economic Dynamics*, 13, 179–208.
- DUSTMANN, C., J. LUDSTECK, AND U. SCHÖNBERG (2009): “Revisiting the German Wage Structure,” *Quarterly Journal of Economics*, 124, 843–881.
- FARHI, E. AND I. WERNING (2012): “Insurance and Taxation over the Life Cycle,” *Review of Economic Studies*, 80, 596–635.
- FITZGERALD, J., P. GOTTSCHALK, AND R. MOFFITT (1998): “An Analysis of Sample Attrition in Panel Data: The Michigan Panel Study of Income Dynamics,” *Journal of Human Resources*, 33, 251–299.
- FRIEDMAN, M. (1957): *A Theory of the Consumption Function*, Princeton, NJ: Princeton University Press.
- FRIEDMAN, M. AND S. KUZNETS (1954): *Income from Independent Professional Practice*, New York: National Bureau of Economic Research.
- GEWEKE, J. AND M. KEANE (2000): “An Empirical Analysis of Earnings Dynamics Among Men in the PSID: 1968–1989,” *Journal of Econometrics*, 96, 293–356.
- GOTTSCHALK, P. AND R. A. MOFFITT (1994): “The Growth of Earnings Instability in the U.S. Labor Market,” *Brookings Papers on Economic Activity*, 2, 217–54.
- GUVENEN, F. (2009): “An Empirical Investigation of Labor Income Processes,” *Review of Economic Dynamics*, 12, 58–79.
- GUVENEN, F., S. OZCAN, AND J. SONG (2014): “The Nature of Countercyclical Income Risk,” *Journal of Political Economy*, 122, 621–660.
- HAGEDORN, M., J. LUO, I. MANOVSKII, AND K. MITMAN (2019a): “Forward Guidance,” *Journal of Monetary Economics*, 102, 1–23.
- HAGEDORN, M., I. MANOVSKII, AND K. MITMAN (2019b): “The Fiscal Multiplier,” National Bureau of Economic Research, Working Paper 25571.
- HEATHCOTE, J. (2005): “Fiscal Policy with Heterogeneous Agents and Incomplete Markets,” *Review of Economic Studies*, 72, 161–188.
- HEATHCOTE, J., F. PERRI, AND G. L. VIOLANTE (2010a): “Unequal We Stand: An Empirical Analysis of Economic Inequality in the United States: 1967–2006,” *Review of Economic Dynamics*, 13, 15–51.

- HEATHCOTE, J., K. STORESLETTEN, AND G. L. VIOLANTE (2010b): “The Macroeconomic Implications of Rising Wage Inequality in the United States,” *Journal of Political Economy*, 118, 681–722.
- (2014): “Consumption and Labor Supply with Partial Insurance: An Analytical Framework,” *American Economic Review*, 104, 2075–2126.
- HOFFMANN, E. B. AND D. MALACRINO (2019): “Employment Time and the Cyclical Growth of Earnings,” *Journal of Public Economics*, 169, 160–171.
- HRYSKO, D. (2012): “Labor Income Profiles Are Not Heterogeneous: Evidence From Income Growth Rates,” *Quantitative Economics*, 3, 177–209.
- HRYSKO, D. AND I. MANOVSKII (2017): “How Much Consumption Insurance in the U.S.?” mimeo, University of Alberta and University of Pennsylvania.
- (2018): “Income Dynamics and Consumption Insurance,” mimeo, University of Alberta and University of Pennsylvania.
- HUBBARD, R. G., J. SKINNER, AND S. P. ZELDES (1995): “Precautionary Saving and Social Insurance,” *Journal of Political Economy*, 103, 360–399.
- HUGGETT, M., G. VENTURA, AND A. YARON (2011): “Sources of Lifetime Inequality,” *American Economic Review*, 101, 2923–2954.
- KAPLAN, G., B. MOLL, AND G. VIOLANTE (2018): “Monetary Policy According to HANK,” *American Economic Review*, 108, 697–743.
- KAPLAN, G. AND G. L. VIOLANTE (2010): “How Much Consumption Insurance Beyond Self-Insurance?” *American Economic Journal: Macroeconomics*, 2, 53–87.
- KRUEGER, D. AND F. PERRI (2006): “Does Income Inequality Lead to Consumption Inequality? Evidence and Theory,” *Review of Economic Studies*, 73, 163–193.
- KRUEGER, D., F. PERRI, L. PISTAFERRI, AND G. L. VIOLANTE (2010): “Cross-sectional Facts for Macroeconomists,” *Review of Economic Dynamics*, 13, 1–14.
- LIVSHITS, I., J. MACGEE, AND M. TERTILT (2007): “Consumer Bankruptcy: A Fresh Start,” *American Economic Review*, 97, 402–418.
- MACURDY, T. E. (1982): “The Use of Time Series Processes to Model the Error Structure of Earnings in a Longitudinal Data Analysis,” *Journal of Econometrics*, 18, 83–114.
- (2007): “A Practitioner’s Approach to Estimating Intertemporal Relationships Using Longitudinal Data: Lessons from Applications in Wage Dynamics,” in *Handbook of Econometrics*, Vol. 6A, ed. by J. J. Heckman and E. E. Leamer, Amsterdam: Elsevier, chap. 62, 4057–4167.

- MEGHIR, C. AND L. PISTAFERRI (2004): “Income Variance Dynamics and Heterogeneity,” *Econometrica*, 72, 1–32.
- (2011): “Earnings, Consumption, and Life Cycle Choices,” in *Handbook of Labor Economics*, Vol. 4B, ed. by D. Card and O. Ashenfelter, Amsterdam: Elsevier, chap. 9, 773–854.
- MOFFITT, R. A. AND P. GOTTSCHALK (2012): “Trends in the Transitory Variance of Male Earnings: Methods and Evidence,” *Journal of Human Resources*, 47, 204–236.
- POSTEL-VINAY, F. AND H. TURON (2010): “On-the-Job Search, Productivity Shocks, and the Individual Earnings Process,” *International Economic Review*, 51, 599–629.
- STORESLETTEN, K., C. TELMER, AND A. YARON (2001): “How Important Are Idiosyncratic Shocks? Evidence from Labor Supply,” *American Economic Review*, 91, 413–417.

TABLE 1: DANISH DATA, 1981–2006. SUMMARY STATISTICS FOR SELECTED YEARS.

	9 consec.	20 not nec. consec.	Balanced
Number of individuals	102,825	90,668	67,008
Number of observations	2,367,552	2,298,429	1,742,208
Education			
Less than high school	0.227	0.222	0.206
High school degree	0.032	0.031	0.029
Vocational training	0.505	0.521	0.542
Two-year university degree	0.046	0.046	0.047
Bachelors degree	0.125	0.122	0.124
Master or Ph.D.	0.065	0.059	0.051
Earnings			
1985	40,157 (12,831)	40,227 (12,889)	41,383 (12,278)
1995	48,197 (20,562)	48,444 (20,462)	50,004 (19,954)
2005	52,656 (26,635)	51,511 (26,279)	53,298 (25,917)
Spell counts			
Start 1981, end 2006	67,008	80,787	67008
Start after 1981, end 2006	13,439	4,376	0
Start in 1981, end before 2006	17,723	5,210	0
Start after 1981, end before 2006	4,655	295	0
Total	102,825	90,668	67008
Number of spells with 20 or more not nec. consec. observations, by length			
[Proportion of missing observations within spell in square brackets]			
20		1,634 [0.144]	
21		2,009 [0.119]	
22		2,665 [0.096]	
23		3,296 [0.079]	
24		4,486 [0.054]	
25		9,570 [0.030]	
26		67,008 [0.00]	

Notes: Earnings are expressed in 2005 Euros; the standard deviation of earnings is given in parentheses. See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. “Balanced” sample contains individual earnings spells that cover the entire period.

TABLE 2: GERMAN DATA, 1984–2008. SUMMARY STATISTICS FOR SELECTED YEARS.

	9 consec.	20 not nec. consec.	Balanced
Number of individuals	18,130	13,635	9,452
Number of observations	379,080	330,748	236,300
Education			
Middle school or no degree	0.05	0.04	0.04
Vocational training	0.72	0.74	0.76
High school degree	0.06	0.05	0.05
College	0.17	0.17	0.15
Earnings			
1985	33,626 (15,876)	33,930 (13,323)	34,559 (12,881)
1995	45,309 (24,702)	47,180 (24,295)	47,965 (24,463)
2005	49,121 (36,473)	51,289 (37,106)	52,457 (37,666)
Spell counts			
Start 1984, end 2008	9,452	11,179	9,452
Start after 1984, end 2008	3,136	1,007	0
Start in 1984, end before 2008	4,463	1,393	0
Start after 1984, end before 2008	1,079	56	0
Total	18,130	13,635	9,452
Number of spells with 20 or more not nec. consec. observations, by length [Proportion of missing observations within spell in square brackets]			
20		575 [0.054]	
21		509 [0.054]	
22		623 [0.05]	
23		871 [0.037]	
24		1,605 [0.027]	
25		9,452 [0.00]	

Notes: Earnings are expressed in 2005 Euros; the standard deviation of earnings is given in parentheses. See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. “Balanced” sample contains individual earnings spells that cover the entire period.

TABLE 3: ESTIMATES OF THE EARNINGS PROCESS IN ADMINISTRATIVE DATA.

	9 consec.				20 not nec. consec.				Balanced sample			
	German data		Danish data		German data		Danish data		German data		Danish data	
	Levs. (1)	Diff. (2)	Levs. (3)	Diff. (4)	Levs. (5)	Diff. (6)	Levs. (7)	Diff. (8)	Levs. (9)	Diff. (10)	Levs. (11)	Diff. (12)
$\hat{\phi}_p$	0.976 (0.001)	0.992 (0.0008)	0.955 (0.0008)	0.987 (0.0004)	0.999 (0.001)	0.991 (0.001)	0.964 (0.0007)	0.982 (0.0006)	1 (0.001)	0.998 (0.002)	0.969 (0.0007)	0.970 (0.0009)
$\hat{\sigma}_\xi^2$	0.008 (0.0002)	0.019 (0.0003)	0.008 (0.0001)	0.013 (0.0001)	0.0048 (0.0001)	0.009 (0.0002)	0.007 (0.0001)	0.012 (0.0001)	0.0031 (0.0001)	0.0033 (0.0001)	0.005 (0.00004)	0.005 (0.00005)
$\hat{\theta}$	0.129 (0.005)	0.153 (0.009)	0.204 (0.002)	0.209 (0.003)	0.119 (0.008)	0.192 (0.008)	0.137 (0.003)	0.217 (0.003)	0.278 (0.011)	0.258 (0.012)	0.212 (0.003)	0.209 (0.003)
$\hat{\sigma}_\epsilon^2$	0.024 (0.0003)	0.009 (0.0002)	0.019 (0.0001)	0.012 (0.0001)	0.016 (0.0003)	0.009 (0.0003)	0.022 (0.0002)	0.013 (0.0001)	0.008 (0.0002)	0.0078 (0.0002)	0.009 (0.0001)	0.009 (0.0001)
$\hat{\sigma}_\alpha^2$	0.024 (0.002)	—	0.020 (0.0004)	—	0.027 (0.002)	—	0.023 (0.0004)	—	0.024 (0.001)	—	0.017 (0.0004)	—
χ^2 (d.f.)	1024.43 320	725.21 296	7104.42 346	4527.83 321	1562.51 320	1393.16 296	7002.05 346	5836.32 321	1205.84 320	935.52 296	6950.05 346	5855.67 321

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. “Balanced” sample contains individual earnings spells that cover the entire period. The estimated earnings process is: $y_{it} = \alpha_i + p_{it} + \tau_{it}$, where $p_{it+1} = \phi_p p_{it} + \xi_{it+1} + \theta \epsilon_{it}$. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. German data span the period 1984–2008, while Danish data span the period 1981–2006.

TABLE 4: RESIDUAL EARNINGS AND SQUARED RESIDUAL EARNINGS SURROUNDING MISSING OBSERVATIONS.

	Means			Variances			
	9 consec.	20 not nec. consec.	9 consec.	20 not nec. consec.	9 consec.	20 not nec. consec.	9 consec.
	German data (1)	Danish data (2)	German data (3)	Danish data (4)	German data (5)	Danish data (6)	German data (7)
Year observed: first	-0.57*** (-64.24)	-0.48*** (-119.58)	-0.65*** (-34.49)	-0.47*** (-58.12)	0.23*** (41.80)	0.19*** (70.08)	0.29*** (21.12)
Year observed: second	-0.11*** (-22.57)	-0.17*** (-54.50)	-0.14*** (-11.38)	-0.19*** (-27.24)	0.04*** (11.54)	0.08*** (34.16)	0.09*** (7.91)
Year observed: third	-0.07*** (-16.71)	-0.09*** (-36.26)	-0.09*** (-8.82)	-0.10*** (-16.87)	0.02*** (7.93)	0.04*** (21.78)	0.05*** (5.80)
Year observed: second-to-last	-0.03*** (-8.34)	-0.05*** (-22.93)	-0.05*** (-6.71)	-0.08*** (-15.56)	0.01*** (5.62)	0.02*** (11.23)	0.02*** (4.31)
Year observed: next-to-last	-0.06*** (-13.99)	-0.08*** (-33.78)	-0.10*** (-9.97)	-0.11*** (-20.98)	0.03*** (10.47)	0.04*** (20.42)	0.06*** (7.25)
Year observed: last	-0.43*** (-59.93)	-0.28*** (-85.30)	-0.47*** (-30.92)	-0.32*** (-43.58)	0.20*** (42.37)	0.15*** (58.04)	0.25*** (21.38)
3 years before earn. miss.			-0.03*** (-4.59)	-0.03*** (-10.06)		0.02*** (2.97)	0.02*** (9.53)
2 years before earn. miss.			-0.05*** (-7.76)	-0.04*** (-15.33)		0.04*** (5.08)	0.04*** (13.13)
1 year before earn. miss.			-0.27*** (-27.16)	-0.26*** (-78.38)		0.15*** (15.31)	0.12*** (38.05)
1 year after earn. miss.			-0.39*** (-34.56)	-0.43*** (-112.15)		0.23*** (20.82)	0.19*** (50.92)
2 years after earn. miss.			-0.15*** (-19.29)	-0.14*** (-46.64)		0.05*** (6.99)	0.07*** (22.16)
3 years after earn. miss.			-0.13*** (-16.97)	-0.11*** (-36.71)		0.03*** (4.94)	0.04*** (15.55)
Adj. R sq.	0.126	0.057	0.102	0.079	0.053	0.025	0.033
No. obs.	379080	2367552	330748	2298429	379080	2367552	2298429
No. indiv.	18130	102825	13635	90668	18130	102825	13635

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Log residual earnings are residuals from cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. German data span the period 1984–2008, while Danish data span the period 1981–2006. The dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1984 in German data and 1981 in Danish data, zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 2008 in German data and 2006 in Danish data, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE 5: EARNINGS, WAGE, AND DAYS WORKED RESIDUALS. GERMAN DATA.

