

# Robust Inference for the Frisch Labor Supply Elasticity

By MICHAEL KEANE AND TIMOTHY NEAL\*

December 12, 2022

*The Frisch labor supply elasticity plays a key role in many policy debates, but its magnitude remains controversial. Many studies use 2SLS regressions of hours changes on wage changes to estimate the Frisch elasticity. But a little appreciated power asymmetry property of 2SLS causes estimates to appear (spuriously) imprecise when they are shifted away from the OLS bias. This makes it difficult for a 2SLS t-test to detect a true positive Frisch elasticity. We illustrate this problem in an application to NLSY97 data. We obtain a Frisch estimate of 0.60 for young men, but the t-test indicates it is insignificant. In contrast, the Anderson-Rubin (AR) test – which avoids the power asymmetry problem – implies the estimate is highly significant. The same power asymmetry issue that afflicts the t-test here will arise in many IV applications. Thus, we argue the AR test should be widely adopted in lieu of the t-test.*

**Keywords:** Frisch elasticity, labor supply, weak instruments, 2SLS, Anderson-Rubin test, LIML, Continuously Updated GMM

**JEL:** J22, D15, C12, C26

## I. Introduction

The Frisch elasticity measures the response of labor supply to predictable wage changes. It plays a key role in many economic policy debates because predictable wage changes have pure substitution effects. For example, Conesa, Kitao and Krueger (2009) show that a higher Frisch elasticity implies a higher optimal tax rate on capital income. And macro models where real shocks play a key role in business cycles require the Frisch to be large to match observed fluctuations in work hours over the cycle, see Prescott (2006). Because of its importance, a large literature attempts to estimate the Frisch elasticity, as in classic papers by MaCurdy (1981) and Altonji (1986) and surveyed in Keane (2011, 2021).

Classic micro data studies in the style of MaCurdy (1981) typically find the Frisch elasticity is small and insignificant,<sup>1</sup> while macro economists using DSGE models typically calibrate it to be in the 0.50 to 2.0 range. This has led to a long-standing “macro-micro controversy” over the magnitude of the Frisch elasticity.

We argue the classic studies were inherently biased against finding the Frisch is both large and significant, due to a combination of (i) weak instrument problems and (ii) use of biased testing procedures. The weak IV problem is well understood,

\* CEPAR & School of Economics, UNSW. Corresponding author: m.keane@unsw.edu.au

<sup>1</sup>For example, MaCurdy (1981), Altonji (1986), Browning, Deaton and Irish (1985) and Ziliak and Kniesner (1999) obtain estimates for men ranging from 0.09 to 0.17.

but the testing problem arises from a little appreciated power asymmetry property of two-stage least squares (2SLS)  $t$ -tests that we explain and illustrate. Due to this power asymmetry, it is highly unlikely that any of the classic studies could have found a significant positive Frisch elasticity. Furthermore, simply using stronger instruments does not solve this problem, as it persists even when instruments are strong according to conventional weak IV testing thresholds.

The idea behind the classic studies is as follows: Given panel data on workers, one may run an OLS regression of changes in log hours on changes in log wages. This gives a downward biased estimate of the elasticity of hours with respect to *predictable* wage changes – as some changes are surprises. Instead, the approach pioneered by MaCurdy (1981) involves running a 2SLS regression, using an instrument for the change in log wages with two properties: First, it predicts wage growth at the individual level. Second, it is known at the start of the time period over which changes in wages are calculated (so it is uncorrelated with wage surprises). Then, fitted values from the first-stage of 2SLS give us predictable changes in wages, and the second stage delivers an estimate of the elasticity of hours with respect to these predictable wage changes, which is the Frisch concept.

A little appreciated property of the 2SLS estimator is that it generates a strong association between the 2SLS estimates and their standard errors – see Keane and Neal (2022). This association is positive if the OLS bias is negative. Hence, positive Frisch elasticity estimates have artificially inflated standard errors. As a result, a 2SLS  $t$ -test has little power to detect a true positive Frisch elasticity – unless instruments are far stronger than is typical in this literature.

We further show that the Anderson-Rubin (AR) test does not suffer from the power asymmetry that afflicts the  $t$ -test. Hence it is a far more reliable guide to inference in 2SLS applications. Using NLSY97 data, we obtain a Frisch elasticity of 0.60, which is quite large compared to typical estimates in the classic studies. Nevertheless, the conventional 2SLS  $t$ -test indicates it is not significantly different from zero. In contrast, the AR test indicates it is highly significant ( $p=.018$ ). Thus, application of a superior inferential procedure reveals clear evidence to support a fairly large Frisch elasticity value for young men.

It is well-known that the literature on estimating the Frisch elasticity using 2SLS has been hampered by weak instrument problems, as it is hard to find instruments that are strong predictors of wage growth. This issue was explored by Lee (2001), who shows how weak instrument problems in classic papers like MaCurdy (1981) and Altonji (1986) biased their estimates of the Frisch elasticity towards zero (i.e., towards OLS). Those authors used over-identified models, where the instruments used to predict wage growth were primarily age and schooling (including linear, quadratic and interaction terms). As Lee (2001) shows, given PSID samples of the size they had available, and the type of instruments they used, first-stage  $F$ -statistics in their models would have been no higher than one. This is far below conventional weak instrument testing thresholds, such as the  $F > 10$  rule advocated by Staiger and Stock (1997). Thus, the instruments were very weak in

the classic studies. It is well-known that with multiple weak instruments 2SLS is seriously biased towards OLS – see Bound, Jaeger and Baker (1995).

Using more recent data, Lee (2001) constructs a PSID sample five times larger than in the classic studies, and obtains a first-stage  $F$  of 26.3 and a Frisch estimate of .503 (se=.092) for 25-60 year old men.<sup>2</sup> In contrast, in a Monte Carlo exercise where he uses randomly drawn 1/5th sub-samples, he obtains an average Frisch estimate of only .253 (se=.227), illustrating the downward bias in the classic studies.<sup>3</sup> This bias towards OLS when instruments are weak is widely appreciated in the IV literature, although Lee (2001)’s result appears to have done little to alter the conventional wisdom that the Frisch elasticity is small.

We emphasize a different issue – completely unrelated to bias – that has been largely overlooked in the prior literature. In contrast to Lee (2001), we start by focusing on the just identified case where instrument strength exceeds conventional weak IV thresholds. In this case 2SLS is approximately median unbiased, so bias towards OLS is not a concern (see Keane and Neal 2022; Angrist and Kolesár 2021). Nevertheless, we demonstrate that the 2SLS  $t$ -test continues to exhibit very poor behavior in this supposedly benign context.

Specifically, we show that the association between 2SLS estimates and their standard errors that we document is a serious problem for inference using  $t$ -tests even when instruments are strong by conventional standards – such as the  $F > 10$  criterion. The consequence is that the 2SLS  $t$ -test has little power to detect a true positive Frisch elasticity even if instruments are “strong.” This will bias studies that rely on 2SLS  $t$ -tests against concluding the Frisch is large.

The implications of our results go well beyond the present application: Theorists have often advocated using the AR test when instruments are weak, because it is robust to weak instrument problems. But the association between 2SLS estimates and standard errors renders 2SLS standard errors and  $t$ -tests highly misleading even when instrument strength is well above conventional weak IV thresholds. Hence, we argue the AR test should replace the  $t$ -test in 2SLS applications, not only when instruments are weak but even when they are strong.

In the over-identified case two-step GMM suffers from a similar power asymmetry problem. But as we also show, LIML and continuously updated GMM, in conjunction with the conditional likelihood ratio test, resolve the problem.

The outline of the paper is as follows: Section II discusses our NLSY79 data and estimating equation. Section III presents our 2SLS estimates and  $t$ -test results. Section IV presents AR test results. Section V compares the behavior of the AR and  $t$ -tests, and Section VI interprets the evidence in light of that comparison. Section VII presents results with multiple instruments. Section VIII concludes.

<sup>2</sup>Substantively, our estimate of 0.6 is large relative to Lee’s estimate of 0.5, as our sample is much younger and there is growing evidence that the Frisch increases substantially with age - see Keane (2021). For example, French and Jones (2012) obtain an estimate of 0.36 for 40 year old men, increasing to 1.28 for 60 year olds. Imai and Keane (2004) obtain 0.36 for 25 year old men, increasing to 1.96 at age 60.

<sup>3</sup>Similarly, if he limits his analysis to the 1/5 of the PSID data available to the classic studies, he obtains a first stage  $F$  of only 1.06 and a Frisch estimate of .258 which is insignificant (se=.172).

## II. Estimating the Frisch Labor Supply Elasticity

We estimate the Frisch elasticity using data from the National Longitudinal Survey of Youth 1997 (NLSY97). The NLSY97 follows a sample of American youth born in 1980-84. The 8,984 respondents were aged 12-17 when first interviewed in 1997.<sup>4</sup> We use data from rounds 10 through 15, which contain information on labor income and work hours in 2005 to 2010. The regression we run is:

$$(1) \quad \Delta \ln H_{it} = \alpha + \beta \Delta \ln W_{it} + \gamma \mathbf{C}_{it} + \epsilon_{it}$$

where  $H_{it}$  is annual hours worked for respondent  $i$  in year  $t$ ,  $W_{it}$  is the wage, and  $\mathbf{C}_{it}$  is a vector of control variables which includes year dummies (to capture business cycle effects on hours worked) as well as respondent age and race/ethnicity.

Our hours measure is “Total annual hours worked at all civilian jobs during the year in question” while our income measure is “Annual income from wages, salary, commissions, and tips before tax deductions.” We obtain an annual wage measure by taking the ratio of annual income to annual hours. Regressions that involve percentage changes can be quite sensitive to measurement error and outliers, as these can generate extreme percentage changes. So, as is typical in this literature, we implement a number of sample screens designed to eliminate outliers.<sup>5</sup>

Obviously, OLS estimation of (1) fails to identify the Frisch elasticity, as predictable and unpredictable wage changes have different effects on labor supply. A surprise wage increase has both substitution and income effects. In contrast, a predictable wage increase has no income effect (precisely because it was predictable), so it induces a pure substitution effect that increases labor supply. It is this Frisch substitution effect of predictable wage changes we want to estimate.

Our key task then is to choose an instrument that is known to workers at the start of each year, and that generates predictable wage growth during the year. MaCurdy (1981) and many subsequent papers use education as the primary instrument for wage growth. The motivation is that annual wage growth tends to be faster for more educated workers.<sup>6</sup> We adopt a closely related approach: The NLSY97 administered an aptitude test called the Armed Services Vocational Aptitude Battery (ASVAB) to respondents when they were 13 to 18 years old.<sup>7</sup> We find that the ASVAB percentile score is a stronger predictor of wage growth than education, so we use that as our instrument. But the idea is similar: Not surprisingly, wage growth is predictably faster for higher ability workers.

<sup>4</sup>Of that, 6748 is a random sample of the birth cohort while 2236 is an over-sample of minority groups.

<sup>5</sup>Observations were excluded if income was less than \$3,000, the annual wage was less than \$2.70 per hour worked, the total number of hours worked was less than 400 or above 4,160 (roughly 80 hours a week), or if the percentage change in wages from the last year was below -50% or above 70%.

<sup>6</sup>He also used interactions of education and age, to allow the effect of education to differ by age.

<sup>7</sup>The ASVAB measures aptitude in several areas including mathematics, general science, paragraph comprehension, and mechanical skills. It was administered in summer 1997 to spring 1998, when the youth were aged 13 to 18 (those aged 13 to 14 were given an easier version of the test). The NLS grouped respondent’s into three-month age windows and calculated a youth’s percentile rank within his age group.

