

FINANCIAL FRICTIONS AND THE NON-DISTORTIONARY EFFECTS OF DELAYED TAXATION

Andreas Fagereng and Marius A. K. Ring*

Preliminary

February 23, 2021

Abstract

Financially constrained agents discount future cash flows at above-market rates. In this paper, we present the hypothesis that delaying tax payments can materially reduce distortions when agents are credit constrained. We test this hypothesis in the context of the labor supply decisions of young workers in Norway, where a kinked income-contingent student-debt conversion scheme replicates an income tax with delayed payments. Bunching analyses reveal elasticities that are an order of magnitude below those we find at a regular income tax threshold, and which are increasing in ex-ante financial resources. These findings underline the potential for delayed taxation to be a powerful new component of optimal tax policy.

JEL: H21, G51, D15

Keywords: deferred taxation, delayed taxation, credit constraints, income taxation

*Fagereng is at BI Norwegian Business School. Ring (Corresponding author) is at the University of Texas at Austin. E-mail: mariuskallebergring[at]gmail.com. See www.mariusring.com for the most recent version. This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 851891).

1 Introduction

The optimal design of any tax depends crucially on its distortionary effects. How responsive labor supply is to the net-of-tax wage has seen tremendous empirical attention. While there are offsetting substitution and income effects at play in the response to labor income taxes, it is generally considered that the substitution effect dominates: A reduction in take-home wages lowers labor supply. However, little attention has been given to the effects of introducing a substantial time delay between tax accrual and tax payment.¹ The purpose of this paper is to propose and test the hypothesis that delaying the payment of labor income taxes may reduce their distortionary effects.

The intuition for this hypothesis comes from basic finance theory. Financial frictions render agents with steep earnings profiles unable to borrow against higher future incomes at the market rate. This creates a wedge between the market rate and their personal discount rate (Carvalho et al. 2016, Epper 2017). Therefore, constrained agents may be indifferent between a high income tax that is payable in the future and a lower income tax that is payable today. Assuming that tax authorities may borrow at the market rate, financial frictions therefore introduce a wedge between the tax authorities and constrained agents' net present value (NPV) of a given tax liability. If the net-wage elasticity of labor supply is positive, this wedge may allow tax authorities to increase tax revenues by allowing delayed payments, as constrained agents choose to supply more labor.

Delayed taxation amounts to government-provided income tax financing. Even if providing financing at the market rate, the government effectively earns an interest premium through higher tax revenues. Whether such a system of delayed taxation is optimal from a partial-equilibrium revenue-maximization perspective thus depends on whether the effective interest premium exceeds potential administrative and default costs. A potential upside occurs when consumption-smoothing agents increase their future labor supply as well in response to higher debt, which is the case in our model.

This paper focuses on empirically testing the core partial-equilibrium mechanism: In the presence of financial frictions, labor income taxes whose payments can be delayed become substantially less distortionary. Performing such a test is challenging, since there is little variation in the timing of tax payments. Taxes are typically paid either immediately (through withholding) or a year later when tax returns are due. We overcome this challenge by studying the

¹The concept that allowing tax payments (and not accrual) to be delayed may reduce the distortionary effects of labor income taxation in the presence of financial frictions is new, but the notion of delaying tax payments is not. For example, capital gains are typically only taxed at realization. Similarly, taxing entrepreneurial dividends rather than profits allow entrepreneurs to delay their tax payments (see Dávila and Hébert 2019).

effects of a student debt forgiveness scheme in Norway. This scheme creates a large jump in the effective marginal income tax rate where marginally accrued taxes can be financed with the same generous terms as subsidized student loans.

More specifically, Norwegian students are eligible to receive a yearly loan of around \$10,000, of which roughly half may be forgiven. However, if the student has labor earnings above approximately \$15,000, each additional dollar of earnings reduces the amount forgiven by 50 cents. This produces a hike in the effective marginal income tax rate of 50 percentage points, where the marginal tax may be financed at the same attractive terms as subsidized student loans.

This setting is well suited to investigate how financial frictions may render delayed taxation less distortionary. First, students are, almost by definition, highly constrained. Just a few years later they face significantly higher incomes against which it is hard to borrow. The dramatic increase in the effective tax rate at the earnings threshold is also more than significant enough for any student to be cognizant of it: At the threshold, the marginal net-of-tax (and debt increase) wage drops from 75 to 25 cents.² Despite this drastic reduction in the marginal (effective) wage, students are astoundingly irresponsive. While there is clear visual evidence that students do respond, these responses pale in comparison to the effective after-tax wage reduction that occurs. Our bunching analysis offers an implied labor earnings elasticity to the after-tax wage of only 0.016. While this estimate is highly statistically significant, it is an order of magnitude below most existing estimates (Keane, 2011). Labor market frictions are unlikely to explain our relatively low elasticity. Our sample is limited to students near the debt-conversion threshold, which is substantially below a full-time salary in Norway. This ensures that students are part-time workers who likely face flexible work arrangements

To shed light on the observed non-bunching behavior, we examine how student characteristics covary with their position relative to the debt-conversion threshold. These analyses suggest that non-bunchers (and their parents) have significantly less liquid wealth, but not lower future earnings. This is precisely what we would expect to see if irresponsiveness to the threshold is driven by financially constrained agents. We further find no evidence that the educational attainment of students' parents change in a manner consistent with these characteristics driving the differences in bunching behavior. Informed by these analyses, we study heterogeneity in bunching by the ex-ante financial situation of students and their parents. Students who have below median liquidity (and their parents as well) exhibit an implied labor earnings elasticity less than half as large as those above the median.