	9 consec.			20 not nec. consec.		
	Earn. (1)	Days (2)	Daily Wages (3)	Earn. (4)	Days (5)	Daily Wages (6)
Year observed: first	-0.57*** (-64.24)	-0.43*** (-57.32)	-0.14*** (-32.99)	-0.67*** (-35.59)	-0.49*** (-31.91)	-0.17*** (-16.48)
Year observed: second	-0.11*** (-22.57)	-0.03*** (-8.79)	-0.09*** (-23.34)	-0.16*** (-12.50)	-0.03*** (-4.37)	-0.11*** (-11.56)
Year observed: third	-0.07*** (-16.71)	-0.02*** (-6.76)	-0.05*** (-16.48)	-0.11*** (-10.04)	-0.02*** (-3.08)	-0.08*** (-8.86)
Year observed: second-to-last	-0.03*** (-8.34)	-0.02*** (-6.66)	-0.01*** (-5.48)	-0.05*** (-6.59)	-0.02*** (-3.19)	-0.03*** (-6.52)
Year observed: next-to-last	-0.06*** (-13.99)	-0.03*** (-10.61)	-0.03*** (-9.93)	-0.10*** (-10.28)	-0.04*** (-5.66)	-0.06*** (-9.60)
Year observed: last	-0.43*** (-59.93)	-0.38*** (-59.94)	-0.05*** (-15.03)	-0.48*** (-31.25)	-0.38*** (-30.14)	-0.09*** (-11.49)
3 years before earn. miss.				-0.03*** (-4.19)	-0.03*** (-6.64)	-0.00 (-0.22)
2 years before earn. miss.				-0.05*** (-7.74)	-0.04*** (-8.09)	-0.01*** (-2.91)
1 year before earn. miss.				-0.27*** (-27.10)	-0.23*** (-27.08)	-0.04*** (-8.23)
1 year after earn. miss.				-0.39*** (-34.42)	-0.27*** (-29.88)	-0.12*** (-23.85)
2 years after earn. miss.				-0.15*** (-19.07)	-0.05*** (-9.39)	-0.10*** (-20.62)
3 years after earn. miss.				-0.12*** (-16.42)	-0.03*** (-7.26)	-0.09*** (-17.54)
Adj. R sq.	0.126	0.193	0.013	0.103	0.150	0.021
No. obs.	379080	379080	379080	330748	330748	330748
No. indiv.	18130	18130	18130	13635	13635	13635

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Log residual earnings (days) [daily wages] are residuals from cross-sectional regressions of log earnings (log days worked) [log daily wages] on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. The dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1984, zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 2008, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE 6: VARIANCES OF PERMANENT AND TRANSITORY SHOCKS IN THE PERMANENT-TRANSITORY DECOMPOSITION OF EARNINGS.

	9 consec.				20 not nec. consec.			
	German data		Danish data		German data		Danish data	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)
Perm. var., full sample, $\hat{\sigma}_{\xi}^2$	0.013	0.024	0.016	0.019	0.0096	0.018	0.013	0.019
Perm. var., outliers, $\hat{\sigma}_{\xi,o}^2$	0.034	0.158	0.053	0.124	-0.009	0.137	-0.004	0.133
Perm. var., net of outliers, $\hat{\sigma}_{\xi,n}^2$	0.010	0.010	0.013	0.013	0.0097	0.0097	0.013	0.013
Trans. var., full sample, $\hat{\sigma}_{\epsilon}^2$	0.020	0.008	0.014	0.009	0.018	0.007	0.019	0.009
Trans. var., outliers, $\hat{\sigma}_{\epsilon,o}^2$	0.143	0.011	0.104	0.022	0.162	0.011	0.173	0.030
Trans. var., net of outliers, $\hat{\sigma}_{\epsilon,n}^2$	0.007	0.007	0.008	0.008	0.007	0.007	0.008	0.008

Notes: The variances are calculated as in Eq. (2)–(4). See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations.

TABLE 7: ESTIMATES OF THE EARNINGS PROCESS IN UNBALANCED SAMPLES.
GERMAN DATA.

	Full sample		Drop first & last three obs.		Model obs. around miss.	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)
A: 9 consec. sample						
$\hat{\phi}_p$	0.976 (0.001)	0.992 (0.0008)	0.982 (0.001)	0.994 (0.001)	0.973 (0.001)	1.0 (0.002)
$\hat{\sigma}_\xi^2$	0.0078 (0.0002)	0.019 (0.0003)	0.006 (0.0001)	0.005 (0.0001)	0.008 (0.0002)	0.005 (0.0002)
$\hat{\theta}$	0.129 (0.005)	0.153 (0.009)	0.197 (0.007)	0.186 (0.007)	0.142 (0.004)	0.141 (0.006)
$\hat{\sigma}_\epsilon^2$	0.024 (0.0003)	0.009 (0.0002)	0.010 (0.0002)	0.009 (0.0002)	0.01 (0.0002)	0.01 (0.0002)
$\hat{\sigma}_\alpha^2$	0.024 (0.002)	— —	0.019 (0.002)	— —	0.022 (0.001)	— —
B: 20 not nec. consec. sample						
$\hat{\phi}_p$	0.999 (0.001)	0.991 (0.001)	0.992 (0.001)	0.995 (0.001)	0.999 (0.001)	0.994 (0.002)
$\hat{\sigma}_\xi^2$	0.0048 (0.0001)	0.009 (0.0002)	0.0047 (0.0001)	0.0046 (0.0001)	0.0046 (0.0001)	0.0057 (0.0002)
$\hat{\theta}$	0.119 (0.008)	0.192 (0.008)	0.204 (0.007)	0.190 (0.008)	0.176 (0.007)	0.167 (0.006)
$\hat{\sigma}_\epsilon^2$	0.016 (0.0003)	0.009 (0.0002)	0.009 (0.0002)	0.008 (0.0002)	0.0095 (0.0003)	0.0096 (0.0002)
$\hat{\sigma}_\alpha^2$	0.027 (0.002)	— —	0.021 (0.001)	— —	0.029 (0.001)	— —

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Columns (1) and (2) reproduce the corresponding estimates from Table 3. In columns (1)–(4), the estimated earnings process is: $y_{it} = \alpha_i + p_{it} + \tau_{it}$, where $p_{it+1} = \phi_p p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta \epsilon_{it}$. In columns (5)–(8), the estimated earnings process is in Eq. (5) in the text. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. German data span the period 1984–2008, while Danish data span the period 1981–2006. Full estimation results for the models in columns (5)–(8) are contained in Table A-11.

TABLE 8: ESTIMATES OF THE EARNINGS PROCESS IN UNBALANCED SAMPLES.
DANISH DATA.

	Full sample		Drop first & last three obs.		Model obs. around miss.	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)
A: 9 consec. sample						
$\hat{\phi}_p$	0.955 (0.001)	0.987 (0.0004)	0.957 (0.001)	0.980 (0.0006)	0.955 (0.0004)	0.990 (0.0003)
$\hat{\sigma}_\xi^2$	0.008 (0.0001)	0.013 (0.0001)	0.007 (0.0001)	0.007 (0.0001)	0.008 (0.00004)	0.007 (0.00004)
$\hat{\theta}$	0.204 (0.002)	0.209 (0.003)	0.226 (0.003)	0.220 (0.003)	0.220 (0.001)	0.204 (0.001)
$\hat{\sigma}_\epsilon^2$	0.019 (0.0001)	0.012 (0.0001)	0.012 (0.0001)	0.011 (0.0001)	0.014 (0.0001)	0.013 (0.0001)
$\hat{\sigma}_\alpha^2$	0.020 (0.0004)	— —	0.019 (0.0004)	— —	0.027 (0.0002)	— —
B: 20 not nec. consec. sample						
$\hat{\phi}_p$	0.964 (0.001)	0.982 (0.0006)	0.963 (0.001)	0.979 (0.0007)	0.961 (0.0004)	0.983 (0.0004)
$\hat{\sigma}_\xi^2$	0.007 (0.0001)	0.012 (0.0001)	0.006 (0.0001)	0.007 (0.0001)	0.007 (0.00003)	0.008 (0.0001)
$\hat{\theta}$	0.137 (0.003)	0.217 (0.003)	0.225 (0.003)	0.221 (0.003)	0.198 (0.001)	0.191 (0.002)
$\hat{\sigma}_\epsilon^2$	0.022 (0.0002)	0.013 (0.0001)	0.013 (0.0001)	0.011 (0.0001)	0.018 (0.0001)	0.013 (0.0001)
$\hat{\sigma}_\alpha^2$	0.023 (0.0004)	— —	0.020 (0.0004)	— —	0.031 (0.0002)	— —

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Columns (1) and (2) reproduce the corresponding estimates from Table 3. In columns (1)–(4), the estimated earnings process is: $y_{it} = \alpha_i + p_{it} + \tau_{it}$, where $p_{it+1} = \phi_p p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta \epsilon_{it}$. In columns (5)–(8), the estimated earnings process is in Eq. (5) in the text. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. German data span the period 1984–2008, while Danish data span the period 1981–2006. Full estimation results for the models in columns (5)–(8) are contained in Table A-11.

TABLE 9: ESTIMATES OF THE EARNINGS PROCESS IN UNBALANCED SAMPLES.
SIMULATED “GERMAN” DATA.

	9 consec.				20 not nec. consec.			
	Full sample		Drop		Full sample		Drop	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)
$\hat{\phi}_p$	0.979 (0.001)	0.988 (0.001)	0.980 (0.0009)	0.980 (0.001)	0.997 (0.0007)	0.995 (0.001)	0.999 (0.0006)	0.999 (0.0007)
$\hat{\sigma}_\xi^2$	0.008 (0.0001)	0.016 (0.0006)	0.008 (0.0001)	0.008 (0.0001)	0.005 (0.0001)	0.009 (0.0003)	0.005 (0.0001)	0.005 (0.0001)
$\hat{\theta}$	0.133 (0.003)	0.143 (0.005)	0.170 (0.004)	0.170 (0.004)	0.152 (0.007)	0.189 (0.004)	0.20 (0.006)	0.20 (0.003)
$\hat{\sigma}_\epsilon^2$	0.018 (0.0005)	0.009 (0.0001)	0.01 (0.0001)	0.01 (0.0001)	0.014 (0.0003)	0.01 (0.0001)	0.01 (0.0002)	0.01 (0.0001)
$\hat{\sigma}_\alpha^2$	0.025 (0.002)	— —	0.025 (0.002)	— —	0.024 (0.002)	— —	0.024 (0.003)	— —

Notes: The true earnings process is in Eq. (5). In columns (1)–(4), $\sigma_\alpha^2 = 0.025$, $\phi_p = 0.98$, $\sigma_\xi^2 = 0.008$, $\theta = 0.170$, $\sigma_\epsilon^2 = 0.01$, while in columns (5)–(8), $\sigma_\alpha^2 = 0.025$, $\phi_p = 0.999$, $\sigma_\xi^2 = 0.005$, $\phi_\tau = 0.20$, $\sigma_\epsilon^2 = 0.01$. In columns (3)–(4) the first 3 (last 3) observations are dropped if an individual’s earnings spell starts (ends) later (earlier) than in 1984 (2008); in columns (7) and (8), in addition, three observations before and after missing earnings records are dropped. The results are the averages across 100 simulations. Our simulated samples replicate the respective samples from German data in terms of the number of individuals, the number of individual-year observations, and the age distribution. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. The model is estimated using the optimal weighting minimum distance method. Standard errors, calculated as the standard deviations of the estimates across simulations, are in parentheses.

APPENDIX

FOR ONLINE PUBLICATION

I Evolution of various earnings percentiles in balanced and unbalanced samples. Administrative data from Germany and Denmark

In Figure A-1 we plot the time series of selected (residual) earnings percentiles for balanced and unbalanced samples.³⁷ The top panel is based on the German data, whereas the bottom panel utilizes the Danish data.³⁸

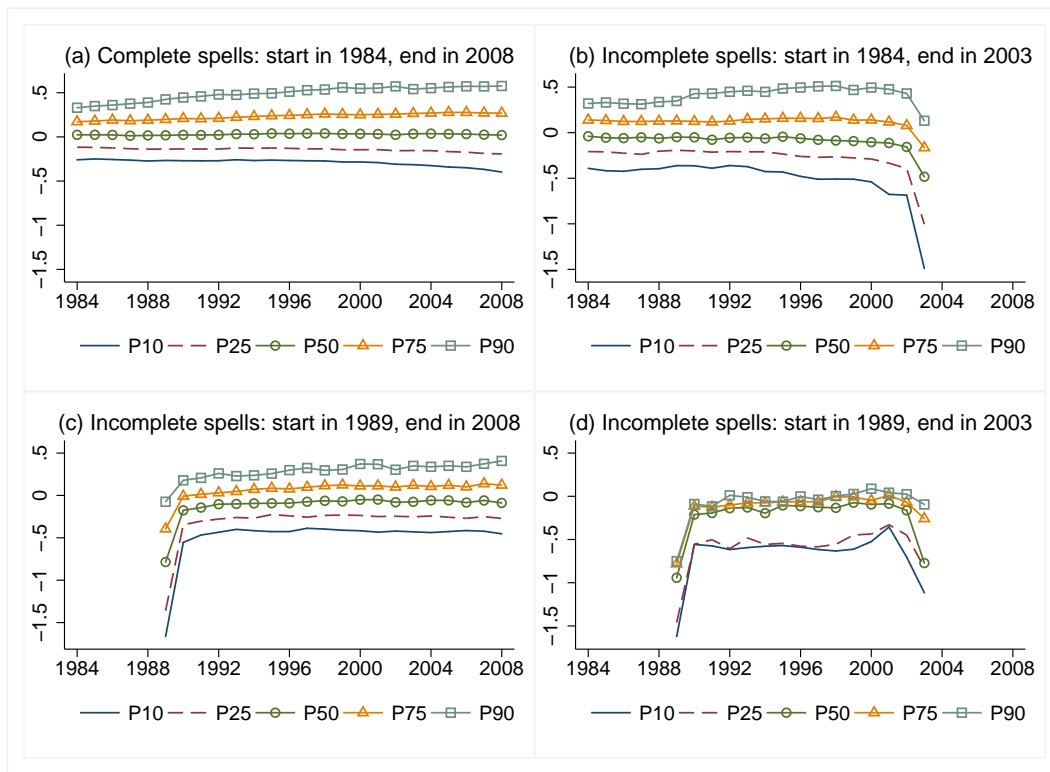
Panel (a) considers the balanced sample. For the vast majority of individuals in the balanced sample, their first year in the sample does not coincide with the first year of their earnings history. Similarly, their last year in the sample mechanically truncates earnings histories, implying that it is not the last year of the earnings spell of individuals in the sample. Thus, earnings of these individuals in the first and the last sample years are not expected to differ systematically from the neighboring observations in the interior of the sample window.

This stands in sharp contrast to the earnings histories of individuals entering and/or exiting the data in the interior of the sample window. For example, panel (b) plots earnings percentiles for individuals leaving the sample early, in 2003. Panel (c) plots earnings percentiles for individuals entering the sample late – in 1989 (1984) for the German (Danish) data. Finally, panel (d) plots earnings percentiles for individuals who enter the sample late and also leave it early. Clearly, the earnings at the start and/or the end of incomplete earnings spells are considerably lower on average and more volatile than typical earnings observations. Earnings reductions in the beginning and end of incomplete spells are visible for every percentile plotted in the figure.

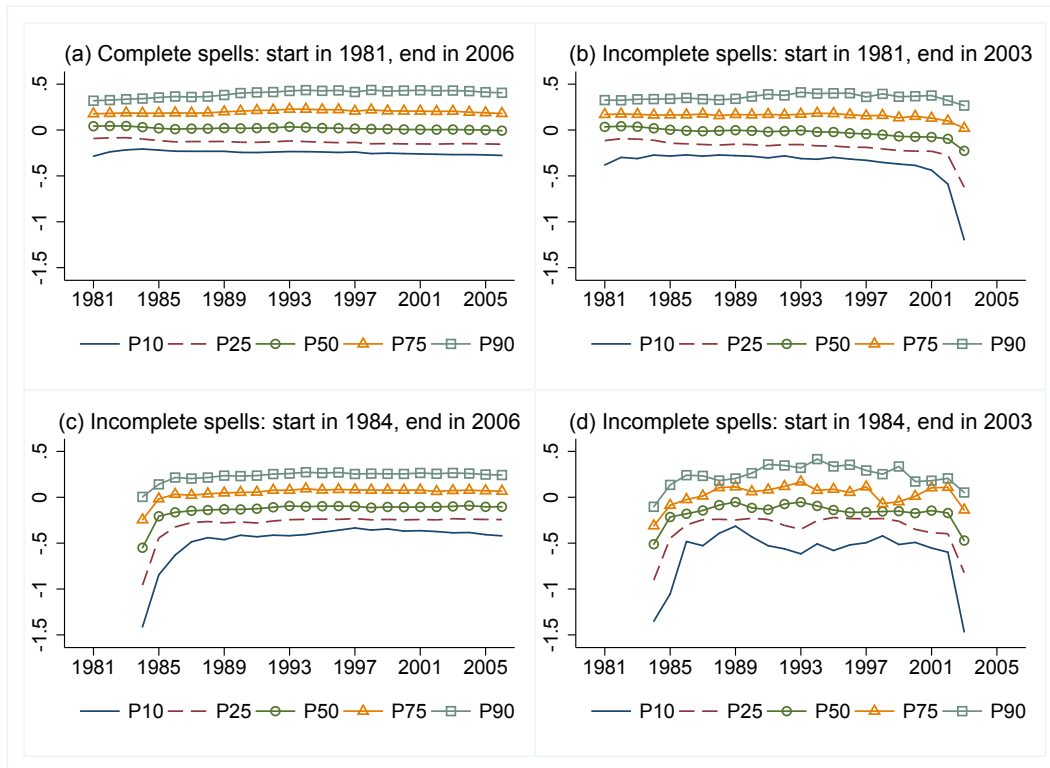
³⁷We thank an anonymous referee for proposing the design of this figure.