We did the analysis separately for men and women, as prior literature has shown that their labor supply behavior differs in important ways. Interestingly, the ASVAB score is a much better predictor of wage growth for men than women.<sup>8</sup> For this reason, we decided to focus only on results for men. Our full data set has 5,931 annual observations on 2,100 young men aged 22 to 30 who we observe over 2 to 6 years (the average being 3.8 years).

### III. NLSY97 Estimates of the Frisch Elasticity

Table 1 presents regressions of changes in log hours on changes in log wages, as in equation (1). The first column reports OLS results. The coefficient on the log wage change is -0.42 and highly significant, with a standard error of 0.015.<sup>9</sup> This implies a 10% wage increase is associated with a 4.2% reduction in work hours. There are two reasons for a negative relationship: Of course surprise wage changes may generate income effects that reduce labor supply. But it is implausible that income effects alone could generate such a large negative effect.

A second key factor driving the OLS estimate negative is “denominator bias” arising because the wage is measured as the ratio of earnings to hours. If hours in the denominator are measured with error, it causes a worker’s measured wage to be too low precisely when his measured hours are too high. This induces an (artificial) negative covariance between measured hours and measured wages that drives the estimated elasticity negative. As a result, the OLS estimate cannot be interpreted causally. A second virtue of instrumenting for wage changes is that it also deals with this measurement error problem - see Altonji (1986).

Next we consider the 2SLS results. The second column of Table 1 reports the first stage, where we regress log wage changes on the ASVAB percentile score to construct predictable wage changes. The coefficient is 0.039 and highly significant (standard error 0.012). The effect size is substantial: A male worker in the 100th percentile of ability is predicted to have annual wage growth 3.9 pp higher than a male worker in the 1st percentile. The heteroskedasticity robust F-test for significance of ability in the first stage regression is 10.12, which gives a  $p$ -value of 0.002. So the ASVAB instrument is significant at well above the 1% level, and passes the Staiger and Stock (1997)  $F > 10$  rule of thumb for IV strength.

Notably, however, the  $R^2$  of the first stage regression is only .007, implying a correlation between our predictions and actual wage changes of .084. In fact, the partial  $R^2$  that shows the fraction of wage variation explained by the ASVAB test alone is .002, implying a partial correlation of only .041. This illustrates the point that annual wage growth is very hard to predict. It is important to emphasize,

<sup>8</sup>It is not clear if this is because wages grow relatively faster for high-ability men than for high-ability women – perhaps due to discrimination – or because the ASVAB is not as good a proxy for labor market skills of women. It would be interesting to explore this issue in future research.

<sup>9</sup>All standard errors and F-statistics reported in this paper are heteroskedasticity robust or cluster robust. The cluster robust standard errors account for both heteroskedasticity and serial correlation. They are always slightly smaller, because the errors in the hours change regression exhibit negative serial correlation. Hence the heteroskedasticity robust statistics are slightly more conservative.

TABLE 1—FRISCH ELASTICITY ESTIMATES - NLSY97

	OLS	2SLS 1 <sup>st</sup> Stage	2SLS 2 <sup>nd</sup> Stage	Reduced Form
Dependent Variable:	$\Delta H$	$\Delta W$	$\Delta H$	$\Delta H$
Wage Change	-0.416 (0.015) [0.015]		0.597 (0.403) [0.363]	
ASVAB Ability Score		0.039 (0.012) [0.011]		0.024 (0.011) [0.010]
F-Stat (Hetero- $\sigma$ Robust) <i>p-value</i>		10.12 0.002		4.47 0.035
F-Stat (Cluster Robust) <i>p-value</i>		12.23 0.001		5.64 0.018
$R^2$	0.210	0.007		0.009

*Note: Heteroskedasticity robust standard errors are in parentheses. Clustered standard errors (by individual) are in square brackets. All regressions controls for year effects, age, and race/ethnicity.  $N = 5,931$*

however, that a higher first-stage  $R^2$  would not necessarily be desirable in this context. Measured wage changes contain both unpredictable and measurement error components that we specifically want to filter out, so we expect the  $R^2$  of the first stage regression to be far less than one.

Now consider the second stage 2SLS results, where we regress log hours changes on log predictable wage changes to obtain an estimate of the Frisch elasticity. This is reported in the third column of Table 1. The 2SLS estimate of the Frisch elasticity is 0.597, implying that a 10% predictable wage increase generates a 6% increase in work hours. So the use of 2SLS flips the sign of the coefficient.

This 2SLS estimate is clearly more reasonable: Economic theory predicts a positive Frisch elasticity, as a predictable wage increase should have a positive substitution effect on labor supply. And a Frisch elasticity of 0.6 is well within the range of estimates surveyed in Keane (2011, 2021), although it is clearly towards the high range of estimates for young men.

Notice however, that the (heteroskedasticity robust) standard error on the 2SLS estimate is a substantial .403, giving a  $t$ -statistic of only 1.48 and a  $p$ -value of 0.138. So, while the estimated Frisch elasticity is a substantial 0.6, it is not even significantly different from zero at the 10% level.<sup>10</sup> This imprecision leaves us in a quandry over what we ought to conclude from the analysis.

<sup>10</sup>The cluster robust standard error is slightly smaller, at 0.363, because the serial correlation in the hours change regression is negative. But even then the  $t$ -stat is only 1.65 ( $p$ -value = 0.099).

The imprecision in our 2SLS estimate is a consequence of the fact that the ASVAB score only explains a small part of the variance of wage changes. Because the partial correlation between the ASVAB score and wage changes is .041, the standard error goes up by a factor of 25 when we go from OLS to 2SLS (i.e.,  $1/.041 \approx 25$ ). This imprecision in 2SLS estimates has plagued much of the literature on estimating the Frisch elasticity using the 2SLS approach.

Should we trust the 2SLS results in this case? The first-stage  $F$ -statistic for the ASVAB instrument exceeds the commonly used weak IV threshold of 10 (if only marginally), suggesting the results may be viewed as reliable.<sup>11</sup> However, weak IV tests of the type developed by Staiger and Stock (1997) and Stock and Yogo (2005) are designed to assess bias in 2SLS estimates and size distortions in  $t$ -tests. They say nothing about whether instruments are strong enough for the  $t$ -test to have acceptable power properties.<sup>12</sup> In the next section we present AR test results, and in Section V we assess the relative performance of the two tests in this environment. We show that the  $t$ -test has very poor power properties in the present application, and that the AR test ought to be relied on instead.

#### IV. The Anderson-Rubin Approach

Anderson and Rubin (1949) developed an alternative approach to inference that can also be used to test if our estimate of the Frisch elasticity is significant. The Anderson-Rubin (AR) test relies on a reduced form regression of the outcome of interest on the instrument itself, along with the control variables. In our case this is a regression of the change in log hours on the ASVAB score itself, along with the controls (time, age, race). The AR test judges the Frisch elasticity estimate to be significant if the ASVAB score is significant in the reduced form regression.

The logic of the AR test is simple: A fundamental assumption of the IV method is that the instrument only affects the outcome of interest indirectly through its effect on the endogenous variable. Hence, if the instrument is significant in the reduced form, it implies that the endogenous variable has a causal impact on the outcome of interest. In our case, if the ASVAB score is significant in the reduced form, it implies that predictable wage changes influence work hours.

Of course the ASVAB score could appear significant in the reduced form merely because it somehow affects hours growth directly (not indirectly via its effect on wage growth). That is, the ASVAB score may be significant because the exclusion restriction is violated. But in that case the ASVAB score is not a valid instrument, so the 2SLS estimate and  $t$ -test results are also invalid. The very assumptions that make the IV approach valid also make the AR test valid.

<sup>11</sup>For example, Stock and Watson (2015, p.490) say: “One simple rule of thumb is that you do need not to worry about weak instruments if the first stage  $F$ -statistic exceeds 10.”

<sup>12</sup>Stock and Yogo (2005) proposed critical values for  $F$  based on “worst-case” size distortion in  $t$ -tests. For example, in the exactly identified case, a sample  $F > 8.96$  gives 95% confidence that a two-tailed 5%  $t$ -test will reject  $H_0:\beta = 0$  at a 15% rate or less when the true  $\beta$  is zero. That is, it has a size distortion of no more than 10%. But passing such a test does not imply the  $t$ -test will have acceptable power.

The last column of Table 1 reports the reduced form results. Here, the ASVAB score is clearly significant, with a  $t$ -stat of 2.18 ( $p$ -value 0.035). So we are left with a quandry: The AR test indicates the 2SLS estimate of the Frisch elasticity is significant, while the  $t$ -test says it isn't. Which result should we believe?

The AR test is recommended by theory as clearly superior to the  $t$ -test when instruments are weak, and no worse when instruments are strong - see Andrews, Stock and Sun (2019). This is because the AR test has three major advantages: First, it is “robust” to weak instrument problems, which means a 5% level AR test rejects a true null hypothesis at the correct 5% rate *regardless* of the strength or weakness of the instruments. In contrast, the  $t$ -test suffers from size distortions: If instruments are weak, a 5%  $t$ -test may reject a true hypothesis at rates well above or below 5%, depending on details of the situation. Second, the AR test is unbiased, meaning its power is appropriately minimized when the null corresponds to the true  $\beta$ . Third, Moreira (2009) shows that in the case of a single instrument (as we have here) the AR test is the most powerful unbiased test: If the null hypothesis is false, the AR test will reject the null, and conclude the parameter of interest is significant, at least as frequently as any other unbiased test.<sup>13</sup>

Despite its clear advantages, the AR test has been widely neglected by applied researchers. In fact, with the exception of Lee (2001), it has never been adopted in the large literature on estimating the Frisch elasticity. In the next section we present a numerical experiment based on our data that shows the performance of the AR test is *dramatically* superior in practice. We also present analytical results that lead to the same conclusion.

## V. Monte Carlo Experiment and Power Analysis

In this section we compare the AR test and the  $t$ -test to see which is a more reliable guide to the statistical significance of our Frisch elasticity estimate. To do this, we conduct the following experiment: We start from the NLS sample of  $N=5,931$  observations that we used to generate the estimates in Table 1. We can then “bootstrap” a new artificial dataset by sampling 5,931 observations with replacement from the original sample. We do this 10,000 times to form 10,000 artificial datasets. We then repeat the analysis of Table 1, applying OLS and 2SLS to all 10,000 datasets, and summarize the results in Table 2.<sup>14</sup>

<sup>13</sup>Andrews, Stock and Sun (2019) argue the AR test should be widely adopted by applied researchers. They state its advantages more formally: “In just-identified models ... Moreira (2009) shows that the AR test is uniformly most powerful unbiased. ... Thus, the AR test has (weakly) higher power than any other size- $\alpha$  unbiased test no matter the true value of the parameters. In the strongly identified case, the AR test is asymptotically efficient in the usual sense and so does not sacrifice power relative to the conventional  $t$ -test. ... Since AR confidence sets are robust to weak identification and are efficient in the just-identified case, there is a strong case for using these procedures in just-identified settings.” Moreira and Moreira (2019) extend this optimality result to models with heteroskedasticity and clustering.

<sup>14</sup>By sampling with replacement from the original 5,931 observations we break the panel structure of the data. As a result, the standard errors and  $F$  statistics in Table 2 will mimic the heteroskedasticity robust statistics in Table 1, not the cluster robust statistics.