²The marginal tax rate around the threshold was approximately 25% during the sample period. This marginal tax applies to all marginal earnings regardless of the increase in student debt.

We continue to examine student bunching behavior at a regular tax threshold. This allows us to compare the implied labor supply elasticities under different payment schemes but among a similar sample of individuals.³ The tax threshold analyzed occurs at around \$6,000, where the marginal income tax rate goes from 0 to 25 percent. Using the same techniques as before, we estimate an implied labor supply elasticity of 0.13. This is more than ten times higher than the elasticity inferred from the delayed-tax threshold, which is consistent with materially reduced distortions when payments are substantially delayed.

We further present a simple model that relates differential responses to regular and delayed taxation to the marginal discount rate of a life-cycle agent. Under the assumption of a homogenous structural labor supply elasticity, an annualized marginal discount rate of 23% can explain the relatively muted responses to delayed taxation. The significantly higher responsiveness to delayed taxation for students with above-median liquidity indicates that these students optimize with a 10 percentage point lower discount rate.

The central contribution of this paper is to propose and test the hypothesis that delaying the payments of income taxes may substantially reduce its distortionary effects in the presence of financial frictions. To our knowledge, there exists neither theoretical nor empirical research on this topic. Our empirical setting fits the bill for testing due to three important features. (i) It effectively replicates a delayed income taxation system with a sizable hike in the marginal tax rate, where marginal taxes at this threshold are subject substantial delay in the payment. In addition, the sample consists of tax-payers where (ii) labor supply is highly flexible and (iii) borrowing constraints play an important role.

Related literature. On the conceptual front, this paper contributes to the literature on dynamic optimal taxation (see, e.g., [Ndiaye 2020](#) and the surveys in [Goloso and Tsyvinski 2015](#) and [Stantcheva 2020](#)). Most closely related is research that considers altering the timing of taxpayments or incorporating financial frictions.⁴ The conceptual novelty of this paper lies in this intersection.

³Ideally, this will keep unobservable factors causing frictions in labor supply optimization constant. An alternative would be to compare our elasticity under delayed taxation with elasticities from other research. However, this raises the concern that differences in labor market frictions are driving the differences in the elasticities.

⁴[Lockwood \(2020\)](#) theoretically examines how hyperbolic discounting affects the optimal timing of tax payments. [Andreoni \(1992\)](#) studies how financial frictions may affect tax policy, but the focus is on enforcement rather than timing. [Lozachmeur \(2006\)](#) studies optimal age-specific income taxation and finds that benefits from alleviating financial frictions lower the optimal tax rate for young (and more constrained) agents, but the analyses do not consider the potential optimality of delaying the payment of the tax (rather than lowering the rate itself) to achieve this benefit. Studying corporate taxation, [Dávila and Hébert \(2019\)](#) find that taxing payouts rather than profits is optimal in the presence of financial frictions. This essentially allows constrained firms with productive investment opportunities to delay when they pay taxes on their profits.

On the empirical front, this paper contributes to the growing literature studying bunching at tax thresholds (see, e.g., [Saez 2010](#), [Bastani and Waldenström 2020](#), [Søgaard 2019](#), [Seim 2017](#), and the review by [Kleven 2016](#)). Our contribution is to study bunching at a threshold where the *payment* of marginally accrued taxes is substantially delayed. This adds an intertemporal dimension to bunching behavior not present in other studies. We further add to the literature using income-contingent transfer schemes to identify labor supply elasticities (see, e.g., [Ong 2020](#) who exploits the income contingency of child support.) Finally, this paper also relates to the emerging literature on the effects of debt on labor supply (see, e.g., [Zator 2019](#), [Brown and Matsa 2020](#), [Donaldson et al. 2019](#), [Bernstein 2016](#)).

This paper proceeds as follows. Section 2 describes the empirical setting. Section 3 presents the empirical analysis. Section 4 introduces a simple two-period life-cycle model with endogenous labor supply and financial frictions that formalizes some of the intuition introduced in this paper. Section 5 briefly discusses aspects related to the implementation and potential trade-offs associated with introducing delayed taxation.