³⁸Statistics Denmark considers percentiles as micro data, and hence, percentiles cannot be exported from the Statistics Denmark servers. In order to comply with Statistics Denmark's data rules, the figures show local averages of these percentiles. Specifically, after percentiles are calculated, averages are constructed using the five observations whose residual is closest to the raw percentile value.

FIGURE A-1: VARIOUS EARNINGS PERCENTILES. GERMAN AND DANISH ADMINISTRATIVE DATA.



1. German data



2. Danish data

II Reasons for missing observations

As incomplete earnings spells are preceded and/or followed by missing observations, it is important to understand the nature of these observations. In survey data, missing records are often due to nonresponse, and the literature typically assumes that such observations are missing at random, i.e., the actual unobserved earnings are a random draw from an individual’s earnings distribution; e.g., Altonji et al. (2013). In administrative datasets, nonresponse is usually not an issue because data come from government records. However, there are numerous other reasons why individual earnings observations might be missing from administrative datasets. It is typically not possible, however, to know why a particular observation is missing. In this appendix, we utilize the exceptional richness of the Danish data enabled by the ability to trace individuals across various administrative registers to document the relative importance of various reasons for why individuals may not be present in the administrative earnings records.

We consider Danish men born between 1941 and 1986 who appeared on the population register at least once between 1980 and 2006. We then tabulate the fractions of those workers who would have no earnings records in typical administrative earnings datasets using the 2006 cross-section in Table A-1.³⁹ Some of these individuals have missing data because they are dead while others are alive but do not appear in the registers (we are able to determine that most of these individuals are not legally residing in Denmark on the status date of the registers). Some other individuals would not have earnings records because they truly had zero earnings, e.g., students or retirees, while others, e.g., self-employed or civil servants (government workers), might have positive earnings but are not included in typical administrative datasets (for example, our administrative German data is built from Social Security records that do not include these individuals). Finally, some missing records are introduced by researchers when conducting the analysis. As detailed in the main text, the literature often drops records with earnings outliers (defined here by an increase in earnings of more than 500 percent or a fall of more than 80 percent in adjacent years), drops records for those who work or earn little in a given year (defined here as working less than 10% of the year), or excludes earnings records because some covariates used in the analysis are missing, with education being the only covariate used in our analysis that has missing values in the Danish data.

Before discussing the patterns we observe, we provide more details on the construction of the categories in Table A-1. The category “Unemployed” contains those individuals who are unemployed and have paid into an unemployment insurance fund and are thought capable of getting a job in a relatively short time (currently set at three months). From 2002, this category also contains those individuals who did not pay into an unemployment insurance fund but claimed cash assistance and were deemed able to find a job in the short term.

The category “Unemployed, in training” contains those unemployed individuals who paid into an unemployment insurance fund and have been unsuccessful in finding a job or who were assessed as unlikely to find a job in the short term. In this case, individuals are considered “activity ready.” Activation can take three forms: job training, a job with salary subsidy, or guidance and qualification. Currently, individuals between 30 and 50 are activated at the latest

³⁹Using 2006 data allows us to extract the most information, but patterns are similar for comparable categories if we use data from other years. Detailed information for the reason for missing observations is mainly derived from the death register and from a variable (“pstill”), generated from the socioeconomic status of an individual on the register-based workforce statistics (RAS), that describes each individual’s main labor market status at the end of November in the calendar year.

when the unemployment spell hits six months; younger and older individuals are activated at the latest when their spell hits thirteen weeks. From 2002, this category also contains unemployed individuals who did not pay into unemployment insurance but instead received cash benefits and were activated. In addition, individuals on leave from unemployment (for instance, for education or childcare) are included in this category as they are generally longer-term unemployed.

The category “Sick or rehabilitation” captures those individuals who are currently receiving sickness or rehabilitation benefits. Rehabilitation allowances are meant to help an individual maintain his connection with the labor market or help that individual enter the labor market. Before 2001, individuals needed to be unemployed prior to receiving sick benefits to be considered in this category. From 2002, this is no longer a requirement.

The category “Pension” includes old-age pensioners, early retirement receivers, those who qualified for transitional allowances, and those receiving civil servant pensions. The right to early retirement depends on several conditions relating to membership in an unemployment insurance fund and restrictions on employment. When the conditions are met, a person is entitled to early retirement starting at the age of 60, until he or she transfers to the national pension. On July 1, 2004, a new old-age pension law came into force. The official retirement age was reduced from 67 to 65. Transitional allowances are given to persons aged 50-59 years who are members of an unemployment insurance fund and, in addition, fulfill several conditions (for instance, are long-term unemployed). In 2006, the scheme was phased out during the year, with the last transitional benefit recipients switching to early retirement in December 2006. Finally, the pension category contains those receiving civil servant pensions. Those 45 and younger who are receiving pensions are solely those who are receiving civil servant pensions.

Individuals whose ability to work has been substantially and permanently reduced to the point that working is not possible are captured in the category “Early retirement, unable to work.” An individual is not eligible for this type of early retirement if his reduced ability to work is temporary or if his ability to work can improve through activation, treatment, or the like.

The category “Welfare (cash assistance)” includes individuals who cannot support themselves because they have experienced a major event such as illness, unemployment, or cessation of cohabitation that has prevented the individual from providing for himself or his family. This type of support is given when no other support is available from others, and the individual is not covered by other benefits such as unemployment benefits, pension, etc.

The category “Outside workforce, other” can, among other things, consist of homemakers; young people supported by their parents; persons who live on undeclared work, and unemployed individuals who have given up looking for work and who are not registered as unemployed. This category also includes unpaid employees on parental or other leave (for instance, for education or childcare).

II.1 Discussion

The results in Table A-1 suggest that it is quite difficult to predict the counterfactual earnings for individuals whose earnings are not observed in the data. It is also difficult to associate transitions into or out of missing earnings observations with particular realizations of the earnings shocks process. For example, about nine percent of the records are missing because individuals are out of registers – this happens predominantly if individuals temporarily or

TABLE A-1: REASONS FOR MISSING OBSERVATION. DANISH DATA, 2006, AGES 20–65.

<i>Reason for missing</i>	Count	Fraction
Self-employed	110550	0.167
Pension	94495	0.143
Student	80822	0.122
Early retirement, unable to work	80093	0.121
Not in registers	60178	0.091
Dead (current year or within 10 years)	57866	0.087
Central government employees	50143	0.076
Working less than 10% of the year	30356	0.046
Outside workforce, other	25280	0.038
Unemployed	15553	0.024
Earnings outlier	13902	0.021
Welfare (cash assistance)	12815	0.019
Unemployed, in training	11551	0.017
Sick or rehabilitation	9589	0.014
Missing education	8605	0.013
Total number of individuals missing	661798	1.000

Notes: Danish men born between 1941 and 1986 and who appeared on the population register at least once between 1980 and 2006 are shown in the 2006 cross-sectional counts above. Observations are missing in this analysis if they have zero earnings, work less than 10% of the year full time, not in needed registers, dead, or are often dropped as part of sample selection or whose earnings are not available in the data (self-employed, a student, missing education, or are an earnings outlier). Further information for the reason for missing observations is mainly derived from the death register and from a variable (pstill) that describes each individual’s main labor market status at the end of November in the calendar year. Earnings outliers are defined as in the paper. The category “Not in registers” captures those individuals who are not on the population, tax, and IDA register in a given year. More than 99.9% of these individuals are not legally residing in Denmark on the status date of the registers.

permanently move out of the country, which could happen for a wide variety of reasons. The fact that these individuals have no earnings record in Denmark contains no information about their actual earnings abroad. Many individuals choose to retire early, some of whom make this transition because their productivity declined (say, due to health shocks) while others retire because they accumulated sufficient wealth due to portfolio choice or because of wealth windfalls such as lottery winnings and bequests, among a variety of other reasons. Nearly nine percent of missing observations are attributable to people who have died, and at least accidental deaths are likely to be mostly orthogonal to innovations in the earnings process. It is also unclear how the decisions to become civil servants or to transition in and out of education or self-employment relate to innovations of the earnings process. This can be potentially studied in the Danish data, but not in the administrative data from most other countries (e.g., our data from Germany), where neither the earnings of these workers can be observed, nor the reasons why they are not observed can be distinguished. Earnings that give rise to growth outliers are typically set to missing in part due to a concern that they might include measurement errors obscuring the true dynamics of earnings. As is seen in the table, earnings missing due to nonparticipation represent a relatively small share of missing observations. Selection into nonparticipation might be induced by unobserved heterogeneity

and/or adverse permanent or transitory shocks. However, there are other pathways in and out of participation, such as, e.g., paternity leaves, the timing of which may not be linked to earnings levels or dynamics.

The multitude of reasons why earnings might be missing in administrative data, in conjunction with the common inability to even determine these reasons, makes understanding their relationship with realizations of the income shocks extremely difficult. However, one thing is clear: there is little reason to expect that economic events that induce missing observations would occur on December 31. People die, become sick, start their businesses, move abroad on various dates throughout a year and they recover, exit self-employment, or return to their countries on a wide range of days. As a consequence, if earnings are missing in year t , but not in, e.g., year $t - 1$ or year $t + 1$, total earnings in years $t - 1$ or $t + 1$ would typically not represent full-year earnings because the events that induced missing earnings in year t likely started during year $t - 1$ and ended sometime during year $t + 1$. In the next section, we show that this is indeed the case and earnings preceding or following missing earnings observations are lower than average and more volatile regardless of the reason for missing earnings documented in this section.

II.2 Missing Observations by Reason and Surrounding Nonmissing Earnings

In Table A-2 we show that, regardless of the reason for why particular observations are missing, earnings residuals in the year before and the year after such observations are low and volatile. Unfortunately, even the exceptionally rich Danish data does not allow us to consider some of the reasons for missing observations, as explained below.

To create the Danish estimation sample, we used merged registers including the Danish Integrated Database for Labor Market Research, IDA, and the population register, FAIN. IDA and FAIN use different inclusion restrictions: an individual's status as of December 31st is used to determine his inclusion in FAIN the *next year*, whereas IDA uses an individual's status as of December 31st to determine his inclusion in the *current year*. As a consequence of merging these two datasets, if an individual is out of the country by the end of the year, he is missing in that year and the following year. As a result, observations adjacent to those missing because an individual is out of the country have already been dropped. Of course, this issue is specific to our merged data, and researchers using raw tax registers that are immune to those merging issues can expect to see the low mean and high variance of earnings in the years when individuals leave or enter the country.

Similarly, an individual is not included in IDA in the year in which he dies, so we do not observe the associated part-year earnings. In the dataset where such earnings are included, one can expect them to be relatively low and volatile.

Our Danish data also does not allow us to decompose total annual earnings into the earnings earned in public vs. private sectors. Thus, for those making a transition between sectors in a given year, we do not observe the effects of working only part of the year in the private sector. In datasets that exclude public sector earnings for at least some government workers only part-year private sector earnings will be observed with the associated lower mean and high variance relative to the typical full-year earnings.

TABLE A-2: RESIDUAL EARNINGS AND SQUARED RESIDUAL EARNINGS IN THE YEAR NEXT TO A MISSING OBSERVATION BY REASON.

	9 consec.		20 not nec. consec.	
	Mean	Variance	Mean	Variance
Reason for missing before the beginning of a spell that started after 1981:				
In education	-0.758 (-54.89)	0.176 (21.88)	-0.738 (-37.45)	0.163 (14.41)
Unemployed	-0.534 (-95.98)	0.151 (40.68)	-0.464 (-39.27)	0.161 (18.92)
Early retirement, unable to work	-0.433 (-10.75)	0.159 (7.027)	-0.398 (-2.796)	0.213 (2.792)
Working less than 10 percent of a full-time year	-0.291 (-37.38)	0.149 (26.05)	-0.240 (-16.57)	0.168 (15.44)
Other reason, outside workforce	-0.475 (-50.73)	0.199 (29.29)	-0.473 (-21.60)	0.194 (11.87)
Self-employed	-0.374 (-65.82)	0.182 (43.56)	-0.378 (-26.49)	0.206 (18.06)
Reason for missing after the end of a spell that ended before 2006:				
Unemployed	-0.391 (-59.61)	0.146 (33.45)	-0.437 (-18.94)	0.185 (10.64)
Early retirement, unable to work	-0.373 (-29.74)	0.176 (19.25)	-0.402 (-21.13)	0.195 (13.65)
Working less than 10 percent of a full-time year	-0.222 (-27.47)	0.127 (18.84)	-0.362 (-15.24)	0.201 (9.408)
Other reason, outside workforce	-0.364 (-56.01)	0.167 (34.17)	-0.439 (-33.01)	0.199 (18.24)
Self-employed	-0.173 (-34.98)	0.113 (28.20)	-0.300 (-19.11)	0.210 (12.51)
No. obs.	2,419,321	2,419,321	2,262,368	2,262,368
No. indiv.	113,657	113,657	90,532	90,532

Notes: Danish data spans the period 1981–2006. The 9 consecutive sample contains individuals whose maximum spell contains at least nine consecutive periods with nonmissing earnings. The 20 not necessarily consecutive sample considers individuals who have at least 20 not necessarily consecutive periods with nonmissing earnings. Individuals have missing observations in a year if they have zero earnings, work less than 10% of the year full time, not in the registers or are dead. The reasons for missing observations are derived from a variable that describes each individual's main labor market status of at the end of November in the calendar year.

III Supplemental tables for Section 3.4

This appendix includes the following tables:

- Table A-3 provides a benchmark for interpreting the results in Table 4 in the main text. Table 4 documents residual earnings and squared residual earnings at the start and end of earnings spells inside the sample window. In Table A-3, we repeat the same analysis but focus on individuals whose earnings spells begin in the first sample year or end in the last sample year. For the vast majority of these individuals, such cutoffs do not represent an actual start or end of their earnings spells; instead, the sample window mechanically truncates earnings spells in progress.
- Table A-4 helps illustrate the limited impact of the first-stage regression on our conclusions. Specifically, it shows that results are similar to those in Table 4 even when the analysis is performed on raw earnings using the regressions that do not control for education and predictable time variation.
- Table A-5 documents the low predictive power of observables – earnings growth rates before and after missing earnings records, together with education dummies and age – on the incidence of missing earnings.
- Table A-6 documents that years at the start and at the end of earnings spells are associated with a somewhat elevated probability of occupation, industry, and employer change.

TABLE A-3: RESIDUAL EARNINGS AND SQUARED RESIDUAL EARNINGS AT THE START AND END OF SAMPLES.