TABLE 2—RESULTS FROM MONTE CARLO BOOTSTRAP SAMPLES

	OLS		2SLS		First Stage	Reduced Form	
	$\hat{\beta}$	S.E.	$\hat{\beta}$	S.E.	$F$ Statistic	$\hat{\pi}$	S.E.
Median	-0.4163	0.0146	0.5998	0.4013	10.1314	0.0238	0.0112
Mean	-0.4164	0.0146	0.7185	4.7202	11.0923	0.0237	0.0112
Std. Dev.	0.0148	0.0004	3.8636	251.4631	6.4577	0.0111	0.0003

*Note:  $N = 5,931$  for each of the 10,000 samples used to form the results.*

#### A. OLS Estimates and Standard Errors

In Table 2 we see that both the median and mean OLS estimates of  $\beta$  (across all 10,000 datasets) are roughly equal to the (downward biased) value of -0.416 we obtained using the original NLS sample. This is as expected, as our 10,000 “bootstrap” datasets mimic the covariances of the variables in the original NLS sample. The third row of Table 2 reports that the standard deviation of the OLS estimates across the 10,000 artificial samples is 0.015, which equals (to three decimal places) the OLS standard error estimate reported in Table 1. Thus, the estimated OLS standard error is a very good guide to how the OLS estimates actually vary across the different samples.<sup>15</sup>

#### B. 2SLS Estimates and Standard Errors

Now we examine how the 2SLS estimates and standard errors behave. The first thing to note in Table 2 is that the median 2SLS estimate of the Frisch elasticity (across all 10,000 datasets) is 0.600, which is very close to the 2SLS estimate 0.597 we obtained using the original NLS dataset. This is exactly as expected: As our artificial datasets are constructed from our original NLS sample, we can think of the NLS sample as the “population” from which all 10,000 datasets are drawn. In this population, 0.597 is in fact the true value of the Frisch elasticity. We see that the median 2SLS estimate accurately uncovers the true Frisch elasticity value. As Keane and Neal (2022) and Angrist and Kolesár (2021) show, the median bias of 2SLS is negligible at this level of instrument strength.

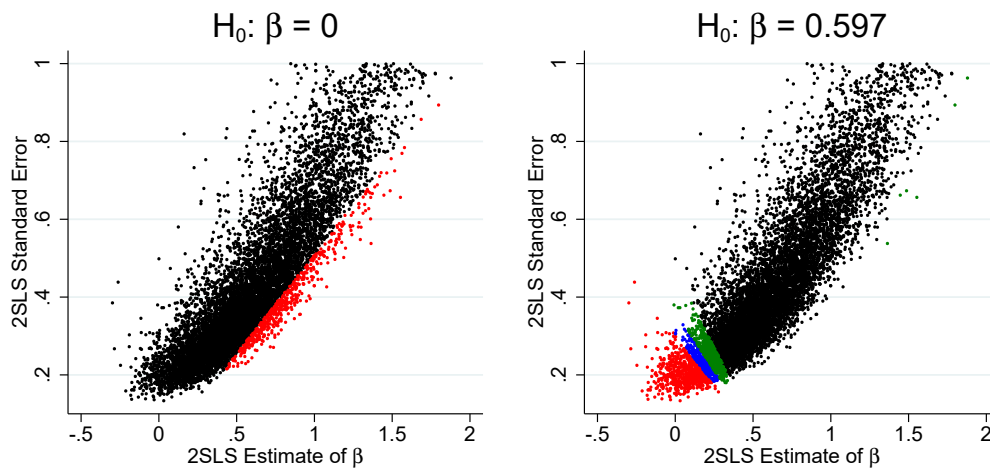
Second, note that the median of the estimated 2SLS standard errors, reported in the first row of Table 2, is 0.401. This agrees closely with the 2SLS standard error estimate of 0.403 in Table 1. However, the actual empirical standard deviation of the 2SLS estimates across the 10,000 data sets is 3.864. In contrast to OLS, the actual variation of the 2SLS estimates bears no resemblance to the estimated 2SLS standard errors. This is our first indication that the 2SLS standard errors are not

<sup>15</sup>Table 2 also reports the mean and median of the estimated OLS standard error across the 10,000 artificial datasets. These are again 0.015. And the variation across samples of this standard error estimate is trivially small. So the estimated standard error in each individual sample is a good guide to the actual variability of the OLS estimates across all samples.

a good guide to the actual variability of the 2SLS estimates across samples.<sup>16</sup> This in turn means that 2SLS  $t$ -statistics – which rely on those standard error estimates – will not be a useful guide to significance of 2SLS estimates.

To further explore the behavior of the 2SLS standard error, Figure 1 plots the 2SLS standard errors against the 2SLS estimates of the Frisch elasticity from each of the 10,000 samples. A striking aspect of the figure is the strong positive association between the 2SLS estimates and their standard errors: The Spearman correlation is an extraordinarily large 0.905. This means that in samples where the estimated Frisch elasticity is larger, the standard error is also larger.<sup>17</sup> As we will see, this pattern has extremely important empirical implications.

FIGURE 1. STANDARD ERROR OF  $\hat{\beta}_{2SLS}$  PLOTTED AGAINST  $\hat{\beta}_{2SLS}$  ITSELF



*Note:* Runs with standard error  $> 1$  are not shown. In the left panel, red dots indicate  $H_0: \beta = 0$  is rejected at the 5% level by a 2SLS  $t$ -test. In the right panel red dots indicate  $H_0: \beta = 0.597$  is rejected at the 5% level. Blue and green indicate 10% and 20%.

The association between 2SLS estimates and their standard errors is not specific to this application. It is a generic but little-appreciated property of the 2SLS estimator. Exact finite sample theory can shed light on this phenomenon. Phillips (1989) derives two key properties of 2SLS in the unidentified case. First, the 2SLS estimator converges in distribution to a scale mixture of normals centered on  $E(\hat{\beta}_{OLS})$ . Second, the 2SLS variance estimator ( $\hat{\sigma}^2$ ) converges in distribution to a quadratic function of  $\hat{\beta}_{2SLS}$ , with a minimum at  $E(\hat{\beta}_{OLS})$ . Thus, the standard error of regression ( $\hat{\sigma}$ ) is minimized when  $\hat{\beta}_{2SLS}$  is close to  $E(\hat{\beta}_{OLS})$ . Of course, the standard error of the regression ( $\hat{\sigma}$ ) is a fundamental driver of the standard

<sup>16</sup>Of course, in the single instrument case the mean and variance of the 2SLS estimator do not exist, which means that if we did many more than 10,000 runs the mean and variance wouldn't converge. Hence, the standard deviation of the 2SLS standard error cannot be bootstrapped.

<sup>17</sup>A graph of OLS standard errors vs. estimates is a spherical cloud, as they have zero correlation.

error of  $\hat{\beta}_{2SLS}$ . Thus, in the unidentified case, the standard error of  $\hat{\beta}_{2SLS}$  tends to be minimized when the estimate is near  $E(\hat{\beta}_{OLS})$ .

Importantly, these properties of 2SLS in the unidentified case still influence the behavior of 2SLS estimates and standard errors in strongly identified models. In fact, Phillips (1989) calls this the “leading case” as it provides the leading term of the series expansion of the density of the estimator in the general case. As a result, even in strongly identified models, the standard error of  $\hat{\beta}_{2SLS}$  tends to be minimized when the estimate is near  $E(\hat{\beta}_{OLS})$ , as we see in Figure 1.

Intuitively, in finite samples where the exogenous instrument happens to be *positively* correlated with the structural error, the 2SLS estimate is shifted toward OLS and the instrument appears stronger, so the estimate seems more precise.<sup>18</sup> Appendix A.1 provides additional mathematical detail on why this pattern arises. For our present purposes, it suffices to note the following: Because of this pattern, large positive 2SLS estimates of the Frisch elasticity will have relatively large standard errors, while estimates near zero will have small standard errors.

This brings us to our key point: The positive association between 2SLS estimates of the Frisch elasticity and their standard errors has important implications for statistical inference. As we now show, this mechanical relationship makes it very difficult for a 2SLS  $t$ -test to detect a true positive Frisch elasticity.

Recall that our 10,000 simulated data sets are constructed so the true value of the Frisch elasticity in these data sets is 0.597. Thus, if the 2SLS  $t$ -test is reliable it should have two properties: First, if we run 5%  $t$ -tests of the hypothesis that the true Frisch elasticity is zero we should reject that false hypothesis at a high rate (indicating the test has good power). Second, if we run 5%  $t$ -tests of the true hypothesis that the Frisch is equal to 0.597 (the true value) we should reject that hypothesis approximately 5% of the time (indicating the test has correct size). Furthermore, those rejections should be evenly split between cases where the estimated Frisch elasticity is above and below the true value.

In the left panel of Figure 1 we shade in red the cases where the 2SLS  $t$ -test rejects the false null hypothesis that the true Frisch elasticity is equal to zero. These are the cases where the ratio of the estimate to the standard error exceeds the 5% critical level of 1.96 (in absolute value). Notice that the red shaded area

<sup>18</sup>A more detailed intuitive explanation for the association between 2SLS estimates and their standard errors is as follows: As we discussed in Section III, in the original NLS sample the partial correlation between the ASVAB score and wage growth is 0.04. But the correlation fluctuates across our 10,000 subsamples due to sampling variation. Two things happen in samples where it is relatively high:

First, the 2SLS standard error estimate is smaller: The higher the correlation between the instrument and the endogenous variable, the stronger the instrument seems to be, and hence the smaller is the 2SLS standard error. 2SLS standard errors are suspect due to this pattern.

Second, the 2SLS estimate is more shifted in the direction of the OLS bias (which is negative). This is because, as we discussed in Section III, if the predictable part of wage growth is small, then a high correlation between the instrument and the endogenous variable is not really a good thing. In samples where that correlation rises above 0.04, the instrument is picking up some of the endogenous part of wage growth that arises due to measurement error and surprise wage growth. This in turn means the 2SLS estimate will be shifted in the direction of the OLS bias (negative).

Putting these two facts together, we get that 2SLS estimates that are most shifted in the direction of the OLS bias (negative) appear to be more precise. This is exactly the pattern we see in Figure 1.

is quite small. In fact, the false null hypothesis is only rejected 5.1% of the time. This is an extremely low level of power. It scarcely exceeds the 5% rate at which a well-behaved 5% level test should reject a *true* null hypothesis.

In the right panel of Figure 1 we shade in red the cases where the 2SLS  $t$ -test rejects the null hypothesis that the true Frisch elasticity is equal to the true value of 0.597. The test rejects the null hypothesis 6.6% of the time. This is not bad when viewed in isolation, as it is not too far from the correct rate of 5%, so the  $t$ -test size distortion is small. But more importantly, the rate of rejecting the true hypothesis that the Frisch equals 0.597 is actually *greater* than the rate of rejecting the false hypothesis that the Frisch equals 0. This is extremely poor behavior for a statistical test: The fact that size exceeds power means the  $t$ -test is uninformative about the true parameter value.

These results illustrate our key point: The population  $F$  easily passes conventional weak instrument testing thresholds of the type developed by Staiger and Stock (1997) and Stock and Yogo (2005). So, as expected, the median bias of 2SLS and the size distortion in the  $t$ -test are trivial. Nevertheless, the 2SLS  $t$ -test results are completely uninformative, as size exceeds power.

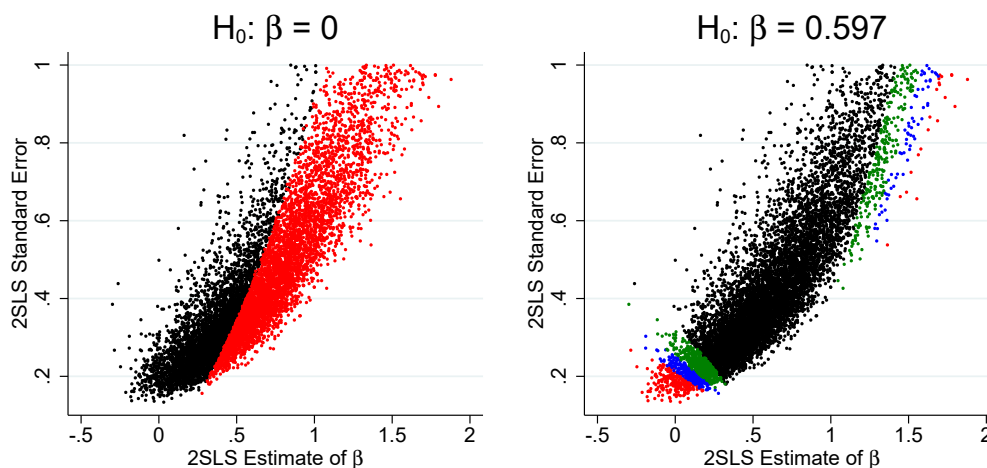
Another notable aspect of the right panel of Figure 1 is that the cases where we reject the null of Frisch = 0.597 are not evenly split between cases where the estimate is above and below the true value. In fact, all the rejections occur when the estimated Frisch elasticity is very small (near zero). This is a direct consequence of the positive association between 2SLS estimates and their standard errors. As large positive estimates of the Frisch elasticity have large standard errors, there is very little chance of concluding a large positive estimate is significant.