	9 consec.			20 not nec. consec.			9 consec.			20 not nec. consec.		
	Means						Variances					
	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)				
Year observed: first	-0.03*** (-13.35)	-0.02*** (-23.35)	-0.01*** (-2.98)	-0.01*** (-10.62)	0.04*** (17.44)	0.04*** (57.41)	0.06*** (20.36)	0.05*** (63.72)				
Year observed: second	-0.02*** (-7.68)	-0.01*** (-8.16)	0.00** (2.05)	0.00*** (4.26)	0.01*** (10.09)	0.03*** (43.70)	0.03*** (16.58)	0.03*** (49.39)				
Year observed: third	-0.01*** (-5.93)	-0.00*** (-0.38)	0.01*** (5.67)	0.01*** (5.85)	0.00*** (3.94)	0.02*** (33.90)	0.01*** (10.13)	0.02*** (39.13)				
Year observed: second-to-last	0.02*** (7.98)	0.00*** (3.26)	0.01*** (5.29)	0.01*** (8.22)	0.00* (1.71)	0.01*** (13.57)	0.02*** (9.63)	0.01*** (21.67)				
Year observed: next-to-last	0.01*** (4.49)	-0.01 (-1.23)	0.01*** (2.96)	0.01*** (6.02)	0.00*** (3.26)	0.01*** (18.17)	0.02*** (10.42)	0.02*** (26.79)				
Year observed: last	-0.00 (-0.55)	-0.01*** (-7.04)	-0.00* (-1.89)	-0.00*** (-3.06)	0.02*** (10.74)	0.02*** (24.21)	0.04*** (16.30)	0.03*** (32.67)				
3 years before earn. miss.			-0.03*** (-4.20)	-0.02*** (-9.08)			0.02*** (2.94)	0.02*** (9.84)				
2 years before earn. miss.			-0.06*** (-7.73)	-0.04*** (-16.00)			0.04*** (5.36)	0.04*** (14.46)				
1 year before earn. miss.			-0.27*** (-27.09)	-0.27*** (-80.08)			0.15*** (14.75)	0.12*** (38.21)				
1 year after earn. miss.			-0.40*** (-34.52)	-0.43*** (-111.93)			0.24*** (20.60)	0.20*** (51.48)				
2 years after earn. miss.			-0.15*** (-18.90)	-0.14*** (-46.20)			0.06*** (7.21)	0.07*** (22.69)				
3 years after earn. miss.			-0.12*** (-16.40)	-0.11*** (-36.11)			0.04*** (5.13)	0.04*** (16.26)				
Adj. R sq.	0.001	0.001	0.051	0.064	0.002	0.004	0.029	0.032				
No. obs.	379080	2367552	330748	2298429	379080	2367552	330748	2298429				
No. indiv.	18130	102825	13635	90668	18130	102825	13635	90668				

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Log residual earnings are residuals from cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. German data span the period 1984–2008, while Danish data span the period 1981–2006. The dummies “Year observed: first” – “Year observed: third” are equal to one if an individual’s first earnings record is in 1984 in German data and 1981 in Danish data, zero otherwise; “Year observed: second-to-last” – “Year observed: last” are equal to one if an individual’s last earnings record is in 2008 in German data and 2006 in Danish data, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-4: RAW EARNINGS AND EARNINGS SQUARED SURROUNDING MISSING OBSERVATIONS.

	9 consec.			20 not nec. consec.			9 consec.			20 not nec. consec.		
	Means									Variances		
	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data	German data	Danish data
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
Year obs.: first	-0.56*** (-62.93)	-0.50*** (-122.40)	-0.71*** (-37.02)	-0.53*** (-63.33)	0.23*** (40.74)	0.19*** (68.28)	0.27*** (19.29)	0.18*** (32.05)				
Year obs.: second	-0.10*** (-20.07)	-0.18*** (-59.49)	-0.18*** (-13.98)	-0.25*** (-34.90)	0.05*** (11.90)	0.08*** (32.94)	0.08*** (6.69)	0.10*** (17.37)				
Year obs.: third	-0.06*** (-14.31)	-0.11*** (-41.31)	-0.11*** (-10.40)	-0.15*** (-25.63)	0.03*** (9.01)	0.04*** (21.56)	0.04*** (4.88)	0.04*** (9.22)				
Year obs.: two years before last	-0.04*** (-12.13)	-0.05*** (-26.84)	-0.05*** (-6.84)	-0.08*** (-16.15)	0.02*** (8.20)	0.02*** (16.18)	0.02*** (3.33)	0.03*** (7.37)				
Year obs.: next-to-last	-0.08*** (-17.48)	-0.09*** (-38.18)	-0.10*** (-10.52)	-0.12*** (-22.65)	0.04*** (12.23)	0.04*** (24.10)	0.06*** (6.27)	0.06*** (11.53)				
Year obs.: last	-0.45*** (-61.98)	-0.30*** (-87.93)	-0.48*** (-31.39)	-0.33*** (-44.77)	0.20*** (43.15)	0.16*** (60.75)	0.24*** (20.07)	0.17*** (26.51)				
3 years before earn. miss.			-0.03*** (-4.68)	-0.03*** (-10.33)			0.03*** (4.13)	0.02*** (10.73)				
2 years before earn. miss.			-0.06*** (-7.83)	-0.04*** (-15.70)			0.05*** (5.82)	0.04*** (13.61)				
1 year before earn. miss.			-0.28*** (-27.72)	-0.27*** (-78.52)			0.15*** (14.87)	0.12*** (37.90)				
1 year after earn. miss.			-0.40*** (-34.85)	-0.44*** (-113.47)			0.23*** (20.95)	0.19*** (51.25)				
2 years after earn. miss.			-0.15*** (-18.90)	-0.14*** (-48.07)			0.06*** (7.52)	0.07*** (22.87)				
3 years after earn. miss.			-0.12*** (-16.17)	-0.11*** (-36.22)			0.04*** (5.79)	0.05*** (17.04)				
Adj. R sq.	0.291	0.219	0.292	0.216	0.059	0.033	0.059	0.041				
No. obs.	379080	2367552	330748	2298429	379080	2367552	330748	2298429				
No. indiv.	18130	102825	13635	90668	18130	102825	13635	90668				

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. The results are from the fixed effects regressions that, in addition to the listed variables, control for a third polynomial in age. German data span the period 1984–2008, while Danish data span the period 1981–2006. The dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1984 in German data and 1981 in Danish data, zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 2008 in German data and 2006 in Danish data, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-5: DEPENDENT VARIABLE: THE INCIDENCE OF A MISSING EARNINGS OBSERVATION; LINEAR PROBABILITY MODEL.
20 NOT NEC. CONSEC. OBSERVATIONS.

	German data					Danish data				
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Earn. growth $t - 4$ to $t - 3$					-0.43*** (-3.34)					-0.58*** (-10.65)
Earn. growth $t - 3$ to $t - 2$				-1.02*** (-6.49)	-1.16*** (-7.04)				-1.06*** (-16.34)	-1.23*** (-18.11)
Earn. growth $t - 2$ to $t - 1$	-3.21*** (-12.24)		-3.10*** (-12.14)	-3.32*** (-12.18)	-3.41*** (-12.20)	-3.15*** (-32.02)		-3.01*** (-31.24)		-3.33*** (-32.57)
Earn. growth $t + 1$ to $t + 2$		3.58*** (12.67)	3.49*** (12.60)	3.88*** (12.80)	4.03*** (12.96)		4.38*** (38.88)	4.28*** (38.54)	4.77*** (40.48)	4.98*** (41.34)
Earn. growth $t + 2$ to $t + 3$				1.32*** (7.81)	1.63*** (8.80)				1.89*** (24.08)	2.29*** (27.41)
Earn. growth $t + 3$ to $t + 4$					0.85*** (6.45)					1.27*** (19.03)
Adj. R sq.	0.005	0.007	0.012	0.013	0.014	0.007	0.012	0.017	0.020	0.021
No. obs.	210641	210641	210641	210641	210641	1486308	1486308	1486308	1486308	1486308
No. indiv.	13635	13635	13635	13635	13635	90584	90584	90584	90584	90584

Notes: See Section 3.1 for the details on the construction of the samples. "20 not nec. consec." sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Education dummy variables and age are also included in the regressions. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-6: SPELL YEARS AND JOB MOBILITY. GERMAN DATA.

	9 consec.			20 not nec. consec.		
	Chg. occ. (1)	Chg. ind. (2)	Chg. emp. (3)	Chg. occ. (4)	Chg. ind. (5)	Chg. emp. (6)
Year obs.: first	4.00*** (10.35)	4.19*** (10.94)	5.13*** (11.46)	3.88*** (4.80)	4.12*** (5.15)	5.15*** (5.52)
Year obs.: second	6.19*** (14.20)	5.76*** (13.72)	7.80*** (15.68)	5.68*** (6.36)	4.53*** (5.53)	6.15*** (6.36)
Year obs.: third	3.33*** (9.36)	3.06*** (8.98)	4.45*** (10.56)	2.00*** (2.93)	2.50*** (3.60)	3.99*** (4.62)
Year obs: sec.-to-last	0.95*** (4.50)	1.86*** (7.65)	3.02*** (9.51)	0.26 (0.85)	2.84*** (5.42)	3.27*** (5.32)
Year obs.: next-to-last	1.32*** (5.95)	1.86*** (7.76)	3.06*** (9.68)	0.82** (2.29)	2.50*** (5.03)	3.44*** (5.60)
Year obs.: last	1.79*** (7.56)	1.98*** (8.19)	2.48*** (8.29)	2.83*** (5.50)	2.91*** (5.59)	3.63*** (5.86)
3 years before miss.				0.66* (1.87)	0.57* (1.69)	2.00*** (4.29)
2 years before miss.				1.16*** (2.92)	1.37*** (3.46)	2.19*** (4.47)
1 year before miss.				2.55*** (5.54)	1.80*** (4.20)	2.66*** (4.95)
1 year after miss.				3.04*** (6.18)	3.36*** (6.86)	4.31*** (7.32)
2 years after miss.				3.43*** (7.63)	2.87*** (6.61)	3.86*** (7.41)
3 years after miss.				1.34*** (3.35)	1.34*** (3.51)	2.57*** (5.17)
Adj. R sq.	0.028	0.026	0.042	0.026	0.024	0.038
No. obs.	376793	376817	377324	328824	328920	329302
No. indiv.	18034	18032	18035	13573	13572	13573

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. “Chg. occ.” in year t (“Chg. ind.”/“Chg. emp.”) equals 100 if an individual had changed occupation (industry/employer) within year t , zero otherwise. The dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1984, zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 2008, zero otherwise. We also control for the full set of age dummies in the regressions, and for the dummies that equal one in the first three years of individual spells starting in the beginning of the sample, and the dummies that equal one in the last three years of individual spells ending in the last sample year. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

IV Alternative samples selected on the length of employment duration

The discrepancy between the estimated earnings processes in differences and levels for the unbalanced samples that we found is largely driven by the high variability of earnings surrounding missing observations that is due to an incomplete work year. In this section, we explore the sensitivity of our results to making our estimation samples more comprehensive or more restrictive with respect to annual work time.

First, we consider more comprehensive samples in which we do not drop any observations based on the employment duration, that is, we keep all nonzero earnings. They are labelled “0%” samples in Tables A-7–A-10. Second, instead of dropping observations when individuals work full-time less than 10% of the year (35 days in the German data),⁴⁰ we apply more restrictive criteria and drop observations when individuals work full-time less than k percent of the year, where $k=30\%$, or 100% (less than 110 and 365 days, respectively, in the German data).

Tables A-7 and A-8 show that the less restrictive samples feature bigger earnings cuts and variances at the start and end of incomplete earnings spells as well as the larger discrepancy between the estimated variance of permanent and transitory shocks when targeting the moments in levels and growth rates. More restrictive samples produce the opposite effect. At the extreme, when we keep only observations when individuals work the entire year, the results are similar to the balanced sample – this is because earnings around missing observations are similar to the interior observations both in level and volatility. Results for the Danish data are qualitatively similar and can be found in Tables A-9 and A-10.

⁴⁰Our baseline selection with respect to the time worked in a year is similar to that employed by Huggett et al. (2011), Storesletten et al. (2001), and Hoffmann and Malacrino (2019).

TABLE A-7: RESIDUAL EARNINGS AND SQUARED RESIDUAL EARNINGS SURROUNDING MISSING OBSERVATIONS. GERMAN DATA.

	9 consec.						20 not nec. consec.					
	Means			Variances			Means			Variances		
	0% (1)	30% (2)	100% (3)	0% (4)	30% (5)	100% (6)	0% (7)	30% (8)	100% (9)	0% (10)	30% (11)	100% (12)
Year observed: first	-0.59*** (-62.34)	-0.44*** (-58.46)	-0.09*** (-22.88)	0.24*** (40.62)	0.15*** (36.49)	0.03*** (14.23)	-0.68*** (-33.80)	-0.49*** (-31.67)	-0.09*** (-11.60)	0.30*** (20.14)	0.18*** (18.35)	0.05*** (9.02)
Year observed: second	-0.11*** (-20.90)	-0.10*** (-21.49)	-0.06*** (-17.80)	0.04*** (10.74)	0.03*** (11.08)	0.02*** (11.77)	-0.15*** (-10.83)	-0.13*** (-11.57)	-0.06*** (-9.21)	0.10*** (7.49)	0.07*** (7.69)	0.03*** (7.85)
Year observed: third	-0.07*** (-15.56)	-0.07*** (-15.91)	-0.04*** (-13.45)	0.02*** (7.58)	0.02*** (9.15)	0.01*** (8.37)	-0.09*** (-8.31)	-0.08*** (-8.79)	-0.05*** (-7.39)	0.04*** (5.68)	0.04*** (6.41)	0.02*** (6.32)
Year observed: two years before last	-0.03*** (-7.96)	-0.02*** (-7.04)	-0.01*** (-3.81)	0.01*** (4.48)	0.01*** (5.29)	0.00*** (4.60)	-0.06*** (-6.66)	-0.04*** (-6.32)	-0.02*** (-4.23)	0.03*** (3.97)	0.01*** (2.85)	0.01*** (3.37)
Year observed: next-to-last	-0.06*** (-13.75)	-0.04*** (-11.13)	-0.01*** (-4.63)	0.04*** (9.89)	0.02*** (9.14)	0.01*** (6.99)	-0.10*** (-9.67)	-0.07*** (-8.17)	-0.02*** (-3.72)	0.07*** (6.65)	0.03*** (5.68)	0.02*** (4.65)
Year observed: last	-0.46*** (-59.10)	-0.33*** (-52.73)	-0.03*** (-9.81)	0.20*** (38.44)	0.12*** (37.04)	0.02*** (10.08)	-0.50*** (-29.89)	-0.35*** (-27.75)	-0.04*** (-7.09)	0.27*** (16.13)	0.15*** (17.89)	0.02*** (5.47)
3 years before earn. miss., dummy							-0.03*** (-3.86)	-0.01** (-2.36)	0.00 (1.61)	0.01 (1.31)	0.01* (1.69)	0.00*** (2.97)
2 years before earn. miss., dummy							-0.06*** (-7.51)	-0.03*** (-5.59)	-0.00* (-1.91)	0.05*** (4.21)	0.01*** (3.15)	0.00*** (3.74)
1 year before earn. miss., dummy							-0.29*** (-24.04)	-0.20*** (-23.77)	-0.02*** (-7.01)	0.21*** (12.24)	0.09*** (13.46)	0.01*** (6.87)
1 year after earn. miss., dummy							-0.40*** (-30.22)	-0.28*** (-30.08)	-0.04*** (-17.86)	0.29*** (15.37)	0.14*** (17.64)	0.01*** (6.87)
2 years after earn. miss., dummy							-0.16*** (-17.40)	-0.13*** (-19.27)	-0.03*** (-15.98)	0.07*** (6.22)	0.04*** (7.33)	0.01*** (5.23)
3 years after earn. miss., dummy							-0.13*** (-14.76)	-0.11*** (-16.24)	-0.04*** (-14.56)	0.05*** (4.31)	0.03*** (5.96)	0.00*** (3.40)
Adj. R sq.	0.130	0.091	0.007	0.052	0.038	0.005	0.101	0.070	0.009	0.052	0.035	0.005
No. obs.	367079	363761	323514	367079	363761	323514	321759	317899	274704	321759	317899	274704
No. indiv.	17379	17241	16052	17379	17241	16052	13232	13088	11566	13232	13088	11566

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Log residual earnings are residuals from cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. German data span the period 1984–2008, while Danish data span the period 1981–2006. The dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1984 in German data and 1981 in Danish data, zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 2008 in German data and 2006 in Danish data, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-8: ESTIMATES OF THE EARNINGS PROCESS, VARIOUS SAMPLES SELECTED ON FRACTION OF THE YEAR WORKED.
GERMAN DATA.

	9 consec. sample						20 not nec. consec. sample					
	0%		30%		100%		0%		30%		100%	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)	Levs. (9)	Diffs. (10)	Levs. (11)	Diffs. (12)
$\hat{\phi}_p$	0.977 (0.001)	0.993 (0.001)	0.978 (0.001)	0.993 (0.001)	0.987 (0.001)	0.998 (0.001)	0.994 (0.001)	0.990 (0.001)	0.994 (0.001)	0.991 (0.001)	0.995 (0.001)	0.998 (0.001)
$\hat{\sigma}_\xi^2$	0.008 (0.0002)	0.019 (0.0003)	0.007 (0.0002)	0.013 (0.0002)	0.005 (0.0001)	0.004 (0.0001)	0.005 (0.0001)	0.008 (0.0002)	0.004 (0.0001)	0.006 (0.0001)	0.003 (0.0001)	0.003 (0.0001)
$\hat{\theta}$	0.131 (0.005)	0.158 (0.009)	0.131 (0.005)	0.153 (0.009)	0.154 (0.019)	0.143 (0.017)	0.133 (0.009)	0.195 (0.008)	0.134 (0.008)	0.181 (0.008)	0.188 (0.015)	0.170 (0.017)
$\hat{\sigma}_\epsilon^2$	0.023 (0.0004)	0.009 (0.0002)	0.016 (0.0002)	0.007 (0.0002)	0.003 (0.0001)	0.002 (0.0001)	0.014 (0.0003)	0.009 (0.0002)	0.011 (0.0002)	0.007 (0.0002)	0.003 (0.0001)	0.002 (0.0001)
$\hat{\sigma}_\alpha^2$	0.024 (0.002)	– –	0.024 (0.002)	– –	0.020 (0.001)	– –	0.027 (0.002)	– –	0.024 (0.001)	– –	0.019 (0.001)	– –
χ^2 (d.f.)	977.1 320	689.5 296	908.8 320	752.2 296	899 320	896.8 296	1530.6 320	1381.7 296	1548.6 320	1294.2 296	1176.6 320	932.9 296

Notes: The estimated earnings process is: $y_{it} = \alpha_i + p_{it} + \tau_{it}$, where $p_{it+1} = \phi_p p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta \epsilon_{it}$. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. See Section 3.1 for the details on the construction of the samples. German data span the period 1984–2008, while Danish data span the period 1981–2006. In the Danish data, columns headed with 0% display results for the sample in which no observations are dropped according to the fraction of the year worked. Columns headed with 30% (100%) present results when observations in which a person worked less than 30 (100) percent of the year full time are dropped. The 9 consecutive sample contains individuals whose maximum spell contains at least nine consecutive periods with nonmissing earnings. The balanced sample contains individuals whose spell covers the entire 25 periods. The 20 not necessarily consecutive sample considers individuals who have at least 20 not necessarily consecutive periods with nonmissing earnings.

TABLE A-9: RESIDUAL EARNINGS AND SQUARED RESIDUAL EARNINGS SURROUNDING MISSING OBSERVATIONS. DANISH DATA.

	9 consec.						20 not nec. consec.					
	Means			Variances			Means			Variances		
	0% (1)	30% (2)	100% (3)	0% (4)	30% (5)	100% (6)	0% (7)	30% (8)	100% (9)	0% (10)	30% (11)	100% (12)
Year observed: first	-0.52*** (-94.47)	-0.33*** (-111.19)	-0.05*** (-51.41)	0.25*** (51.66)	0.11*** (66.30)	0.01*** (39.17)	-0.51*** (-36.64)	-0.34*** (-61.37)	-0.03*** (-16.82)	0.38*** (16.54)	0.12*** (33.96)	0.02*** (18.39)
Year observed: second	-0.19*** (-44.87)	-0.11*** (-51.74)	-0.03*** (-37.47)	0.11*** (29.04)	0.04*** (34.63)	0.01*** (30.95)	-0.22*** (-20.65)	-0.12*** (-27.35)	-0.02*** (-8.95)	0.17*** (14.34)	0.05*** (23.44)	0.01*** (20.20)
Year observed: third	-0.11*** (-32.09)	-0.06*** (-33.50)	-0.02*** (-27.44)	0.06*** (19.75)	0.02*** (21.28)	0.00*** (18.15)	-0.13*** (-15.43)	-0.06*** (-17.05)	-0.00** (-2.19)	0.08*** (9.73)	0.02*** (12.76)	0.01*** (14.78)
Year observed: two years before last	-0.06*** (-21.61)	-0.03*** (-17.16)	0.01*** (6.77)	0.03*** (10.94)	0.01*** (9.86)	0.00*** (12.61)	-0.10*** (-14.99)	-0.05*** (-11.55)	-0.02*** (-6.56)	0.06*** (5.95)	0.01*** (5.36)	0.01*** (8.00)
Year observed: next-to-last	-0.10*** (-30.41)	-0.05*** (-26.08)	0.00*** (2.58)	0.05*** (17.66)	0.02*** (17.95)	0.01*** (20.88)	-0.15*** (-19.26)	-0.08*** (-16.82)	-0.03*** (-9.60)	0.12*** (9.35)	0.03*** (9.96)	0.02*** (10.72)
Year observed: last	-0.40*** (-83.40)	-0.19*** (-67.06)	0.00*** (2.98)	0.22*** (49.80)	0.09*** (48.47)	0.01*** (27.73)	-0.46*** (-38.66)	-0.21*** (-34.80)	-0.04*** (-12.70)	0.43*** (15.40)	0.10*** (21.43)	0.03*** (13.34)
3 years before earn. miss., dummy							-0.02*** (-4.39)	-0.01*** (-3.60)	0.01*** (8.79)	0.03*** (3.83)	0.01*** (7.92)	0.00*** (7.40)
2 years before earn. miss., dummy							-0.03*** (-6.60)	-0.03*** (-14.12)	0.00*** (3.09)	0.07*** (8.04)	0.02*** (11.40)	0.00*** (10.41)
1 year before earn. miss., dummy							-0.43*** (-64.13)	-0.19*** (-72.94)	0.00*** (4.77)	0.31*** (20.12)	0.07*** (35.29)	0.01*** (17.27)
1 year after earn. miss., dummy							-0.61*** (-80.92)	-0.29*** (-100.95)	-0.04*** (-43.20)	0.45*** (25.63)	0.10*** (45.57)	0.01*** (16.20)
2 years after earn. miss., dummy							-0.14*** (-24.92)	-0.09*** (-42.69)	-0.02*** (-29.28)	0.16*** (13.08)	0.03*** (18.21)	0.00*** (11.05)
3 years after earn. miss., dummy							-0.11*** (-21.50)	-0.07*** (-31.42)	-0.02*** (-25.58)	0.07*** (7.89)	0.02*** (13.58)	0.00*** (6.66)
Adj. R sq.	0.047	0.038	0.005	0.019	0.021	0.005	0.069	0.052	0.006	0.038	0.026	0.005
No. obs.	2481879	2360611	1348153	2481879	2360611	1348153	2387468	2314898	1288165	2387468	2314898	1288165
No. indiv.	104305	101502	72911	104305	101502	72911	93320	91265	53517	93320	91265	53517

Notes: See Section 3.1 for the details on the construction of the samples. "9 consec." sample contains individual earnings spells with nine or more consecutive earnings observations. "20 not nec. consec." sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. Log residual earnings are residuals from cross-sectional regressions of log earnings on educational dummies, a third polynomial in age, and the interactions of the age polynomial with the educational dummies. German data span the period 1984-2008, while Danish data span the period 1981-2006. The dummies "Year observed: first"-"Year observed: third" are equal to one if an individual's first earnings record is later than in 1984 in German data and 1981 in Danish data, zero otherwise; "Year observed: second-to-last"-"Year observed: last" are equal to one if an individual's last earnings record is earlier than in 2008 in German data and 2006 in Danish data, zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-10: ESTIMATES OF THE EARNINGS PROCESS, VARIOUS SAMPLES SELECTED ON FRACTION OF THE YEAR WORKED.
DANISH DATA.

	9 consec. sample						20 not nec. consec. sample					
	0%			30%			100%			0%		
	Levs. (1)	Diffs. (2)		Levs. (3)	Diffs. (4)		Levs. (5)	Diffs. (6)		Levs. (7)	Diffs. (8)	
$\hat{\phi}_p$	0.933 (0.001)	0.990 (0.0004)		0.960 (0.001)	0.986 (0.0004)		0.982 (0.001)	0.990 (0.001)		0.963 (0.001)	0.978 (0.007)	
										0.968 (0.001)	0.979 (0.001)	
$\hat{\sigma}_\xi^2$	0.012 (0.0002)	0.014 (0.0002)		0.007 (0.0001)	0.009 (0.0001)		0.004 (0.0001)	0.004 (0.0001)		0.008 (0.0001)	0.012 (0.0001)	
										0.006 (0.0001)	0.008 (0.0001)	
$\hat{\theta}$	0.200 (0.003)	0.228 (0.003)		0.179 (0.002)	0.182 (0.003)		0.118 (0.005)	0.110 (0.006)		0.178 (0.004)	0.230 (0.003)	
										0.142 (0.003)	0.189 (0.003)	
$\hat{\sigma}_\epsilon^2$	0.018 (0.0001)	0.013 (0.0001)		0.011 (0.0001)	0.008 (0.0001)		0.002 (0.0001)	0.002 (0.0001)		0.020 (0.0002)	0.012 (0.0001)	
										0.013 (0.0001)	0.008 (0.0001)	
$\hat{\sigma}_\alpha^2$	0.045 (0.001)	–		0.018 (0.001)	–		0.017 (0.001)	–		0.036 (0.001)	–	
										0.019 (0.0003)	–	
χ^2 (d.f.)	6162.8 346	5128.7 321		7225.7 346	4827.7 321		3652.3 346	3865.2 321		5618.5 346	6965.9 321	
										7711.5 346	6686.4 321	
										3155.0 346	4003.4 321	

Notes: The estimated earnings process is: $y_{it} = \alpha_i + p_{it} + \tau_{it}$, where $p_{it+1} = \phi_p p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta \epsilon_{it}$. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. See Section 3.1 for the details on the construction of the samples. German data span the period 1984–2008, while Danish data span the period 1981–2006. In the Danish data, columns headed with 0% display results for the sample in which no observations are dropped according to the fraction of the year worked. Columns headed with 30% (100%) present results when observations in which a person worked less than 30 (100) percent of the year full time are dropped. The 9 consecutive sample contains individuals whose maximum spell contains at least nine consecutive periods with nonmissing earnings. The balanced sample contains individuals whose maximum spell covers the entire 26 periods. The 20 not necessarily consecutive sample considers individuals who have at least 20 not necessarily consecutive periods with nonmissing earnings.

V Full estimates of an extended earnings process

TABLE A-11: ESTIMATES OF THE EARNINGS PROCESS IN UNBALANCED SAMPLES.

	9 consec.				20 not nec. consec.			
	German data		Danish data		German data		Danish data	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)
$\hat{\phi}_p$	0.973 (0.001)	1.0 (0.002)	0.954 (0.0004)	0.985 (0.0004)	0.999 (0.001)	0.994 (0.002)	0.963 (0.0004)	0.981 (0.0004)
$\hat{\sigma}_\xi^2$	0.008 (0.0002)	0.005 (0.0002)	0.008 (0.0001)	0.006 (0.00004)	0.0046 (0.00004)	0.0057 (0.0002)	0.007 (0.00004)	0.008 (0.0001)
$\hat{\theta}$	0.142 (0.004)	0.141 (0.006)	0.176 (0.002)	0.233 (0.002)	0.176 (0.007)	0.167 (0.006)	0.219 (0.002)	0.204 (0.002)
$\hat{\sigma}_\epsilon^2$	0.01 (0.0002)	0.01 (0.0002)	0.013 (0.00006)	0.014 (0.00006)	0.0095 (0.0003)	0.0096 (0.0002)	0.014 (0.0001)	0.013 (0.0001)
$\hat{\sigma}_\alpha^2$	0.022 (0.001)	— —	0.024 (0.0002)	— —	0.029 (0.001)	— —	0.022 (0.0002)	— —
$\hat{\mu}_{\nu_t^f}$	-0.55 (0.007)	-0.56 (0.007)	-0.39 (0.002)	-0.47 (0.003)	-0.55 (0.014)	-0.57 (0.014)	-0.43 (0.004)	-0.47 (0.005)
$\hat{\mu}_{\nu_t^l}$	-0.41 (0.005)	-0.42 (0.006)	-0.29 (0.002)	-0.30 (0.002)	-0.33 (0.011)	-0.22 (0.01)	-0.31 (0.004)	-0.20 (0.004)
$\hat{\sigma}_{\nu_t^f}^2$	0.25 (0.005)	0.23 (0.005)	0.17 (0.002)	0.19 (0.002)	0.28 (0.01)	0.31 (0.01)	0.18 (0.003)	0.19 (0.004)
$\hat{\sigma}_{\nu_t^l}^2$	0.21 (0.004)	0.21 (0.004)	0.18 (0.001)	0.15 (0.001)	0.14 (0.009)	0.03 (0.006)	0.00 (0.00)	0.00 (0.002)
$\hat{\mu}_{\nu_{t+1}^m}$					-0.36 (0.008)	-0.33 (0.007)	-0.49 (0.002)	-0.44 (0.002)
$\hat{\mu}_{\nu_{t-1}^m}$					-0.24 (0.007)	-0.21 (0.008)	-0.27 (0.002)	-0.26 (0.002)
$\hat{\sigma}_{\nu_{t+1}^m}^2$					0.20 (0.008)	0.13 (0.008)	0.24 (0.002)	0.11 (0.002)
$\hat{\sigma}_{\nu_{t-1}^m}^2$					0.11 (0.007)	0.07 (0.007)	0.14 (0.002)	0.09 (0.002)

Notes: See Section 3.1 for the details on the construction of the samples. “9 consec.” sample contains individual earnings spells with nine or more consecutive earnings observations. “20 not nec. consec.” sample contains individual earnings spells with at least twenty not necessarily consecutive earnings observations. The estimated earnings process is in Eq. (5) in the text. Models are estimated using the optimally weighted minimum distance method. Asymptotic standard errors are in parentheses. German data span the period 1984–2008.

VI Additional simulations

Using simulated data in Section 4.2.1, we found that the variance of the transitory shock is unbiased when estimation is based on the moments for earnings growth rates, and that the variance of the permanent shock is recovered fairly well when estimation targets the moments in levels. In this section, we present additional simulations in which we vary the variance and persistence of simulated permanent and transitory shocks. We find that this pattern remains robust across various parametrizations.

Columns (1)–(2) of Table A-12 reproduce the results of Table 9 based on the simulation of the unbalanced German sample with 9 or more consecutive observations using the earnings process in Eq. (5), the estimated parameters of which are reported in Table A-11. In the rest of the columns of Table A-12 and continuing through Tables A-13 and A-14 we vary true values for the variance and persistence of permanent and transitory shocks as described in the heading of each column. In all of these experiments, we find that the estimated variance of permanent shocks is substantially higher when we use the moments in growth rates and that the estimated variance of transitory shocks is higher when we use the moments in levels. Moreover, estimations in levels (growth rates) recover reasonably well the permanent (transitory) component of the earnings process. This is true for both the 9 or more consecutive and the 20 not necessarily consecutive samples (results for the latter samples are not reported for brevity but available upon request).

TABLE A-12: SIMULATED DATA. 9 OR MORE CONSEC. OBS. HIGHER VAR. OF SHOCKS.