We refer to the low power of the 2SLS  $t$ -test to detect a true effect opposite in sign to the OLS bias as the “power asymmetry” problem.

### C. Anderson-Rubin Test Results

Figure 2 reports the same results for the AR test. The contrast with the  $t$ -test is dramatic. We again plot the 2SLS standard errors against the 2SLS estimates, as in Figure 1. But now we plot in red the cases where the AR test rejects the false null hypothesis that  $\beta = 0$  at the 5% level. In the left panel we see that the red region is quite large. The AR test rejects the false null that the Frisch is equal to zero 56.5% of the time. This is a good level of power that is more than ten times greater than the 5.1% rate achieved by the  $t$ -test.

The right panel of Figure 2 shows the rate of rejecting the true null hypothesis that the Frisch equals 0.597. The AR test rejects 4.9% of the time, which is almost exactly equal to the correct 5% rate. Furthermore, we plot in blue and green the cases where 10% and 20% AR tests reject. These rates are 10% and 19.5%, so again almost perfect. This illustrates how the AR test is “robust” in the sense that it has correct size (rejection rates) regardless of the strength or weakness of the instruments. Thus we see that *the AR test has correct size and ten times the power of the  $t$ -test.*

FIGURE 2. STANDARD ERROR OF  $\hat{\beta}_{2SLS}$  PLOTTED AGAINST  $\hat{\beta}_{2SLS}$  ITSELF (AR TEST)

Note: Runs with standard error  $> 1$  are not shown. In the left panel, red dots indicate  $H_0: \beta = 0$  is rejected at the 5% level by the AR test. In the right panel red dots indicate  $H_0: \beta = 0.597$  rejected at the 5% level. Blue and green indicate 10% and 20%.

The only limitation of the AR test is that it doesn't quite generate symmetric rejections when the estimates are above and below the true value. For example, of the 4.9% rejections in the 5% test, 3.6% occur when the estimate is below 0.597 and 1.3% occur when it is above. This is because, like the  $t$ -test, the AR test tends to attribute greater precision to estimates shifted in the (negative) direction of the OLS bias, and less precision to large positive estimates. But this power asymmetry problem is much less severe for the AR test than the  $t$ -test.<sup>19</sup>

These results make it very obvious that in the data environment of our empirical application the AR test provides a far more reliable guide to the significance of the estimate of the Frisch elasticity than does the  $t$ -test. The consequence is that prior work that relied on 2SLS  $t$ -tests will have tended to obtain insignificant results even if the true Frisch elasticity is well above zero.

The superiority of the AR test over the  $t$ -test is not specific to this example. In Keane and Neal (2022) we show that the superiority of the AR test is evident across a wide range of contexts. 2SLS  $t$ -tests perform poorly in general due to the strong association between 2SLS estimates and their standard errors. As a result of this pattern,  $t$ -tests have difficulty detecting true negative (positive) effects when the OLS endogeneity bias is positive (negative). This power asymmetry problem is relevant across a wide range of empirical applications, including cases where instruments are much stronger than here. The AR test is much less susceptible to this problem.

<sup>19</sup>In Keane and Neal (2022) we show the power asymmetry in the AR test vanishes quickly as instruments become stronger. But the  $t$ -test asymmetry remains substantial even with very strong instruments.

*D. Analytical Power Function Comparison: AR vs t-test*

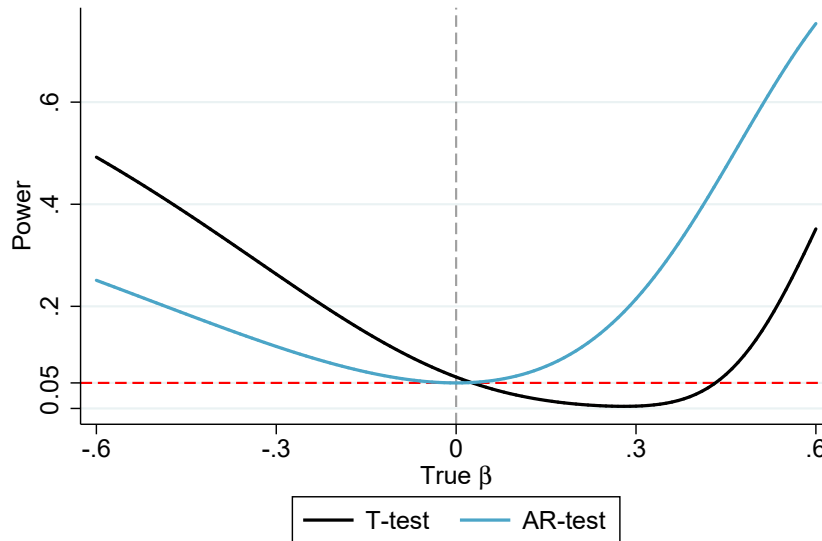
To show the generality of the problem we describe, we now compare the analytical power functions of the AR and  $t$ -tests in an exactly identified linear IV model with *iid* normal errors:

$$(2) \quad \begin{aligned} y_i &= \beta x_i + u_i \\ x_i &= \pi z_i + e_i \quad \text{where} \quad e_i = \rho u_i + \sqrt{1 - \rho^2} \eta_i \end{aligned}$$

where  $u_i \sim iidN(0, 1)$ ,  $\eta_i \sim iidN(0, 1)$ , and  $z_i \sim iidN(0, 1)$ . Thus the instrument  $z$  satisfies  $cov(z, u) = 0$  and  $cov(z, \eta) = 0$ . The parameter  $\rho \in (-1, 1)$  determines the severity of the endogeneity problem, while  $\pi$  determines the strength of the instrument. This *iid* normal setup is not as restrictive as it may first appear, as Andrews, Stock and Sun (2019) show that for any heteroskedastic DGP, there exists a homoskedastic DGP yielding equivalent behavior of 2SLS estimates and test statistics. Any exogenous covariates can be partialled out of  $y$  and  $x$  without changing anything of substance.

Figure 3 compares the power functions of the AR and  $t$ -tests, obtained via the procedure described in the Appendix, with population  $F=10.12$  and  $\rho = -.70$  to mimic our empirical application. The power function is the probability a 5% level test rejects  $H_0: \beta = 0$ , conditional on each alternative true  $\beta$  listed on the  $x$ -axis.

FIGURE 3. POWER OF THE T-TEST VS. AR-TEST WHEN POP.  $F = 10$  ( $\rho = -0.7$ )



*Note: Probability a 5% level test rejects  $H_0 : \beta = 0$ , conditional on each alternative true  $\beta$  listed on the  $x$ -axis. In order to most closely match our application, we set  $\rho = -0.7$  and the population first-stage  $F$ -statistic to 10.12.*

Several features of the power functions are notable. First, the power of the 2SLS 5% level  $t$ -test to reject  $H_0 : \beta = 0$  is close to 5% if the true  $\beta = 0$ . So as expected the  $t$ -test has approximately correct size when population  $F = 10$ . Both Angrist and Kolesár (2021) and Keane and Neal (2022) argue that size distortions in two-tailed 2SLS  $t$ -tests are modest unless instruments are very weak and endogeneity is very severe, and that fact is reflected here.<sup>20</sup>

Second, the severe bias of the  $t$ -test is evident: Its power dips below 5% for a wide range of positive values of  $\beta$ , dipping to near zero when true  $\beta$  is around 0.25 to 0.30. In the model in (2),  $\beta$  is approximately the standard deviation change in  $y$  induced by a one standard deviation change in  $x$ . So effect sizes of 0.25 to 0.30 would be substantial in most applications. Yet the 2SLS  $t$ -test has essentially no power to detect effect sizes in this range. Third, the unbiasedness of the AR test is also evident, as its power is appropriately minimized at  $\beta = 0$ .

Fourth, and most importantly, the AR test clearly has far superior power to detect true positive values of  $\beta$  in this environment. The lack of power of the  $t$ -test in the positive  $\beta$  range is due to the positive association between 2SLS estimates and their standard errors that exists because the OLS bias is negative. This causes larger positive estimates of  $\beta$  to have spuriously inflated standard errors. Conversely, the  $t$ -test *appears* to have better power than AR when the true  $\beta$  is negative, but this is problematic, as it occurs because larger negative estimates of  $\beta$  tend to have spuriously small standard errors.

The Appendix contains additional discussion of analytical power results for different levels of instrument strength. The patterns we discuss here are more pronounced when instruments are weak, but persist even when instruments are far above conventional weak instrument testing thresholds.

## VI. Interpreting the Empirical Results in Light of the Experiment

We now return to the empirical results in Table 1, and assess them based on what we learned from the Monte Carlo and analytical results in Section V. Recall that our 2SLS estimate of the Frisch elasticity is 0.597, but the 2SLS  $t$ -test indicates it is not significantly different from zero at the 5% level. However, the analysis of Section V shows that the  $t$ -test has little power to detect true positive effects of plausible magnitude in this data environment, that is characterized by (i) a first stage  $F$  slightly above 10 and (ii) a correlation between the reduced form residuals of -0.70, so that the OLS bias is strongly negative.

The analysis of Section V showed the AR test is a far more reliable method of inference in this context. The AR test is based on the significance of the instrument (ASVAB) in the reduced form equation for hours changes. The AR test is widely recommended by theorists for use when instruments are weak, because it

<sup>20</sup>Stock and Yogo (2005) derived critical values that the sample  $F$  must surpass in order for a researcher to have high confidence the “worst-case” size distortion in the two-tailed  $t$ -test is modest. But the “worst case” occurs when  $\rho$  is near one or minus one, so endogeneity is extremely severe. For smaller values of  $\rho$  – that are more typical of applications – much lower levels of  $F$  will suffice to render size distortions quite modest. Keane and Neal (2022) and Angrist and Kolesár (2021) also make this point.

is robust to weak instrument problems.<sup>21</sup> But the analysis in Section V revealed that the AR test has greatly superior power properties to the  $t$ -test even when instruments are strong by conventional standards (i.e.,  $F > 10$ ).

The AR test indicates that our Frisch elasticity estimate of 0.597 is significant at the 3.5% or 1.8% level, depending on whether we rely on the heteroskedasticity or cluster robust version of the test. We can also invert the AR test to obtain a weak instrument robust confidence interval, as discussed in Anderson and Rubin (1949).<sup>22</sup> Using cluster robust statistics we obtain a 95% confidence interval for the Frisch elasticity of 0.082 to 2.03, which is clearly bounded above zero, and covers most of the range often used to calibrate macro models.

Our Frisch estimate of 0.60 is well above those obtained in the classic studies. As we discuss in the conclusion, it is in line with more recent estimates obtained for young men using a variety of data sources and methods. But going beyond this particular estimation result, our broader point is that the power problems that afflict the 2SLS  $t$ -test make it difficult to detect a true positive Frisch elasticity. Thus, we argue it is important that future work on estimating the Frisch elasticity should rely on the AR test rather than the  $t$ -test.