True parameters, all simulations \Rightarrow								
$\sigma_\alpha^2 = 0.025, \phi_p = 0.98, \theta = 0.170$								
True shock variances \Rightarrow	$\sigma_\xi^2 = 0.007$ $\sigma_\epsilon^2 = 0.01$		$\sigma_\xi^2 = 0.007$ $\sigma_\epsilon^2 = 0.015$		$\sigma_\xi^2 = 0.01$ $\sigma_\epsilon^2 = 0.01$		$\sigma_\xi^2 = 0.01$ $\sigma_\epsilon^2 = 0.015$	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)
$\hat{\phi}_p$	0.978 (0.001)	0.988 (0.001)	0.978 (0.001)	0.988 (0.001)	0.978 (0.001)	0.987 (0.001)	0.978 (0.001)	0.987 (0.0009)
$\hat{\sigma}_\xi^2$	0.007 (0.0001)	0.015 (0.0006)	0.007 (0.0001)	0.016 (0.0006)	0.01 (0.0001)	0.019 (0.0006)	0.01 (0.0001)	0.02 (0.0006)
$\hat{\theta}$	0.128 (0.004)	0.130 (0.006)	0.135 (0.003)	0.139 (0.004)	0.122 (0.004)	0.128 (0.006)	0.133 (0.003)	0.139 (0.0035)
$\hat{\sigma}_\epsilon^2$	0.017 (0.0004)	0.009 (0.0001)	0.023 (0.0005)	0.013 (0.0001)	0.018 (0.0004)	0.009 (0.0001)	0.023 (0.0001)	0.013 (0.0001)
$\hat{\sigma}_\alpha^2$	0.024 (0.002)	— —	0.024 (0.002)	— —	0.024 (0.003)	— —	0.024 (0.002)	— —

Notes: The true earnings process is in Eq. (5). The results are the averages across 100 simulations. Our simulated samples replicate the respective samples from German data in terms of the number of individuals, the number of individual-year observations, and the age distribution. All samples contain individual earnings spells with nine or more consecutive earnings observations. The model is estimated using the optimal weighting minimum distance method. Standard errors, calculated as the standard deviations of the estimates across simulations, are in parentheses.

TABLE A-13: SIMULATED DATA. 9 OR MORE CONSEC. OBS. LOWER VAR. OF SHOCKS.

True parameters, all simulations \Rightarrow $\sigma_\alpha^2 = 0.025, \phi_p = 0.98, \theta = 0.170$						
True shock variances \Rightarrow	$\sigma_\xi^2 = \frac{1}{2} \cdot 0.007$ $\sigma_\epsilon^2 = 0.01$		$\sigma_\xi^2 = 0.007$ $\sigma_\epsilon^2 = \frac{1}{2} \cdot 0.01$		$\sigma_\xi^2 = \frac{1}{2} \cdot 0.007$ $\sigma_\epsilon^2 = \frac{1}{2} \cdot 0.01$	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)
$\hat{\phi}_p$	0.976 (0.002)	0.990 (0.001)	0.978 (0.001)	0.988 (0.001)	0.977 (0.001)	0.989 (0.001)
$\hat{\sigma}_\xi^2$	0.0036 (0.0001)	0.01 (0.001)	0.007 (0.0001)	0.013 (0.0008)	0.0036 (0.0001)	0.007 (0.0006)
$\hat{\theta}$	0.134 (0.004)	0.134 (0.005)	0.113 (0.004)	0.105 (0.009)	0.126 (0.005)	0.118 (0.007)
$\hat{\sigma}_\epsilon^2$	0.016 (0.0006)	0.009 (0.0001)	0.011 (0.0006)	0.0042 (0.0001)	0.009 (0.0005)	0.0044 (0.0001)
$\hat{\sigma}_\alpha^2$	0.024 (0.001)	— —	0.024 (0.002)	— —	0.024 (0.002)	— —

Notes: The true earnings process is in Eq. (5). The results are the averages across 100 simulations. Our simulated samples replicate the respective samples from German data in terms of the number of individuals, the number of individual-year observations, and the age distribution. All samples contain individual earnings spells with nine or more consecutive earnings observations. The model is estimated using the optimal weighting minimum distance method. Standard errors, calculated as the standard deviations of the estimates across simulations, are in parentheses.

TABLE A-14: SIMULATED DATA. 9 OR MORE CONSEC. OBS. CHANGING PERSISTENCE OF SHOCKS.

True parameters, all simulations \Rightarrow $\sigma_\alpha^2 = 0.025, \sigma_\xi^2 = 0.007, \sigma_\epsilon^2 = 0.01$								
True shock pers. \Rightarrow	$\phi_p = 0.98$ $\theta = 0.00$		$\phi_p = 0.98$ $\theta = 0.34$		$\phi_p = 1.00$ $\theta = 0.17$		$\phi_p = 0.96$ $\theta = 0.17$	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)	Levs. (7)	Diffs. (8)
$\hat{\phi}_p$	0.979 (0.001)	0.988 (0.001)	0.978 (0.001)	0.989 (0.001)	0.998 (0.001)	0.997 (0.001)	0.957 (0.004)	0.982 (0.001)
$\hat{\sigma}_\xi^2$	0.007 (0.0001)	0.014 (0.0007)	0.007 (0.0001)	0.014 (0.0006)	0.007 (0.0001)	0.015 (0.0006)	0.007 (0.0002)	0.015 (0.0007)
$\hat{\theta}$	-0.01 (0.004)	-0.07 (0.008)	0.274 (0.004)	0.349 (0.005)	0.127 (0.003)	0.130 (0.007)	0.126 (0.004)	0.132 (0.006)
$\hat{\sigma}_\epsilon^2$	0.017 (0.0005)	0.009 (0.0001)	0.017 (0.0004)	0.008 (0.0001)	0.017 (0.0004)	0.009 (0.0001)	0.017 (0.0006)	0.009 (0.0001)
$\hat{\sigma}_\alpha^2$	0.024 (0.002)	— —	0.024 (0.002)	— —	0.022 (0.002)	— —	0.025 (0.005)	— —

Notes: The true earnings process is in Eq. (5). The results are the averages across 100 simulations. Our simulated samples replicate the respective samples from German data in terms of the number of individuals, the number of individual-year observations, and the age distribution. All samples contain individual earnings spells with nine or more consecutive earnings observations. The model is estimated using the optimal weighting minimum distance method. Standard errors, calculated as the standard deviations of the estimates across simulations, are in parentheses.

VII U.S. Survey Data from the PSID

Until recently, the research on earnings and wage dynamics has been largely based on survey data from the PSID. In this appendix, we show that the findings in this paper also apply to these data. We consider the PSID sample from Blundell et al. (2008) containing married male heads of household of ages 30–65 observed during the period 1979–1993. We drop (just a few) topcoded male earnings observations and drop data for 1993.⁴¹ Hourly wages are obtained as the ratio of male annual earnings and annual hours worked. We use the same specification and controls as in Blundell et al. (2008) to extract male earnings and hourly wage residuals. We begin by documenting that the low mean and high variance of earnings and wage observations surrounding the missing ones is also a feature of PSID data.

Properties of observations surrounding missing observations in PSID data. Replicating the analysis of administrative datasets in Section 3.4, Table A-15, columns (1)–(4), contains the results of two regressions: earnings residuals and squared earnings residuals on dummies for the first and last observations of earnings spells and observations surrounding missing records. Columns (5)–(8) of the table contain the results of analogous regressions based on hourly wages. In odd-numbered columns, the dummy “Year observed: first” equals one if an individual’s first earnings (wage) record is after 1979, while the dummy “Year observed: last” equals one if an individual’s last earnings (wage) record is prior to 1992. For comparison, in even-numbered columns, these dummies equal 1 in the first and last year of the sample window, i.e., 1979 and 1992, respectively. Wage (earnings) residuals are about 0.08 (0.10) log points lower in the few first (few first and last) periods (if they differ from the first and last sample years) but are substantially lower in the few periods right after wages (earnings) are missing.

In contrast, earnings and wage residuals are, for the most part, not different from the unconditional mean of zero in the few first and last periods of the sample window – columns (2) and (6) of the table. In columns (3) and (4) (and columns (7) and (8) for wages), we net out the mean effects of irregular observations on the residuals and then regress squared (net) residuals on the same dummies as in columns (1) and (2) (and columns (5) and (6) for wages), respectively. Squared residuals are lower in the few first sample years and higher in the last sample years due to the well-known increase in male earnings and wage inequality over the life cycle and over time – columns (4) and (8). The volatility of earnings (wages), however, is much higher in the first and last sample years if individuals’ first earnings (wage) records are not in the first sample year and last earnings (wage) records are not in the last sample year, as can be seen by comparing the first six regressors in columns (3) and (4) (and columns (7) and (8) for wages).⁴² We conclude that the patterns of low mean and high variance of earnings and wage residuals around missing earnings records in the PSID data are qualitatively similar

⁴¹Our results are qualitatively similar if we keep earnings data for 1993. We drop them because the variance of male earnings growth anomalously jumps by about 42% in 1993 when the PSID switches to a computer-assisted telephone interview.

⁴²The coefficients on the dummies measure the variances in the respective periods relative to the average variance in the sample overall measured by a constant. For instance, the estimated constant in the regression of columns (3)–(4) in Table A-15 is 0.33, so that the variance of residual earnings in the first year for the earnings spells that start in the first sample year equals 0.18 ($=0.33 - 0.15$), while the volatility of earnings in the first year for earnings spells that start later than the first sample year equals 0.42. Similarly, the difference in the volatility of earnings residuals in the last year for the spells that end earlier than the last sample year and spells ending in the last sample year equals $0.52=0.63 - 0.11$.

to those in the Danish and German administrative data.

The impact of observations surrounding missing observations on the earnings and wage process estimates in the PSID. Following Blundell et al. (2008), we estimate the earnings and wage processes assuming that the persistence of the permanent component is one and the transitory component is an MA(1). Because variances of log earnings and wages have a nontrivial time trend (see, e.g., Figure A-2), we allow permanent and transitory variances to change over time in estimation and report average variances. To accommodate the pattern of missing observations, we estimate our models by the simulated minimum distance method, assuming that the shocks are drawn from normal distributions. Since our focus is on fitting second moments only, this choice of distribution is inconsequential for the results; we experimented with a Student t-distribution for the shocks fitting, in addition, the fourth moments of the data, but our results for the variances were virtually the same. When estimating the earnings and wage process in levels, we assume that individuals start accumulating shocks at age 25 and estimate the variance of permanent shocks prior to 1979 as an additional parameter, to fit the variance of log male earnings and log male wages in the first year of the sample. The variance of the fixed effects is set to zero as it is not separately identified from the variance of permanent shocks at age 25 when the permanent component is a unit root process. A diagonal weighting matrix calculated by block-bootstrap is used to weight the moments in all estimations.

The results are presented in Tables A-16 and A-17. The variance of permanent shocks estimated using the moments in levels is 0.019 for earnings and 0.017 for wages, but 0.070 for earnings and 0.043 for wages using the moments in differences. The variance of transitory shocks using the moments in levels is estimated at 0.18 for earnings and 0.108 for wages, but only 0.086 for earnings and 0.06 for wages using the moments in differences.

As a first pass for assessing the impact of irregular observations surrounding the missing ones on the estimates, we drop the first three and last three earnings (or wage) observations if an individual's first record is after 1979 and the last record is prior to 1992, as well as the three earnings (or wage) observations before and after missing earnings (or wage) records. In these samples, the average variance of permanent shocks, using the moments in levels, is estimated at about 0.019 for earnings and 0.017 for wages, practically the same as the estimate for the whole sample. The estimated variance of permanent shocks using the moments in differences is, however, substantially reduced from 0.070 to about 0.023 for earnings and from 0.043 to 0.021 for wages. Similarly, while the variance of transitory shocks using the moments in growth rates changes little after dropping these observations, the variance of transitory shocks using the moments in levels is cut substantially. These patterns in the PSID data are once again qualitatively similar to those in the Danish and German administrative data.

Next, we explicitly recognize the presence of ν -shocks in the income and wage process as in Eq. (5) and estimate the parameters of the extended model while retaining all earnings and wage observations. As in Section 4.3, in addition to the standard income moments, we use the regression coefficients in Table A-15 to estimate the mean and variance effects of ν -shocks. We allow for estimation of the mean and variance effects only for earnings and wage observations right next to a missing record since we showed in the main text that this is sufficient to align the estimated variances of permanent and transitory shocks in levels and differences. Specifically, in the full estimation, in addition to all of the autocovariance moments in the original estimation, we also target the regression coefficients in two regressions: residuals and (net) squared earnings (wage) residuals regressed on the right-hand side: two dummies around

missing interior earnings (wage) observations, one dummy for the first earnings (wage) records if the incomplete earnings (wage) spells start later than the first sample year, one dummy for the first earnings (wage) records if spells start in the first sample year, one dummy for the last earnings (wage) records if the incomplete earnings spells end earlier than the last sample year, one dummy for the last earnings (wage) records if earnings (wage) spells end in the last sample year, and a constant. The estimated variances of permanent and transitory shocks in levels and differences reported in columns (5) and (6) are similar to their values in the “dropping experiment” summarized in columns (3) and (4).

The consequences of misspecifying the earnings and wage process are clearly visible in Figures A-2 and A-4. Panel (a) of the figures plots the variance of log earnings (wage) levels in the PSID data over time (solid line) and the fit of various models to these data. The short-dashed line plots the variance implied by the estimates of the model that targets the moments of log earnings (wages) in differences and the moments for observations surrounding the missing ones, while the long-dashed line plots the variance implied by targeting the moments for earnings (wages) in growth rates only. By the last sample year, the implied variance of log earnings (log wages) greatly exceeds the variance in the data when we use the estimates targeting the moments for growth rates only. This is the direct consequence of overestimating the variances of permanent shocks when targeting the moments in differences – panel (c) of Figure A-2 shows that the estimated variances of permanent shocks to earnings are higher for each year when we use the moments in growth rates (dashed line with diamonds) rather than the moments in levels (solid lines with circles). Targeting the same moments in differences as well as the mean and variances of observations adjacent to the missing ones, however, leads to a close match to the variance of earnings and wage levels in the data (not targeted in estimation), as indicated by the short-dashed line. This estimation produces the variances of permanent shocks to earnings (solid line with triangles) that are similar in magnitude to the variances recovered from the estimation relying on the moments in levels (solid lines with circles).

Panel (b) of the figures indicates that the earnings and wage process in Eq. (5) estimated by targeting the moments in levels as well as means and variances of observations adjacent to the missing ones provides a reasonably good fit to the (untargeted) moments for earnings and wages in growth rates (solid line for the data moments and short-dashed line for the moments implied by the estimated model). The estimation of the abbreviated process in Eq. (1) targeting the moments in levels, however, substantially overpredicts the observed earnings and wage growth variances (long-dashed line). This is obviously the consequence of overestimating the variance of transitory shocks using the moments in levels.⁴³ For earnings, this can be clearly seen in panel (d) of Figure A-2 where the solid line with circles depicts the estimated variances using the moments in levels only, and the dashed line with squares shows the estimated variances using the moments in growth rates.

⁴³For completeness, Figures A-3 and A-5 panel (a) show that both the estimated process of Eq. (1) and Eq. (5) targeting the moments in levels fit those moments fairly well. Similarly, Figures A-3 and A-5 panel (b) indicate that both processes Eq. (1) and Eq. (5) estimated by targeting the variances of earnings or wage growth rates provide a reasonably good fit to these variances in the data.

TABLE A-15: MALE EARNINGS AND WAGE RESIDUALS. PSID DATA.