Our results here also provide lessons that are useful in a broader context. In general, 2SLS  $t$ -tests have low power to detect true effects that are opposite in sign to the OLS bias, even if instruments are very strong by conventional standards – see Appendix A for details. For this reason we argue the AR test should replace the  $t$ -test for inference in exactly-identified linear IV models, not only when instruments are weak but even when they are strong. Next we consider over-identified models, where other alternatives to the  $t$ -test are available.

## VII. Results Based on Multiple Instruments

As we discussed in Section II, much of the prior work on estimating the Frisch elasticity used polynomials in education and age as instruments for wage growth, but we rely on the ASVAB score as we find it is a stronger instrument in the first stage of 2SLS. In this section we consider using both education and the ASVB score as instruments. In order to keep the sample identical to that in Table 1, we code education as zero if it is missing, and introduce a dummy for missing education as an additional instrument. As we see in the first column of Table 3 both the ASVAB score and education are significant in the first stage of 2SLS, suggesting they capture somewhat different dimensions of ability.<sup>23</sup>

We also report two versions of the partial  $F$ -statistic for joint significance of the instruments in the first stage (heteroskedasticity and cluster robust), as well as the Oleva-Pfleuger effective  $F$ -test for weak instruments in a non-*iid* setting.

<sup>21</sup>That is, the AR test is guaranteed to have correct size when instruments are weak, yet it is guaranteed to be no less powerful than the  $t$ -test when instruments are strong.

<sup>22</sup>The basic idea of AR test inversion is to run regressions of  $y - xb$  on the instrument and control variables, and find the lower and upper cutoffs for  $b$  where the AR test  $p$ -value is exactly .05.

<sup>23</sup>To be precise, the  $p$ -values for education are .047 or .071 based on the cluster robust or heteroskedasticity robust standard error, respectively.



TABLE 3—FRISCH ELASTICITY - OVER-IDENTIFIED MODELS

Dependent Variable	2SLS	2SLS	Reduced	GMM-2S
	1 <sup>st</sup> Stage	2 <sup>nd</sup> Stage	Form	2 <sup>nd</sup> Stage
	$\Delta W$	$\Delta H$	$\Delta H$	$\Delta H$
Wage Change		1.017 (0.481) [0.442]		0.896 (0.474) [0.433]
ASVAB Ability Score	0.028 (0.014) [0.012]		-0.007 (0.017) [0.016]	
Education	0.002 (0.001) [0.001]		0.006 (0.003) [0.002]	
Education Missing	0.033 (0.034) [0.035]		0.033 (0.044) [0.043]	
F-Stat (Hetero- $\sigma$ Robust)	4.31		4.21	
<i>p-value</i>	0.005		0.006	
F-Stat (Cluster Robust)	5.14		4.75	
<i>p-value</i>	0.002		0.003	
Olea-Pfleuger Effective F	4.57			
Exogeneity Test (Sargan or J)		4.03		3.19
<i>p-value</i>		0.133		0.203
$R^2$	0.008		0.013	

Note: ‘GMM-2S’ refers to 2-step GMM. Heteroskedasticity robust standard errors are in parentheses. Clustered standard errors are in square brackets. All regressions control for year effects, age, and race/ethnicity.  $N = 5,931$ .

These statistics range from 4.3 to 5.1, so they are well below conventional weak instrument testing thresholds.<sup>24</sup> Thus weak instruments are clearly a concern and the 2SLS  $t$ -test cannot be viewed as reliable.

The 2SLS estimate of the Frisch elasticity is 1.017, which is much larger than the estimate of 0.597 we obtained in Table 1. Notably, the heteroskedasticity robust standard error increases from 0.403 to 0.481, so  $t=2.12$  ( $p=.034$ ) and a 5%  $t$ -test judges our estimate significant.<sup>25</sup> It may seem surprising that the 2SLS standard error increases despite the efficiency gain from adding an additional relevant instrument in the first-stage. But the increase in the Frisch estimate from

<sup>24</sup>As Andrews, Stock and Sun (2019) note, in over-identified models it is inappropriate to use a heteroskedasticity-robust or conventional  $F$ -test to assess instrument strength in non-homoskedastic settings. They suggest the Olea and Pflueger (2013) effective first-stage  $F$ -statistic. In the single instrument case we considered in Sections III-VI, this reduces to the conventional heteroskedasticity-robust  $F$ .

<sup>25</sup>The cluster robust standard error increases from 0.363 to 0.442, giving  $t=2.30$  ( $p=.021$ ).

.597 to 1.017 moves us further from the OLS bias. As we discussed in Section V.B, this causes the 2SLS standard error of regression to increase mechanically, which, in turn, tends to inflate the standard error of the 2SLS estimate.

Now consider the AR test, which in the over-identified case is simply the  $F$ -test for joint significance of the three instruments in the reduced form. The cluster robust AR test is 4.75, with a  $p$  value of .0026. Moreover, the AR test is not the most powerful test in the over-identified case: the weak instrument robust conditional likelihood ratio (CLR) test of Moreira (2003) is more efficient.<sup>26</sup> The cluster-robust CLR test is 11.03. It is  $\chi^2(1)$  so the  $p$ -value is .0012. Thus the evidence for a positive Frisch elasticity based on the robust tests is very strong.

Here we see a milder version of the pattern in Table 1: The 2SLS  $t$ -test implies the Frisch elasticity estimate is (just) significant at the 5% level, while the weak instrument robust statistics (the AR and CLR tests) imply much higher levels of confidence. The relative weakness of the 2SLS  $t$ -test result is again attributable to the positive covariance between 2SLS estimates and standard errors, which makes it difficult for 2SLS  $t$ -tests to detect a positive Frisch elasticity.

We can obtain a Sargan test of the 2SLS over-identifying restrictions by regressing the 2SLS residuals on the full set of instruments and exogenous variables. The  $NR^2$  of the regression is distributed  $\chi^2(K)$  under the null  $\beta = 0$ , where  $K$  is the number of over-identifying restrictions. As we see in Table 3, the Sargan test statistic is 4.03. It is distributed  $\chi^2(2)$  so the  $p$ -value is .133. Thus we cannot reject the exogeneity of the instruments. This is important, as a failure of the over-identification test would invalidate both the AR and CLR tests, as it would suggest the instruments may be significant in the reduced form merely because they affect hours changes directly (rather than only indirectly via wages as 2SLS assumes). So the AR, CLR and Sargan statistics should be evaluated in conjunction. Of course, failure of the exogeneity test would invalidate 2SLS  $t$ -test results as well, so this is not a disadvantage of these robust tests relative to the  $t$ -test.

If we invert the AR test (cluster robust  $F$  version) we obtain a 95% confidence interval for the Frisch elasticity of 0.245 to 4.306, while inverting the CLR test gives 0.269 to 4.461.<sup>27</sup> These intervals sit comfortably above zero, and cover the range of values typically used to calibrate macro models.

The last column of Table 3 reports the two-step GMM results. Two-step GMM and 2SLS are equivalent under homoskedasticity. The two-step GMM estimate of 0.896 is slightly smaller than the 2SLS estimate, suggesting that heteroskedasticity

<sup>26</sup>We explain the CLR test in detail in Keane and Neal (2022). It is based on the reduced form system in which the two endogenous variables, hours changes and wage changes, are regressed on all the exogenous variables. The instrument exclusion condition implies the coefficients on the excluded instruments in the hours equation are  $\beta$  times those in the wage equation. The CLR test assesses the deterioration in the log-likelihood of the system when the constraint  $\beta = 0$  is imposed. The likelihood here takes the simple form of a sum of squared residuals, so one is actually assessing how much the residual variance increases. The test is  $\chi^2(1)$ . It was originally proposed by Anderson and Rubin (1949). Moreira (2003) showed how to adjust the critical values of the test so it is robust to weak instruments.

<sup>27</sup>We use the Stata command developed by Finlay and Magnusson (2009) to implement the cluster robust version of the CLR test and to do the inversion.

has only a modest impact on the results. The Hansen-J test has a p-value of 0.203, so again we cannot reject exogeneity of the instruments.

To assess the relative performance of the AR, CLR and  $t$ -tests in the three instrument case, we ran a Monte Carlo analysis like that of Section V, but using the 2SLS estimated model in Table 3 as the data generating process.<sup>28</sup> In terms of power, we find that a 5%  $t$ -test rejects the false null  $H_0:\beta=0$  at a 60.2% rate, compared to 88.4% for the AR test, and 94.7% for CLR. So the ranking is as expected. The extra instruments increase power substantially.

Importantly, the Monte Carlo shows that two-step GMM suffers from the same power asymmetry as 2SLS, due to positive association between GMM estimates and their standard errors. The GMM  $t$ -test is therefore unreliable as well.

TABLE 4—FRISCH ELASTICITY - ALTERNATIVE ESTIMATORS

	LIML	GMM-CU
Wage Change	1.494 (0.802) [0.738]	1.310 (0.548) [0.487]
CLR Test <i>p-value</i>	11.03 0.001	
$S$ Statistic <i>p-value</i>		20.47 0.000
Exogeneity Test (Sargan or $J$ ) <i>p-value</i>	3.22 0.200	3.00 0.224

*Note: 'GMM-CU' refers to continuously updated GMM. Heteroskedasticity robust standard errors are in parentheses and clustered standard errors are in square brackets. All regressions controls for year effects, age, and race/ethnicity.  $N = 5,931$ .*

It is well-known that in the over-identified case 2SLS is seriously biased towards OLS when instruments are weak. Lee (2001) argues that this may have caused the classic studies to obtain estimates of the Frisch elasticity near zero. The limited information maximum likelihood (LIML) estimator of Anderson and Rubin (1949) does not suffer from this bias problem, so it is often recommended with multiple weak instruments. But furthermore, in Keane and Neal (2022) we show that LIML performs much better than 2SLS even when instruments are strong by conventional standards (e.g., when thresholds like  $F > 10$  are satisfied). Table 4 presents LIML results. The LIML estimate of the Frisch elasticity is 1.494, which is indeed larger than the 2SLS estimate of 1.017.

<sup>28</sup>In the one instrument case the instrument is uncorrelated with the 2SLS residuals. So when we treat the full sample as the “population,” the instrument has zero population covariance with the structural error by construction. But in the over-identified case the instruments do have small correlations with the 2SLS residuals. We need to partial out those correlations to set up the experiment.

The LIML standard error is substantial, implying that the LIML estimate is at best marginally significant. But the LIML standard error is no more reliable than the 2SLS standard error, because, as we explain in Keane and Neal (2022), the LIML estimates and standard errors have the same positive association that afflicts 2SLS estimates and standard errors. CLR is the correct test to use in conjunction with LIML, and it gives a highly significant p-value of 0.0012.<sup>29</sup>

Finally, Table 4 also reports continuously updated GMM results. LIML and GMM-CU are identical under homoskedasticity (see Hansen, Heaton and Yaron 1996), so the reason for reporting GMM-CU is to gain more efficient estimates in the presence of heteroskedasticity and clustering. We see that the GMM-CU estimate of 1.310 is slightly smaller than the LIML estimate.

Stock and Wright (2000) develop a weak instrument robust test that generalizes the AR test to the GMM case. This “S-statistic” is the GMM objective function evaluated at  $\hat{\beta}=0$ . For GMM-CU we find  $S=20.47$ . The test is distributed  $\chi^2(3)$  so the p-value is .0001 and the Frisch estimate is highly significant. Finally, we consider Hansen’s test of over-identifying restrictions. As we see in Table 4 the J-test has  $p > 0.20$ , indicating we cannot reject the exogeneity of the instruments. This is important, as a failure of the J-test would invalidate the S test.<sup>30</sup>

In summary, in the over-identified case we obtain larger estimates of the Frisch elasticity than in the exactly-identified case. The estimates range from 0.90 to 1.49 depending on the estimator. The weak instrument robust AR, CLR and S tests indicate that the 2SLS, LIML and GMM estimates of the Frisch elasticity are highly significant. The CLR 95% confidence interval is 0.269 to 4.461, so it sits comfortably above zero but covers a wide range.