Dependent variable	Earnings				Wages			
	Residuals		Squared residuals		Residuals		Squared residuals	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Year obs.: first	-0.12*** (-4.64)	-0.03* (-1.68)	0.09** (2.18)	-0.15*** (-6.78)	-0.08*** (-3.48)	-0.01 (-0.90)	0.04 (1.23)	-0.11*** (-5.75)
Year obs.: second	-0.08*** (-3.39)	-0.00 (-0.06)	-0.01 (-0.19)	-0.10*** (-3.24)	-0.05** (-2.21)	0.00 (0.27)	-0.00 (-0.12)	-0.06** (-2.27)
Year obs.: third	-0.05** (-2.12)	-0.01 (-0.52)	0.02 (0.44)	-0.15*** (-6.70)	-0.03 (-1.47)	-0.00 (-0.06)	-0.03 (-1.08)	-0.07*** (-3.09)
Year obs: sec.-to-last	-0.12*** (-3.15)	0.01 (0.74)	0.28*** (3.60)	0.00 (0.07)	-0.04 (-1.53)	0.00 (0.27)	0.08** (2.04)	0.01 (0.76)
Year obs: next-to-last	-0.01 (-0.42)	-0.02 (-1.34)	0.27*** (4.42)	0.04 (1.63)	-0.03 (-1.13)	-0.01 (-0.55)	0.08** (2.33)	0.02 (1.06)
Year obs: last	-0.07 (-1.60)	-0.02 (-1.09)	0.63*** (6.55)	0.11*** (3.52)	0.05 (1.42)	-0.01 (-0.49)	0.32*** (4.24)	0.10*** (4.10)
3 years before earn. miss.	-0.06 (-0.60)		0.36** (2.04)		-0.10 (-1.17)		0.29* (1.89)	
2 years before earn. miss.	-0.08 (-0.70)		0.72*** (2.81)		-0.07 (-0.75)		0.48*** (2.99)	
1 year before earn. miss.	-0.30** (-2.00)		1.60*** (4.22)		-0.09 (-1.06)		0.46*** (2.96)	
1 year after earn. miss.	-1.10*** (-8.54)		1.26*** (4.75)		-0.49*** (-4.67)		0.86*** (4.22)	
2 years after earn. miss.	-0.51*** (-3.54)		1.42*** (3.50)		-0.28** (-2.49)		0.84** (2.49)	
3 years after earn. miss.	-0.24 (-1.88)		0.20 (0.86)		-0.17* (-1.68)		0.22 (1.38)	
Constant	0.04** (2.45)		0.33*** (20.20)		0.02 (1.19)		0.29*** (19.91)	
No. obs.	15,411		15,411		15,401		15,401	
No. indiv.	1,740		1,740		1,739		1,739	

Notes: We use the same specification and controls as in Blundell et al. (2008) to extract male earnings and hourly wage residuals. PSID male earnings and wage data span the period 1979–1992. Earnings and wages recorded in year t reflect earnings and wages received in year $t - 1$. In columns (1) and (3), the dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is later than in 1979, and are zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is earlier than in 1992, and are zero otherwise. In columns (2) and (4), the dummies “Year observed: first”–“Year observed: third” are equal to one if an individual’s first earnings record is in 1979, and are zero otherwise; “Year observed: second-to-last”–“Year observed: last” are equal to one if an individual’s last earnings record is in 1992, and are zero otherwise. Standard errors are clustered by individual; t-statistics are in parentheses. *** significant at the 1% level, ** significant at the 5% level, * significant at the 10% level.

TABLE A-16: ESTIMATES OF THE EARNINGS PROCESS. PSID DATA.

	Full sample		Drop first & last three obs.		Model obs. around miss.	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)
$\hat{\sigma}_\xi^2$	0.019 (0.002)	0.076 (0.008)	0.019 (0.002)	0.023 (0.004)	0.02 (0.002)	0.026 (0.004)
$\hat{\sigma}_\epsilon^2$	0.18 (0.01)	0.086 (0.007)	0.099 (0.007)	0.079 (0.007)	0.108 (0.008)	0.081 (0.007)
$\hat{\theta}$	0.20 (0.03)	0.01 (0.04)	0.20 (0.05)	0.05 (0.03)	0.28 (0.05)	0.003 (0.04)

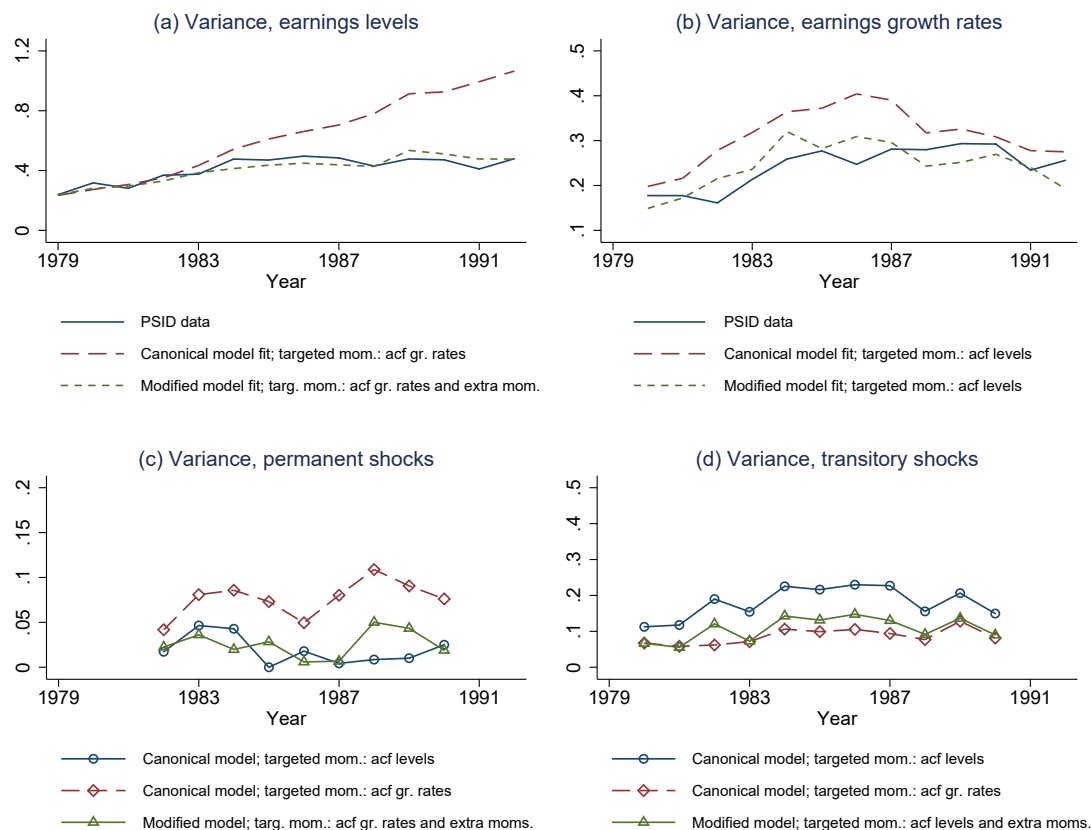
Notes: In columns (1)–(4), the estimated earnings process is: $y_{it} = p_{it} + \tau_{it}$, where $p_{it+1} = p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta\epsilon_{it}$. In columns (5)–(6), the estimated earnings process is as in Eq. (5) with ϕ_p set to 1. Models are estimated using the diagonally weighted simulated minimum distance method. Average variances of permanent and transitory shocks are reported. Asymptotic standard errors are in parentheses. PSID data span the period 1979–1992.

TABLE A-17: ESTIMATES OF THE WAGE PROCESS. PSID DATA.

	Full sample		Drop first & last three obs.		Model obs. around miss.	
	Levs. (1)	Diffs. (2)	Levs. (3)	Diffs. (4)	Levs. (5)	Diffs. (6)
$\hat{\sigma}_\xi^2$	0.017 (0.002)	0.043 (0.005)	0.017 (0.002)	0.021 (0.004)	0.018 (0.002)	0.023 (0.002)
$\hat{\sigma}_\epsilon^2$	0.108 (0.08)	0.06 (0.004)	0.075 (0.005)	0.055 (0.004)	0.072 (0.005)	0.058 (0.004)
$\hat{\theta}$	0.22 (0.04)	0.08 (0.03)	0.25 (0.05)	0.08 (0.03)	0.29 (0.05)	0.07 (0.03)

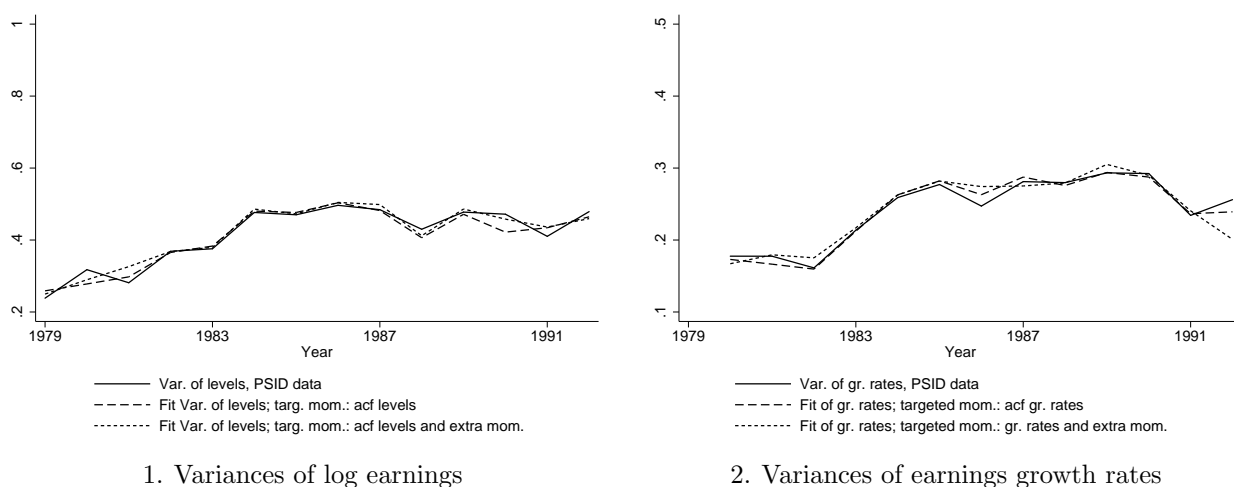
Notes: In columns (1)–(4), the estimated wage process is: $y_{it} = p_{it} + \tau_{it}$, where $p_{it+1} = p_{it} + \xi_{it+1}$ and $\tau_{it+1} = \epsilon_{it+1} + \theta\epsilon_{it}$. In columns (5)–(6), the estimated wage process is as in Eq. (5) with ϕ_p set to 1. Models are estimated using the diagonally weighted simulated minimum distance method. Average variances of permanent and transitory shocks are reported. Asymptotic standard errors are in parentheses. PSID data span the period 1979–1992.

FIGURE A-2: FIT TO THE MOMENTS OF MALE LOG EARNINGS IN LEVELS (TARGETED MOMENTS IN DIFFERENCES) AND DIFFERENCES (TARGETED MOMENTS IN LEVELS). PSID DATA.



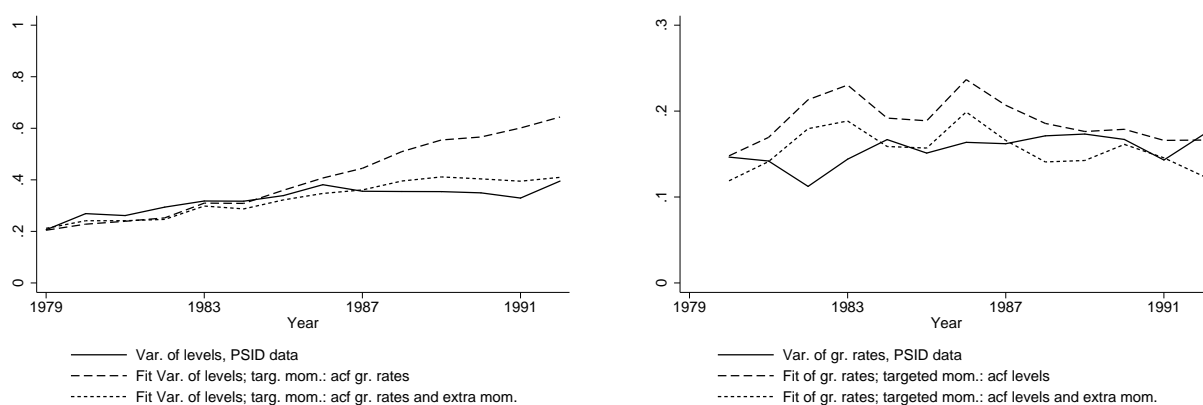
Notes: In panel (a), solid line depicts the variance of log male earnings in PSID data for the period 1979–1992; long dash line depicts the variance of log male earnings implied by the estimates of the model targeting the autocovariance moments for male earnings growth rates; short dash line depicts the variance of log male earnings implied by the model that targets the autocovariance moments for male earnings growth rates and, in addition, the mean and variance of earnings observations surrounding missing records. In panel (b), solid line depicts the variance of male earnings growth in PSID data for the period 1979–1992; long dash line depicts the variance of male earnings growth implied by the estimates of the model targeting the autocovariance moments for male earnings (moments in levels); short dash line depicts the variance of male earnings growth implied by the model that targets the autocovariance moments for male earnings and, in addition, the mean and variance of earnings observations surrounding missing records. In panel (c), solid line with circles depicts the estimated variance of permanent shocks from the model targeting the moments in levels; dash line with diamonds – targeting the moments in growth rates; and dash line with triangles – targeting the moments in growth rates and, in addition, the mean and variance of earnings observations surrounding missing records. In panel (d), solid line with circles depicts the estimated variance of transitory shocks from the model targeting the moments in levels; dash line with diamonds – targeting the moments in growth rates; and dash line with triangles – targeting the moments in growth rates and, in addition, the mean and variance of earnings observations surrounding missing records.

FIGURE A-3: FIT TO THE MOMENTS OF MALE LOG EARNINGS IN LEVELS (TARGETED MOMENTS IN LEVELS) AND DIFFERENCES (TARGETED MOMENTS IN DIFFERENCES). PSID DATA.



Notes: In panel (a), solid line depicts the variance of log male earnings in PSID data for the period 1979–1992; long dash line depicts the variance of log male earnings implied by the estimates of the model targeting the autocovariance moments for male earnings (moments in levels); short dash line depicts the variance of log male earnings implied by the model that targets the autocovariance moments for male earnings and, in addition, the mean and variance of earnings observations surrounding missing records. In panel (b), solid line depicts the variance of male earnings growth in PSID data for the period 1979–1992; long dash line depicts the variance of male earnings growth implied by the estimates of the model targeting the autocovariance moments for male earnings growth rates; short dash line depicts the variance of male earnings growth implied by the model that targets the autocovariance moments for male earnings growth rates and, in addition, the mean and variance of earnings observations surrounding missing records.

FIGURE A-4: FIT TO THE MOMENTS OF MALE LOG WAGES IN LEVELS (TARGETED MOMENTS IN DIFFERENCES) AND DIFFERENCES (TARGETED MOMENTS IN LEVELS). PSID DATA.

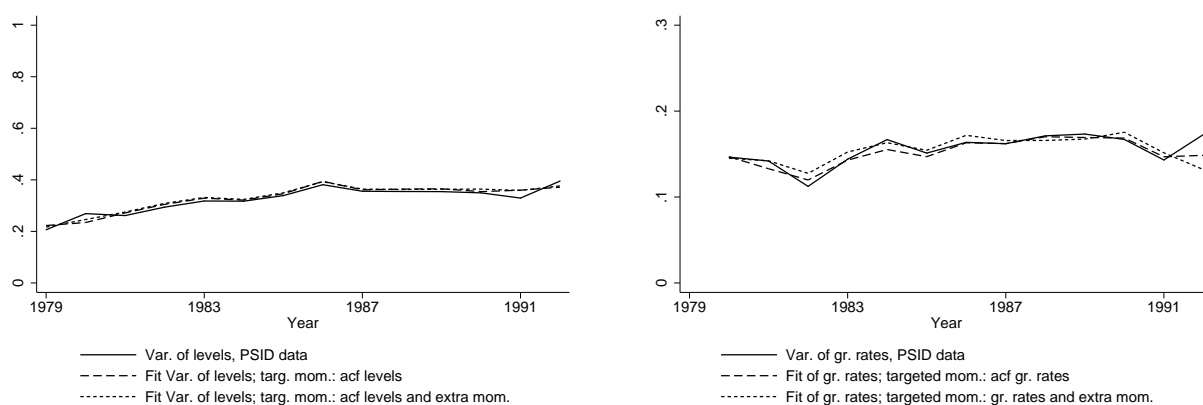


1. Variances of log wages: levels and growth rates

2. Variances of wage growth rates

Notes: In panel (a) solid line depicts the variance of log male wages in PSID data for the period 1979–1992; long dash line depicts the variance of log male wages implied by the estimates of the model targeting the autocovariance moments for male wage growth rates; short dash line depicts the variance of log male wages implied by the model that targets the autocovariance moments for male wage growth rates and, in addition, the mean and variance of wages surrounding missing records. In panel (b) solid line depicts the variance of male wage growth in PSID data for the period 1979–1992; long dash line depicts the variance of male wage growth implied by the estimates of the model targeting the autocovariance moments for male wages (moments in levels); short dash line depicts the variance of male wage growth implied by the model that targets the autocovariance moments for male wages and, in addition, the mean and variance of wages surrounding missing records.