## VIII. Conclusion

The magnitude of the Frisch labor supply elasticity – how work hours respond to predictable wage changes – lies at the center of many economic policy debates, because the pure substitution effect measured by the Frisch is a vital input into tax policy. For example, higher values of the Frisch imply lower optimal tax rates on labor income. Because of its importance, there is a large literature estimating the Frisch elasticity using instrumental variable methods.

Classic studies that attempted to estimate the Frisch elasticity typically found it to be small and insignificant. But these studies suffered from two key problems: First, they were plagued by weak instrument problems, as it is hard to find instruments that are good predictors of wage growth. Second, they relied on estimation methods and inferential procedures (2SLS and the associated  $t$ -test) that are biased towards finding the Frisch is small and insignificant.

<sup>29</sup>The CLR test result for LIML is identical to the result we reported earlier for 2SLS. This is because the CLR test relies on the LIML estimate in either case.

<sup>30</sup>In the single endogenous variable,  $K$  instrument case, the Sargan and J-tests have power to detect if at least one instrument is endogenous, provided the model is over-identified, which means at least two instruments must be relevant. But power of these tests will be low if  $K-1$  instruments are weak.

We revisit this issue, focusing on two improvements: First, our main instrument for wage growth is the ASVAB ability test, which is a much stronger predictor of wage growth than the education and age variables used in classic studies. It generates a first-stage  $F$  statistic of 10.12, which exceeds conventional thresholds for an acceptably strong instrument. Second, we rely on inferential procedures that are robust to weak instrument problems and have better power properties.

We estimate a fairly large Frisch elasticity of 0.597 for young men using data from the NLSY97. But, as is typical of this literature, the 2SLS standard error is 0.403, implying our estimate is very imprecise. As a result, a 2SLS  $t$ -test cannot reject the hypothesis that the Frisch is zero at conventional levels – a result that is typical of many prior papers. However, we present Monte Carlo and analytical results showing the  $t$ -test is a very poor guide to inference in this context.

We show that the 2SLS  $t$ -test has little power to detect a true positive Frisch elasticity due to a strong *positive* association between 2SLS estimates and their standard errors that arises when the OLS bias is *negative* – as it is here. This causes positive estimates of the Frisch elasticity to have artificially inflated standard errors. In fact, we show that the power of the  $t$ -test is so poor that it has only about a 5.1% chance of detecting a true Frisch elasticity as large as 0.597.

The power asymmetry that afflicts the  $t$ -test (poor power to detect effects opposite in sign to the OLS bias) has not been noted in prior literature.<sup>31</sup> Fortunately, the AR test of Anderson and Rubin (1949) avoids this problem. A Monte Carlo experiment shows that in our data environment the AR test has ten times the power of the  $t$ -test to detect a positive Frisch elasticity. The AR test indicates that our estimate of the Frisch elasticity is significant at the 3.5% level.

Theorists commonly recommend using AR rather than the  $t$ -test when instruments are weak; see, e.g., Andrews, Stock and Sun (2019). This is because, unlike the  $t$ -test, the AR test is robust to weak instrument problems, meaning it always has correct size. The new twist here is our finding that the AR test has far superior power properties to the  $t$ -test even in contexts where instrument strength is well above conventional weak IV test thresholds (like the  $F > 10$  rule of thumb).

Our estimated Frisch elasticity of 0.597 for young men is large compared to the median estimate of 0.17 across 11 classic studies of men surveyed in Keane (2011). It is more in line with later studies by Lee (2001) and Ziliak and Kniesner (2005) who obtain estimates of 0.50 to 0.54. More recent studies surveyed in Keane (2021) present accumulating evidence that the Frisch elasticity increases substantially with age.<sup>32</sup> The seven studies surveyed give a mean (median) estimate of 0.58 (0.45) for young men, which is similar to our estimate here.<sup>33</sup>

<sup>31</sup>The  $t$ -test power asymmetry has not been noted because it is unusual to study the power of tests *conditional* on estimates, and conventional power curves do not reveal the sign pattern of rejections. In fact, Angrist and Kolesár (2021) reject the whole concept, stating “[the asymmetry] does not make conventional frequentist inference unreliable. The conventional standard for reliability of inference is the accuracy of confidence interval coverage, gauged without conditioning on parameter estimates.”

<sup>32</sup>See Imai and Keane 2004, Borella, De Nardi and Yang 2019, Erosa, Fuster and Kambourov 2016, French 2005, French and Jones 2012, Iskhakov and Keane 2021, and Keane and Wasi 2016.

<sup>33</sup>This increases substantially to a mean (median) of 1.56 (1.45) for 60 year-old men.

We also estimate over-identified models that use both education and ASVAB as instruments for wage growth. This gives larger estimates of the Frisch elasticity. In the over-identified case, two-step GMM suffers from the same association between estimates and standard errors as 2SLS, rendering  $t$ -tests unreliable. To avoid this problem, we recommend using LIML or continuously updated GMM, in conjunction with the conditional likelihood ratio (CLR) test. CLR is robust to weak instruments, avoids the power asymmetry of the  $t$ -test, and is more efficient than AR in the over-identified case. Our GMM-CU estimate of the Frisch elasticity is 1.31, and the CLR p-value is .0012. If we invert the CLR test we obtain a 95% confidence interval for the Frisch elasticity of 0.269 to 4.461, which covers the range of 0.50 to 2.0 that is typically used to calibrate macro models.

There have been several attempts to reconcile the small and insignificant 2SLS estimates of the Frisch elasticity obtained in classic micro data studies with the larger values used in macro calibrations. As Keane and Rogerson (2012, 2015) discuss, these fall into two broad categories: One set of explanations, exemplified by Imai and Keane (2004) and Domeij and Floden (2006) argues that estimation of equation (1) gives downward biased estimates of the Frisch due to problems created by human capital accumulation or liquidity constraints. The other set of explanations, exemplified by Chang and Kim (2006) and Rogerson and Wallenius (2009), argue that, once one accounts for the participation margin of labor supply and aggregation issues, it is possible for the macro level Frisch elasticity to be large even if the micro level elasticity is small.<sup>34</sup>

Our argument here is complementary but new in that we criticise the micro-econometric literature on its own terms: Suppose the assumptions necessary for a 2SLS regression of hours changes on wage changes to deliver consistent estimates of the Frisch elasticity do in fact hold. Even then, we show that the econometric methods that have been used to draw inferences from those estimates are inherently biased against finding the Frisch is both large and significant.

The power asymmetry problem that afflicts the 2SLS  $t$ -test is relevant beyond the Frisch application. In a different context, where the OLS bias is positive, 2SLS standard errors on *positive* estimates would be spuriously precise, making it difficult for 2SLS  $t$ -tests to detect a true *negative*. In the classic application of instrumental variables to estimate a treatment effect given positive selection into treatment, 2SLS  $t$ -tests have difficulty detecting true negative effects, violating a “first do no harm” principle in policy evaluation.

Furthermore, the association between 2SLS estimates and their standard errors that generates the  $t$ -test power asymmetry problem does not vanish as instrument strength increases. Thus, we argue that researchers ought to adopt the AR test or CLR test in lieu of the  $t$ -test even when instruments are quite strong.

<sup>34</sup>More recently, Gottlieb, Onken and Valladares-Esteban (2021) show a large Frisch elasticity can be reconciled with modest reactions to tax holidays due to a combination of income and equilibrium effects.

## ACKNOWLEDGEMENTS

We thank Peter Phillips, Josh Angrist, Michal Kolesar and Isaiah Andrews for helpful comments. This research was supported by Australian Research Council (ARC) grants DP210103319 and CE170100005.

## REFERENCES

- Altonji, J.G.** 1986. “Intertemporal substitution in labor supply: Evidence from micro data.” *Journal of Political Economy*, 94(3, Part 2): S176–S215.
- Anderson, T.W., and H. Rubin.** 1949. “Estimation of the parameters of a single equation in a complete system of stochastic equations.” *Annals of Mathematical statistics*, 20(1): 46–63.
- Andrews, I., J. Stock, and L. Sun.** 2019. “Weak instruments in instrumental variables regression: Theory and practice.” *Annual Review of Economics*, 11: 727–753.
- Angrist, Joshua, and Michal Kolesár.** 2021. “One instrument to rule them all: The bias and coverage of just-id iv.” National Bureau of Economic Research.
- Borella, M., M. De Nardi, and F. Yang.** 2019. “Are marriage-related taxes and Social Security benefits holding back female labor supply?” National Bureau of Economic Research.
- Bound, J., D. Jaeger, and R. Baker.** 1995. “Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak.” *Journal of the American Statistical Association*, 90(430): 443–450.
- Browning, Martin, Angus Deaton, and Margaret Irish.** 1985. “A profitable approach to labor supply and commodity demands over the life-cycle.” *Econometrica: journal of the econometric society*, 503–543.
- Chang, Y., and S. Kim.** 2006. “From individual to aggregate labor supply: A quantitative analysis based on a heterogeneous agent macroeconomy.” *International Economic Review*, 47(1): 1–27.
- Conesa, J.C., S. Kitao, and D. Krueger.** 2009. “Taxing capital? Not a bad idea after all!” *American Economic Review*, 99(1): 25–48.
- Domeij, D, and M. Floden.** 2006. “The labor-supply elasticity and borrowing constraints: Why estimates are biased.” *Review of Economic dynamics*, 9(2): 242–262.
- Erosa, A., L. Fuster, and G. Kambourov.** 2016. “Towards a micro-founded theory of aggregate labour supply.” *The Review of Economic Studies*, 83(3): 1001–1039.
- Finlay, K., and L.M. Magnusson.** 2009. “Implementing weak-instrument robust tests for a general class of instrumental-variables models.” *The Stata Journal*, 9(3): 398–421.
- French, E.** 2005. “The effects of health, wealth, and wages on labour supply and retirement behaviour.” *The Review of Economic Studies*, 72(2): 395–427.
- French, Eric, and John Jones.** 2012. “Public pensions and labor supply over the life cycle.” *International Tax and Public Finance*, 19(2): 268–287.
- Gottlieb, C., J. Onken, and A. Valladares-Esteban.** 2021. “On the Measurement of the Elasticity of Labour.” *European Economic Review*, 139.
- Hansen, Lars Peter, John Heaton, and Amir Yaron.** 1996. “Finite-sample properties of some alternative GMM estimators.” *Journal of Business & Economic Statistics*, 14(3): 262–280.
- Imai, S., and M.P. Keane.** 2004. “Intertemporal labor supply and human capital accumulation.” *International Economic Review*, 45(2): 601–641.
- Iskhakov, Fedor, and Michael Keane.** 2021. “Effects of taxes and safety net pensions on life-cycle labor supply, savings and human capital: The case of Australia.” *Journal of Econometrics*, 223(2): 401–432.

- Keane, Michael P, and Nada Wasi.** 2016. "Labour supply: the roles of human capital and the extensive margin." *The Economic Journal*, 126(592): 578–617.
- Keane, M.P.** 2011. "Labor supply and taxes: A survey." *Journal of Economic Literature*, 49(4): 961–1075.
- Keane, M.P.** 2021. "Recent Research on Labor Supply: Implications for Tax and Transfer Policy." *Labour Economics*, 102026.
- Keane, M.P., and R. Rogerson.** 2012. "Micro and macro labor supply elasticities: A reassessment of conventional wisdom." *Journal of Economic Literature*, 50(2): 464–76.
- Keane, M.P., and R. Rogerson.** 2015. "Reconciling micro and macro labor supply elasticities: A structural perspective." *Annu. Rev. Econ.*, 7(1): 89–117.
- Keane, M.P., and T. Neal.** 2022. "The Role of Instrument Strength in IV Estimation and Inference: A Guide to Theory and Practice." *Journal of Econometrics*, forthcoming. Available at SSRN: <https://ssrn.com/abstract=4294302>.
- Lee, Chul-In.** 2001. "Finite sample bias in IV estimation of intertemporal labor supply models: is the intertemporal substitution elasticity really small?" *Review of Economics and Statistics*, 83(4): 638–646.
- Lee, David S, Justin McCrary, Marcelo J Moreira, and Jack Porter.** 2022. "Valid t-ratio Inference for IV." *American Economic Review*, 112(10): 3260–90.
- MaCurdy, T.E.** 1981. "An empirical model of labor supply in a life-cycle setting." *Journal of political Economy*, 89(6): 1059–1085.
- Moreira, Humberto, and Marcelo J Moreira.** 2019. "Optimal two-sided tests for instrumental variables regression with heteroskedastic and autocorrelated errors." *Journal of Econometrics*, 213(2): 398–433.
- Moreira, M.J.** 2003. "A conditional likelihood ratio test for structural models." *Econometrica*, 71(4): 1027–1048.
- Moreira, M.J.** 2009. "Tests with correct size when instruments can be arbitrarily weak." *Journal of Econometrics*, 152(2): 131–140.
- Olea, J.L.M., and C. Pflueger.** 2013. "A robust test for weak instruments." *Journal of Business & Economic Statistics*, 31(3): 358–369.
- Phillips, Peter CB.** 1989. "Partially identified econometric models." *Econometric Theory*, 5(2): 181–240.
- Prescott, E.C.** 2006. "Nobel lecture: The transformation of macroeconomic policy and research." *Journal of Political Economy*, 114(2): 203–235.
- Rogerson, R., and J. Wallenius.** 2009. "Micro and macro elasticities in a life cycle model with taxes." *Journal of Economic theory*, 144(6): 2277–2292.
- Staiger, D., and J. Stock.** 1997. "Instrumental variables regression with weak instruments." *Econometrica*, 65(3): 557–586.
- Stock, J., and M. Watson.** 2015. *Introduction to econometrics (3rd global ed.)*. Pearson Education.
- Stock, J., and M. Wright.** 2000. "GMM with Weak Identification." *Econometrica*, 68(5): 1055–96.
- Stock, J., and M. Yogo.** 2005. "Testing for weak instruments in linear IV regression." *Identification and inference for econometric models: Essays in honor of Thomas Rothenberg*, 80–108.
- Ziliak, James P, and Thomas J Kniesner.** 1999. "Estimating life cycle labor supply tax effects." *Journal of Political Economy*, 107(2): 326–359.
- Ziliak, James P, and Thomas J Kniesner.** 2005. "The effect of income taxation on consumption and labor supply." *Journal of Labor Economics*, 23(4): 769–796.



## Appendix: Analytical Power Calculations for the AR and $t$ -tests

Consider the following just-identified *iid*-normal linear IV model:

$$(A1) \quad \begin{aligned} y_i &= \beta x_i + u_i \\ x_i &= \pi z_i + e_i \quad \text{where } e_i = \rho u_i + \sqrt{1 - \rho^2} \eta_i \\ u_i &\sim iidN(0, 1), \eta_i \sim iidN(0, 1), z_i \sim iidN(0, 1) \end{aligned}$$

The power of both the AR and  $t$ -tests depends on three parameters: the true  $\beta$ , the degree of endogeneity  $\rho$ , and the population  $t$ -statistic on  $z$  in the first-stage regression, which we denote  $\lambda$  (= square root of population  $F$ ). The power of the AR test (i.e., rate of rejecting  $H_0: \beta=0$  as a function of the true  $\beta$ ) is simply:

$$(A2) \quad Power_{AR}(\beta|\lambda, \rho) = \Phi(\lambda D - z_{1-\alpha/2}) + \Phi(-z_{1-\alpha/2} - \lambda D)$$

where  $\Phi$  is the standard normal cdf,  $D = \beta/\sqrt{1 + 2\rho\beta + \beta^2}$ , and  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of the standard normal distribution. We set  $\alpha = 0.05$ .

To obtain the power function of the of the  $t$ -test we follow the analysis in Stock and Yogo (2005), Lee et al. (2022) and Angrist and Kolesár (2021). The power of the two-tailed 2SLS  $t$ -test is given by the integral:

$$(A3) \quad Power_t(\beta|\lambda, \rho) = \int_{-\infty}^{\infty} \left( \mathbb{I}\{t^2 \geq (1 - \rho_0^2)z_{1-\alpha/2}^2\} f(t, D, \lambda, \rho_0) + \mathbb{I}\{t^2 \geq z_{1-\alpha/2}^2\} \right) \phi(t - \lambda) dt$$

where  $\phi$  is the standard normal pdf,  $\rho_0 = (\rho + \beta)/\sqrt{1 + 2\rho\beta + \beta^2}$ , and:

$$(A4) \quad f(t, D, \lambda, \rho_0) = \Phi\left(\frac{a_2 - \lambda D - \rho_0(t - \lambda)}{\sqrt{1 - \rho_0^2}}\right) - \Phi\left(\frac{a_1 - \lambda D - \rho_0(t - \lambda)}{\sqrt{1 - \rho_0^2}}\right),$$

$$a_1 = \frac{\rho_0 z_{1-\alpha/2}^2 t - |t| z_{1-\alpha/2} \sqrt{t^2 - (1 - \rho_0^2) z_{1-\alpha/2}^2}}{z_{1-\alpha/2}^2 - t^2},$$

$$a_2 = \frac{\rho_0 z_{1-\alpha/2}^2 t + |t| z_{1-\alpha/2} \sqrt{t^2 - (1 - \rho_0^2) z_{1-\alpha/2}^2}}{z_{1-\alpha/2}^2 - t^2}.$$

The integral in (A3) must be evaluated numerically.

To form Fig. 3 in the text, set  $\lambda = 3.186$  and  $\rho = -0.7$  to mimic the empirical results in Section V, as  $\lambda = 3.186$  corresponds to a population  $F$  of 10.12. The power asymmetry of the  $t$ -test is evident in Fig. 3, as it has little power to detect a wide range of true positive  $\beta$  values because the OLS bias is negative ( $\rho = -0.70$ ).

### A1. Explaining the Power Asymmetry of the $T$ -test

Next we give a simple mathematical explanation of the power asymmetry in the  $t$ -test. Recall that the 2SLS estimator of  $\beta$  is given by:

$$(A5) \quad \hat{\beta}_{2SLS} = \frac{\sum_{i=1}^n z_i y_i}{\sum_{i=1}^n z_i x_i} = \beta + \frac{\sum_{i=1}^n z_i u_i}{\sum_{i=1}^n z_i x_i} = \beta + \frac{\widehat{cov}(z, u)}{\widehat{cov}(z, x)}$$

We assume without loss of generality that the population covariance between the instrument and the endogenous variable is positive,  $cov(z, x) > 0$ . To simplify, we further assume that  $\widehat{cov}(z, x) > 0$ , so the sign of the coefficient on  $z$  in the first-stage regression is correct. Violation of this condition is extremely rare if the instrument is reasonably strong. For instance, in our Monte Carlo experiment in Section V, where  $F=10.12$ , first-stage sign is correct in all 10,000 replications.

Given that  $\widehat{cov}(z, x) > 0$ , equation (A5) makes clear that the sign of  $\widehat{cov}(z, u)$ , the sample covariance between the instrument and the structural error, determines whether  $\hat{\beta}_{2SLS}$  lies above or below the true  $\beta$ .<sup>35</sup> Recall that the sign of  $\rho$  determines the sign of the OLS bias. Thus, if  $\rho \widehat{cov}(z, u)$  is positive the 2SLS estimate is shifted towards the OLS bias, and vice versa.

Our key point is that a larger sample realization of  $\rho \widehat{cov}(z, u)$  also drives up the sample covariance between the instrument and the endogenous variable, making the instrument appear spuriously strong, as is obvious because:

$$(A6) \quad \widehat{cov}(z, x) = \pi \widehat{var}(z) + \rho \widehat{cov}(z, u) + \sqrt{1 - \rho^2} \widehat{cov}(z, \eta)$$

This spurious instrument strength drives down the 2SLS standard error.<sup>36</sup>

Thus, a positive sample realization of  $\rho \widehat{cov}(z, u)$  generates an estimate shifted towards OLS. It also generates a low standard error because a large  $\rho \widehat{cov}(z, u)$  leads to a large  $\widehat{cov}(z, x)$ . Hence, 2SLS will appear spuriously precise in samples where the estimated coefficient is most shifted in the direction of the OLS bias. Conversely, estimates shifted away from OLS appear spuriously imprecise. This generates the power asymmetry in the 2SLS  $t$ -test.

### A2. A Brief Explanation of Weak Instrument Tests

The widely used Stock and Yogo (2005) weak IV tests assess whether instruments are strong enough for  $t$ -test size distortions to be modest. These tests are derived using the  $t$ -test power function in equation (A3). The size of a test is defined as the probability of rejecting  $H_0: \beta = 0$  when the null hypothesis is true. Thus, the size of the  $t$ -test is simply the power function evaluated at  $\beta=0$ .

<sup>35</sup>It follows that 2SLS is approximately median unbiased provided the instrument is strong enough that an incorrect first-stage sign a rare event, as only such events impart median bias.

<sup>36</sup>Note: since  $Var(\hat{\beta}_{2SLS}) = Var(\hat{\beta}_{OLS})/R_{z,x}^2$ , where  $R_{z,x}^2$  is first-stage  $R^2$ , the larger is  $\widehat{cov}(z, x)$  the smaller is the standard error.

A problem arises because the power function also depends on the degree of endogeneity  $\rho$ . To sidestep this issue, Stock and Yogo (2005) focus on the maximum (on the high side) size that arises when endogeneity is very severe,  $\rho = \pm 1$ . For example, by setting  $\beta=0$  and  $\rho = \pm 1$  in the power expression in equation (A3), and doing a grid search for the level of  $F = \lambda^2$  that sets power approximately equal to 15%, they obtain  $F=1.82$ .<sup>37</sup> Thus, if the first-stage  $F$  is 1.82 a 5% level two-tailed  $t$ -test will reject a true null hypothesis at a 15% rate, giving a size distortion of 10%. This is the maximal or “worst-case” size distortion.

A second problem arises as in any given sample we cannot observe population  $F = N \cdot \text{Var}(z\pi)/\sigma_e^2 = N\pi^2\sigma_z^2/\sigma_e^2$ . Rather, we can only observe the sample realization  $\hat{F} = N\hat{\pi}^2\hat{\sigma}_z^2/\hat{\sigma}_e^2$ . The sample  $\hat{F}$  is a draw from the non-central  $F$ -distribution with non-centrality parameter  $F$ . For instance, in the single instrument case, a sample  $\hat{F} > 8.96$  gives 95% confidence that  $F$  is at least 1.82.

Stock and Yogo (2005) weak IV tests give sample  $\hat{F}$  thresholds that give 95% confidence that population  $F$  is above some threshold, where that  $F$  in turn implies a certain maximal size distortion. For example, a sample  $\hat{F} > 8.96$  gives 95% confidence that  $F$  is at least 1.82, which in turn, implies the maximum size of a two-tailed 5%  $t$ -test is 15% (i.e., maximal size distortion of 10%).<sup>38</sup>

Both Keane and Neal (2022) and Angrist and Kolesár (2021) have criticized the focus on the worst-case scenario of  $\rho=\pm 1$ , arguing that for most plausible levels of endogeneity the size distortion is much less. But in addition, we also criticize the exclusive focus on size to the neglect of broader power considerations.