FIGURE A-5: FIT TO THE MOMENTS OF MALE LOG WAGES IN LEVELS (TARGETED MOMENTS IN LEVELS) AND DIFFERENCES (TARGETED MOMENTS IN DIFFERENCES). PSID DATA.



1. Variances of log wages

2. Variances of wage growth rates

Notes: In panel (a) solid line depicts the variance of log male wages in PSID data for the period 1979–1992; long dash line depicts the variance of log male wages implied by the estimates of the model targeting the autocovariance moments for male wages (moments in levels); short dash line depicts the variance of log male wages implied by the model that targets the autocovariance moments for male wages and, in addition, the mean and variance of wages surrounding missing records. In panel (b) solid line depicts the variance of male wage growth in PSID data for the period 1979–1992; long dash line depicts the variance of male wage growth implied by the estimates of the model targeting the autocovariance moments for male wage growth rates; short dash line depicts the variance of male wage growth implied by the model that targets the autocovariance moments for male wage growth rates and, in addition, the mean and variance of wages surrounding missing records.

VIII Quantitative model

In this section, we calibrate the standard lifecycle model of consumption using the canonical earnings process with permanent and transitory shocks modified to include missing observations and ν -shocks with means and variances estimated from the German data. We then generate a panel dataset on income and consumption from this model and treat it as the data available to a researcher. Using this dataset, we then evaluate the consequences of adopting various empirical strategies discussed in the paper. More specifically, in one experiment, we follow the current practice in the literature and estimate the canonical earnings process, which ignores the special nature of observations surrounding the missing ones. We then use that process as an input into the standard lifecycle model of consumption and find that this approach leads to highly misleading conclusions obtained using a calibrated model. In another experiment, we drop the observations surrounding the missing ones in the original simulated dataset, estimate the canonical earnings process, and use this process as an input into a calibrated consumption model. We observe that this experiment allows the researcher to obtain the correct inference about various endogenous outcomes of the model.

VIII.1 Baseline Model

We assume that households value consumption and supply labor inelastically, preferences are CRRA, households start working life at age t_0 , retire at age t_R , and spend time in retirement until age T , when they die with certainty. Households face permanent and transitory income risk and the risk of dying with the probability of survival at age t equalling to s_t .

Household i 's problem is:

$$\max_{\{C_{it}\}_{t=t_0}^T} E_{i,t_0} \sum_{t=t_0}^T \beta^{t-t_0} s_t \frac{C_{it}^{1-\gamma} - 1}{1-\gamma},$$

subject to

$$A_{it+1} = (1+r)(A_{it} + Y_{it} - C_{it}),$$

$$\log Y_{it}^* = g_t + \alpha_i + p_{it} + \tau_{1,it} + \tau_{2,it}$$

$$p_{it} = \phi_p p_{it-1} + \xi_{it}, \quad \xi_{it} \sim \text{i.i.d.}(0, \sigma_\xi^2)$$

$$\tau_{1,it} = \epsilon_{it}, \quad \epsilon_{it} \sim \text{i.i.d.}(0, \sigma_\epsilon^2)$$

$$\tau_{2,it} = \begin{cases} \nu_{it}^b, & \nu_{it}^b \sim \text{i.i.d.}(\mu_{\nu^b}, \sigma_{\nu^b}^2), \quad \mu_{\nu^b} < 0, & \text{w. prob. } p_\nu \\ 0, & & \text{w. prob. } 1 - p_\nu \end{cases}$$

$$E_{it} = \mathbf{1}(\tau_{2,it-1} = 0),$$

$$Y_{it} = Y_{it}^* \exp(\nu_{it}^a)(1 - E_{it-1})E_{it} + Y_{it}^* E_{it-1}E_{it}, \quad \nu_{it}^a \sim \text{i.i.d.}(\mu_{\nu^a}, \sigma_{\nu^a}^2) \quad \mu_{\nu^a} < 0, \quad t = 1, \dots, t_R$$

$$Y_{it} = \kappa \exp(p_{it_R}), \quad t = t_R + 1, \dots, T$$

$$A_{it} \geq 0, \quad t = t_0, \dots, T.$$

Y_{it} is household i 's observed income at age t , stochastic until retirement age t_R , and deterministic afterwards.⁴⁴ Y_{it}^* is household i 's latent earnings at time t , which is either zero

⁴⁴In German data, we do not observe families, so we use male earnings as the sole source of unpredictable

or missing at times when $\tau_{2,it}$ is not equal to zero. g_t is the common lifecycle component of earnings; p_{it} is (log of) the permanent component of earnings, which follows an AR(1) process with persistence ϕ_p ; ξ_{it} is an i.i.d. permanent shock; ϵ_{it} is an i.i.d. transitory shock; ν_{it}^b (ν_{it}^a) is a shock preceding (following) a missing income record. Such an income process is consistent with the data patterns for incomplete earnings spells, as shown in the main text. The earnings process is the same as in Eq. (5) but is cast in terms of a selection model, where, with probability p_ν , there happens a large negative transitory shock ν_{it}^b followed by a missing value, which is then followed by temporarily low earnings affected by a large negative transitory shock ν_{it}^a .

After retirement, household i 's income is proportional to the permanent component at age t_R with a replacement rate κ , as in, e.g., Demyanyk et al. (2017). A_{it} is household i 's wealth at age t , C_{it} is household i 's consumption at age t , and E_{i,t_0} stands for household i 's expectation about future resources based on the information available at age t_0 . Households cannot borrow but can save into a riskfree asset that yields a net interest rate r .

VIII.2 Calibration

We calibrate the time discount factor β to match the wealth-to-income ratio of three. We assume that households start their life at age 26 with zero assets, retire at age 65, and die at age 90, that is, $t_0 = 26$, $t_R = 65$, and $T = 90$. The age-dependent deterministic income profile, g_t , is estimated from the German data. The replacement rate κ is set to 0.70, the interest rate r is set to 4%, and the relative risk aversion parameter γ is set to one.

We take as given the estimated income process for the German data in columns (5)–(6) of Table A-11. Specifically, we set ϕ_p to one (the permanent component is a random walk) and assume that the shocks ξ_{it} and ϵ_{it} are normally distributed with variances 0.005 and 0.01, respectively. We assume that the shocks ν_{it}^b and ν_{it}^a are normally distributed with the respective means of -0.24 and -0.36 and variances of 0.11 and 0.20. The probability of missing values, p_ν , is set to 1%, which is similar to the incidence of missing values in the German nonconsecutive unbalanced sample.

We begin by making the standard assumption that missing observations are missing at random. This means that missing earnings observations are simply not observed by the researcher, but individuals continue to receive earnings according to their standard earnings process. For example, this is a reasonable assumption for individuals with a stint of self-employment, who become civil servants, move to work overseas, or in other similar situations where earnings are there in reality but are not observed in the administrative data because of the restrictions on the populations covered by the data. Under the missing-at-random assumption, the shocks ν_{it}^b and ν_{it}^a represent the fraction of the year when individual earnings are not observed in the data, e.g., when a spell of self-employment starts in the middle of one year, continues through the next year, and ends in the middle of the year after next. We will explain below, however, that none of our conclusions are affected by assuming that for some, or even all, individuals the shocks ν_{it}^b and ν_{it}^a represent genuine cuts in annual earnings.

fluctuations in household budgets in the model.

VIII.3 Quantitative Experiments

We now simulate a dataset from the calibrated model and treat it as genuine data available to a researcher. Having generated the data from the model, we know the true values of various objects an empirical researcher might be interested in. We then assess whether the researcher will make the correct inference about those objects depending on the choice of how to estimate the earnings process. Various statistics computed from the true model-generated data are reported in column (1) of Table A-18. In Panels A, B and C we report various percentiles of income, wealth, and consumption distributions. The rest of the panels will be discussed shortly.

Having obtained this dataset, a researcher will first compute various statistics, including the earnings process. Suppose the researcher follows the standard approach in the literature and estimates the earnings process targeting the moments in differences while ignoring the special properties of observations surrounding the missing ones. As explained in the paper, this approach yields an unbiased estimate of the variance of the transitory shock at 0.01 but a significantly upward biased estimate of the permanent shock at 0.01 as opposed to the true value of 0.005 – see column (1), Panel G1. If, instead, the researcher followed any of the approaches we proposed in this paper, e.g., simply dropping earnings observations surrounding the missing ones before estimating the earnings process, she would have recovered the correct variances of permanent and transitory shocks – see column (1), Panel G2.

Having obtained these estimates of the earnings process, a researcher may use them as an input into a model. Let's assume that the researcher uses the standard model, the same as the baseline model that generated the data but with p_ν set to zero at all times. Thus, the model used by the researcher includes neither missing observations nor additional income shocks that surround missing observations in the true model that generated the data. The researcher then calibrates this model targeting the same statistics as we targeted when calibrating the baseline model (except, of course, those that relate to ν -shocks and missing observations). Having calibrated the model, the researcher computes the same statistics as computed in the baseline model. These are reported in columns (2) and (3) of Table A-18. Column (2) reports the statistics of the model that used as an input the earnings process, which was estimated taking into account the properties of earnings adjacent to the missing ones as we propose in this paper. In contrast, column (3) reports the statistics of the model that used as an input the standard earnings process targeting the moments in differences but ignored the low mean and high variance of the observations surrounding the missing ones in the data.

Comparing income, consumption, and wealth distributions in the true model in column (1) and the model based on the estimated income process that accounted for ν shocks in column (2) indicates that they are essentially identical. The only minor discrepancy is observed when comparing income percentiles because incomes in the true model include ν shocks, while incomes in the model in column (2) exclude them. In Panel D, we observe that individuals in the true model and in the model in column (2) achieve exactly the same welfare and are indifferent between living in the two model worlds.

The corresponding comparisons between the true model in column (1) and the model based on the estimated earnings process that ignored the presence of ν -shocks in the data in column (3) yield very different conclusions. Income, wealth, and consumption distributions differ significantly (despite the fact that the calibration procedure implies that the wealth-to-income ratio is the same across all models). Consumers are willing to give up 3.8% of lifetime consumption to live in the true economy rather than in the economy constructed by

the researcher to interpret the data. The root of the discrepancy is clear. The variances of permanent or transitory shocks are estimated with a bias when the low mean and high variance of observations adjacent to the missing ones are ignored (Panel G1, column (1)), and this incorrectly estimated earnings process is used as an input into the model in column (3). In contrast, when the variances of permanent and transitory shocks are computed correctly by, e.g., simply dropping observations next to the missing ones (Panel G2, column (1)), and this earnings process is used by the researcher in a model as in column (2), the model leads to the correct inference about the various endogenous outcomes considered.

Next, we highlight the importance of correctly measuring the variances of permanent and transitory shocks for assessing the amount of insurance that is available to households, a subject that has attracted considerable attention recently. The true transmission coefficients of income shocks to consumption are measured by estimating the following regression of residual consumption growth on the true permanent and transitory shocks:⁴⁵

$$\Delta \log C_{it} = \phi \xi_{it} + \psi \epsilon_{it} + \text{error}.$$

These estimated transmission coefficients for permanent and transitory shocks, ϕ and ψ respectively, are reported in Panel E. Note that they are essentially identical across the three models.

Suppose the researcher does not know the true amount of insurance but estimates those coefficients and the shock variances from the data generated by our benchmark model in column (1) following the methodology of Blundell et al. (2008) that relies on the moments in growth rates and the minimum-distance method but does not account for the special properties of observations adjacent to the missing ones. This methodology was shown to recover the true ϕ and ψ in the incomplete-markets model environment when the borrowing constraints are not binding frequently; see, e.g., Kaplan and Violante (2010). Households in the benchmark model start working at age 26, and we select the data for ages 30–65 when the borrowing constraints become less important. The transmission coefficients estimated following this methodology are reported in Panel F1, column (1). Specifically, a researcher would find that only 43% of permanent shocks are transmitted to consumption, while the remaining 57% are insured. If the researcher uses the earnings process estimated using this methodology and simulates an incomplete markets model (as is done in column (3)), the researcher would find that the transmission coefficient in the data simulated from this model is 0.83, implying that only 17% of permanent shocks to income are insured. The researcher would conclude that there is much more insurance available to households in the data than in the textbook incomplete markets model, dubbed as an excess insurance puzzle in the literature.

Hryshko and Manovskii (2018) have studied this discrepancy and traced its origin to the biased estimates of the earnings process induced by the failure to control for the properties of observations surrounding missing observations. This is apparent in Panel F2, column (1), when the transmission coefficients and the earnings process parameters are estimated after dropping the observations surrounding the missing ones. This procedure recovers the true amount of insurance available to households in the data, and if this earnings process is used in the incomplete markets model (column(2)), the model also implies the same amount of insurance as measured in the data, eliminating the excess insurance puzzle.

All the experiments we discussed so far in this section were based on the assumption that

⁴⁵In the model, consumption and income data are residualized by subtracting age-specific means from the raw data.

missing observations are missing at random. However, none of our conclusions depend on this assumption. Even making an extreme alternative assumption that earnings are zero for all missing observations and ν -shocks represent genuine changes in earnings (rather than the fraction of earnings that are not observed by the researcher) changes none of our conclusions. Simulating the models in columns (2) or (3) that do not include missing observations or ν -shocks will lead to some underestimation of risk and the corresponding overestimation of welfare. However, the effect is quantitatively small. As missing observations are rare, the effect of ignoring them leads to an overestimation of the welfare equivalent to around one percent of lifetime consumption but leaves the relative differences between the models in columns (2) and (3) unaffected.

These results imply that it is very important to use the correctly specified earnings process for inference based on quantitative incomplete markets models. However, there seems to be relatively little loss from ignoring explicit modeling of missing observations in incomplete markets models if the earnings process that is used as an input into these models is estimated in the data in a way that accounts for the effect of the low mean and high variance of observations surrounding missing observations.

TABLE A-18: QUANTITATIVE EXPERIMENTS.

	RW+iid + miss. + ν -shocks benchmark (1)	RW+iid diffs., drop around miss. (2)	RW+iid diffs., ignore ν -shocks (3)
A: Various income percentiles, in '000s of Euros			
$P10$	24.8	25.1	20.4
$P50$	38.0	38.2	36.4
$P90$	60.6	60.6	67.5
B: Various wealth percentiles, in '000s of Euros			
$P10$	10.2	10.2	10.7
$P50$	111.0	111.0	103.9
$P90$	249.1	249.1	260.3
C: Consumption inequality			
$\log(P90) - \log(P10)$, age 30	0.36	0.36	0.50
$\log(P90) - \log(P10)$, age 65	0.99	0.99	1.38
D: Permanent change in consumption to reach the benchmark's welfare			
Equiv. cons. variation, %	—	0.0	3.8
E: Insurance coefficients, true			
Transm. perm. shock	0.83	0.83	0.83
Transm. trans. shock	0.06	0.06	0.06
F1: Insurance coefficients, minimum distance estimates, ignore ν -shocks			
Transm. perm. shock	0.43	0.84	0.83
Transm. trans. shock	0.05	0.05	0.06
F2: Insurance coefficients, minimum distance estimates, drop around miss.			
Transm. perm. shock	0.83	—	—
Transm. trans. shock	0.05	—	—
G1: Income process parameters, estimates using growth rates moments, ignore ν -shocks			
Var. perm. shock	0.01	0.005	0.01
Var. trans. shock	0.01	0.01	0.01
G2: Income process parameters, estimates using growth rates moments, drop around miss.			
Var. perm. shock	0.005	—	—
Var. trans. shock	0.01	—	—
H: Internally calibrated parameter			
Time disc. factor, β	0.9729	0.9729	0.9713

Notes: See the text in Appendix VIII for an explanation of all entries in this table.