### A3. Power of the AR vs $t$ -test at Different Levels of Instrument Strength

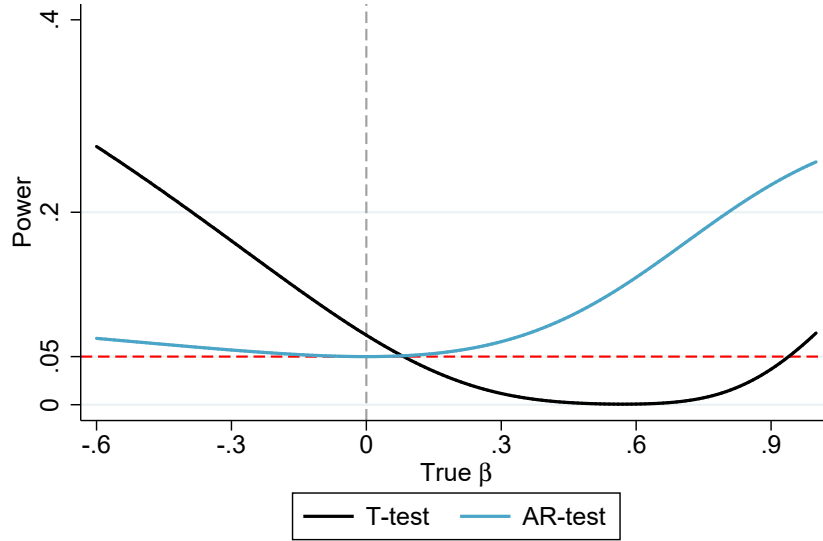
A critical point is that weak IV tests do not reveal if the estimator has acceptable power properties. We now present some comparisons of the power properties of the AR and  $t$ -tests at different levels of instrument strength  $F$ .

Figure A1 presents the case of  $F=1$ , which is indicative of the poor instrument strength in classic studies of the Frisch elasticity. We see the 2SLS  $t$ -test has essentially no power to detect positive Frisch elasticities in the plausible range of 0.1 to 0.9 (power is less than the 5% size of the test throughout this range). The AR test performs better: It does have power greater than size at all levels of  $\beta$  except  $\beta = 0$ , reflecting that it is an unbiased test. But its power is still very low: It doesn't pass 20% until the elasticity exceeds 0.8. This reflects the fact that the data is simply not very informative at this low level of instrument strength.

Another notable feature of Figure A1 is that the  $t$ -test appears to have much better power than the AR test for negative values of true  $\beta$ . This reveals the flip side of the power asymmetry problem: In samples where the 2SLS estimate

<sup>37</sup>Thus, the figure of 1.82 is subject to numerical error, and should not be viewed as exact. The same is true of other weak IV thresholds in this literature.

<sup>38</sup>Similarly, a sample  $\hat{F}$  of 10 gives 95% confidence that population  $F$  is at least 2.3, which implies a maximal size distortion of 8.5%. And a sample  $\hat{F} > 16.4$  gives 95% confidence that  $F$  is at least 5.78, which in turn implies the size distortion is no more than 5%.

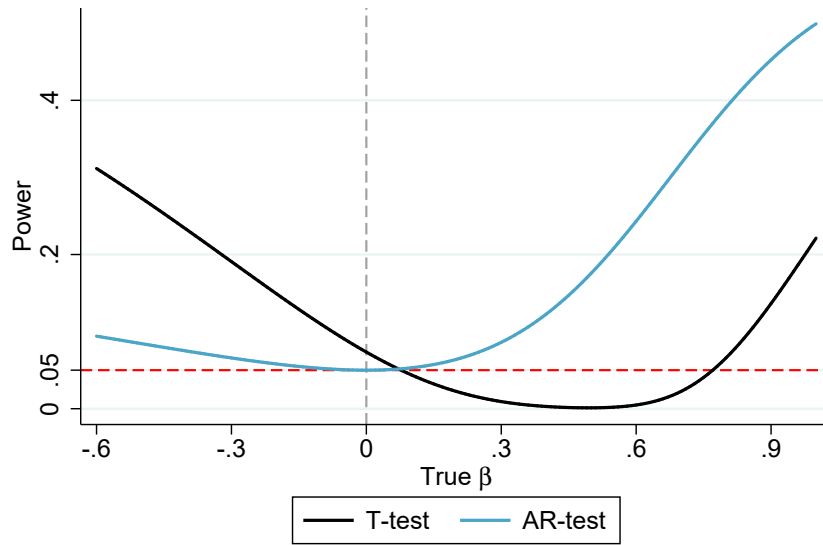
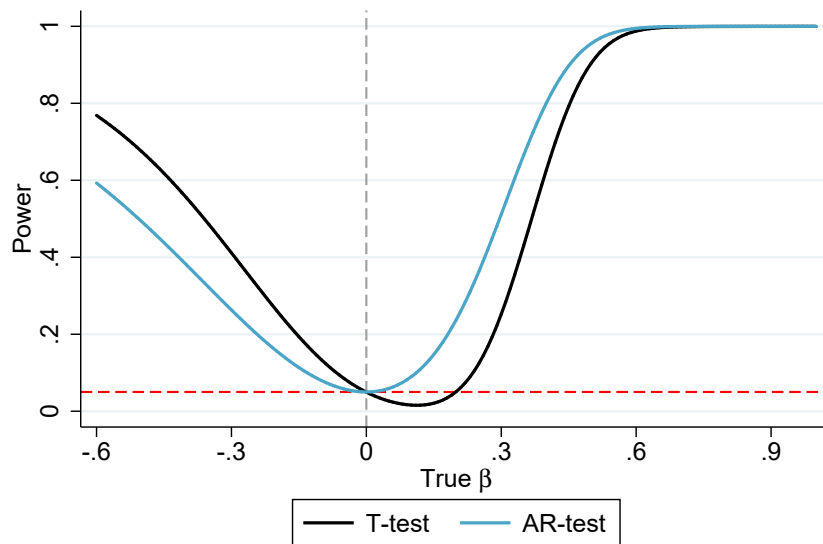
FIGURE A1. POWER OF THE T-TEST VS. AR-TEST WHEN POP.  $F = 1$  ( $\rho = -0.7$ )

is shifted in the direction of OLS, which in this case means it is shifted in the negative direction, the 2SLS standard error is spuriously small, which inflates the power of the  $t$ -test. This is not a desirable property, as the standard error exaggerates the precision of the estimate in such cases.

Figure A2 considers the case of  $F=2.3$ . This is particularly interesting, as a sample  $\hat{F}$  of 10 is required to give 95% confidence that  $F$  is at least 2.3. Thus, this case corresponds to the widely used Staiger-Stock rule of thumb, that a first-stage  $\hat{F}$  of at least 10 indicates an acceptable level of instrument strength. However, Figure A2 reveals that the  $t$ -test has very poor properties in this case. It has essentially no power to detect positive Frisch elasticities in the plausible range of 0.1 to 0.8, as power is less than the 5% size of the test throughout this range. The AR test has much better power to detect a true positive Frisch elasticity, but its power is still uninspiring (e.g., it doesn't pass 20% until the elasticity exceeds 0.5). So the data is not very informative at this level of instrument strength.

The severe bias of the  $t$ -test is also evident in Figure A2. Power is minimized in the vicinity of  $\beta=0.50$  rather than at  $\beta=0$ . This again reflects the power asymmetry of the  $t$ -test, and the fact that it has little power to detect a wide range of plausible positive elasticities because the OLS bias is negative.

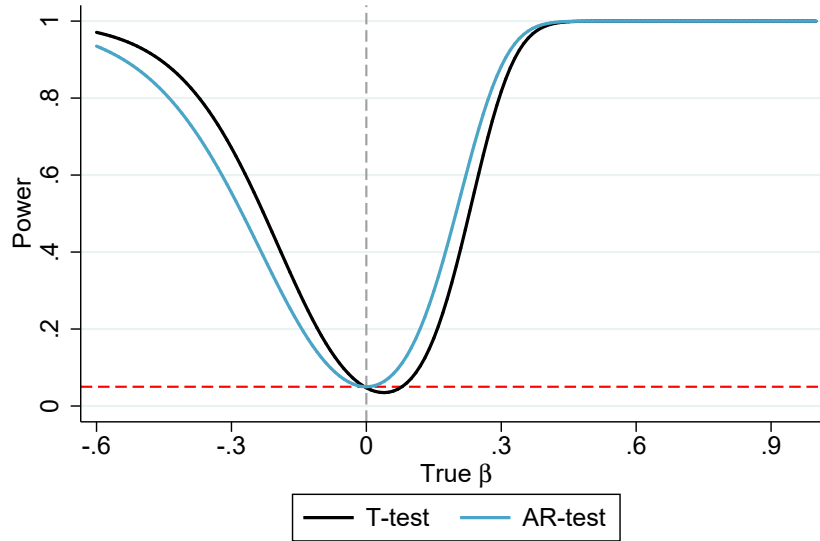
Figure 3 in the main text considers the case of  $F=10.12$ , which is the value we used in the Monte Carlo exercise in Section V. A first-stage  $\hat{F}$  of at least 23.2 is required to have 95% that population  $F$  is at least 10.12. In this case the  $t$ -test has power less than size for true effects in the 0.05 to 0.50 range, so, as we note in the text, it is uninformative over that range. The power of the AR test is far superior, reaching about 60% when true elasticity is 0.50.

FIGURE A2. POWER OF THE T-TEST VS. AR-TEST WHEN POP.  $F = 2.3$  ( $\rho = -0.7$ )FIGURE A3. POWER OF THE T-TEST VS. AR-TEST WHEN POP.  $F = 29.44$  ( $\rho = -0.7$ )

An obvious question is how large  $F$  must be for the  $t$ -test to begin to exhibit acceptable power for plausible elasticity values. Figure A3 reports results for a population  $F$  of 29.44. A first-stage  $\hat{F}$  of at least 50 is required to have 95% confidence that population  $F$  is at least this large. At this level of instrument strength

the power of both tests approaches one when the true elasticity approaches 0.6. However, the  $t$ -test still has power less than size for elasticities in the 0.0 to 0.2 range, and very poor power compared to the  $t$ -test for elasticities in the 0.0 to 0.4 range. The size of the  $t$ -test (i.e., power at  $\beta = 0$ ) is 4.96% in this case, so it is very close to the correct 5%. But bias is still evident as power is minimized at an elasticity of roughly  $\beta = 0.1$ .

FIGURE A4. POWER OF THE T-TEST VS. AR-TEST WHEN POP.  $F = 73.75$  ( $\rho = -0.7$ )



Finally, Figure A4 considers the case of true  $F=73.75$ , which is a very high level of instrument strength. A first-stage  $\hat{F}$  of at least 104.7 is required to have 95% confidence that population  $F$  is at least this large. We examine this case because Lee et al. (2022) show that a first-stage sample  $\hat{F}$  of at least 104.7 is required for the worst-case size distortion in the  $t$ -test to be no more than 5%.<sup>39</sup>

At this high level of instrument strength the power curves of the two tests are much more similar, and power of both tests approaches 1 for  $\beta$  around 0.40. But the power advantage of the AR test is still evident in the  $\beta \in (0.0, 0.30)$  range. For instance, for  $\beta=0.15$  the AR test power is 40% vs. 25% for the  $t$ -test.

In summary, our results clearly show that the power advantage of the AR test over the  $t$ -test is substantial at empirically relevant elasticity values. The power asymmetry of the  $t$ -test (i.e., its low power to detect plausible positive elasticities) is dramatic when instruments are weak but persists even when instruments are very strong.

<sup>39</sup>Their analysis is subtly different from Stock and Yogo (2005), in that their “worst case” refers to the maximum size distortion over all possible values of endogeneity  $\rho$  and all possible values of the true  $F$ . The worst case scenario for  $\rho$  is again  $\pm 1$ , while the worst case for  $F$  is  $[\hat{F}/(\sqrt{\hat{F}} + 1.96)]^2$ .