2 Empirical Setting

The main years of study are 2004–2011. During these years, most Norwegian students faced an earnings threshold ranging from NOK 104,500 in 2004 to NOK 140,823 in 2011. The monthly transfers ranged from NOK 8000 in 2004 to NOK 9785 in 2011. These transfers are initially given as a loan, but 40% may be forgiven (converted to a stipend) to the extent that students pass classes and stay below the earnings threshold above. Students are notified of the amount of transfers in the beginning of the academic year. These notification letters contain a breakdown of the transfers, noting the amount (40% of the total) that is given as a conversion loan, and stating that conversion from loan to stipend is contingent on incomes being below an income limit. The following year, students are notified how much of their loan was converted based on grades reported by educational institutions and income reported by the tax authorities. Loans must typically be paid off within 20 years following graduation. No interest is charged while still receiving support and loan payments may be delayed at the student’s discretion for up to 3 years in total.

This study is facilitated by administrative data hosted by Statistics Norway. The key data is derived from tax returns, including data on individuals’ incomes, assets and debts. The sample consists of students receiving standard student support for full-time studies for at least one full fiscal year during 2004–2011. We limit the sample to students who after conversion received a strictly positive stipend. This eliminates students who are ineligible for any debt-

conversion due to, e.g., living at home with parents. This ensures that close to all students in our sample are subject to income-contingent debt-conversion.

Summary statistics are provided in Table 1. The average student is 23 years old. This is reasonable in light of high school graduation occurring at age 18 and that we condition on students being enrolled for higher education for both semesters within a given year. The summary statistics reveal a substantial spread in the amount of liquid assets available to students. While students at the 25th percentile only hold NOK 8,000 (\$1,300) in liquid assets, students at the 75th percentile hold almost *ten* times more. A similar spread can be observed in the liquid assets of the students' parents. We further see that the average student earns around NOK 100,000 (\$17,000), which is a direct consequence of our sample restrictions aimed at students around the debt-conversion threshold. Four years later, the average student faces considerably higher earnings at around NOK 360,000 (\$60,000).

TABLE 1: SUMMARY STATISTICS

This table provides summary statistics. The main sample period is 2004–2011. Financial variables are denominated in NOK. The USD/NOK exchange rate was around 6 in 2010. The main sample is restricted to students who had labor earnings within 50,000 of the debt-conversion threshold. Liquid Assets are made up of deposits, mutual funds, and ownership in public equity. Labor earnings are censored to be below NOK 1,000,000 in 2010 NOKs. The Bottom Tax Threshold is only considered for the years 2005–2011.

	N	Mean	p25	p50	p75
Liquid Assets $_{t-1}$	230,906	57,522	7,989	29,296	77,099
Liquid Assets $_{t-1}$ (Parents)	214,419	429,326	59,805	176,471	460,545
Age	231,036	23.4	22	23	25
Labor Earnings $_t$	231,036	101,394	81,156	98,536	118,966
Labor Earnings $_{t+4}$	229,027	357,506	226,244	372,615	464,829
Debt-Conversion Threshold $_{tt}$	231,036	120,162	108,680	116,983	128,360
Bottom Tax Threshold $_t$	198,815	36,706	29,600	39,900	39,900

3 Empirical Analysis

3.1 Bunching at the debt-conversion threshold

In this section, we use the nonparametric methodology developed by [Chetty et al. \(2011\)](#) to compute the amount of bunching at the tax-like debt-conversion threshold. This provides a bunching elasticity, b , which is the relative excess mass of students near the threshold. We translate this into a more informative bunching estimate, B , which is the estimated relative

reduction in earnings caused by the presence of the threshold.⁵ This estimate can be used to infer the implied compensated labor supply elasticity via the following formula introduced by Saez (2010):

$$e = \frac{B}{d\tilde{\tau}/(1 - \tilde{\tau})}, \quad (1)$$

where $d\tilde{\tau}$ is the change in the marginal income tax rate at the threshold.

In the presence of delayed taxation, the tax rate $\tilde{\tau}$ —according to which the agent optimizes—differs from the nominal tax rate, τ . This means that the standard result that e equals the Frisch elasticity of labor supply as in Saez (2010) does not necessarily hold. In Section 4, we outline a simple two-period model, in which the agent faces a labor income tax where only a fraction δ is payable in the current period. In this model, the agent behaves as if facing a standard (payable-today) labor income tax of

$$\tilde{\tau} = \delta\tau + \frac{1 - \delta}{\tilde{R}}\tau, \quad (2)$$

where \tilde{R} is the relevant marginal (gross) interest rate faced by the agent.

If taxes are payable today ($\delta = 1$), e may provide an estimate of the Frisch elasticity of labor supply. However, in the current setting, $\delta = 0$, as non-converted student loans are paid in the future. Thus, in order to relate our empirical estimate of e to the Frisch elasticity of labor supply, we need to know the applicable marginal gross discount rate of bunchers, \tilde{R} .

Since this is unobservable, we proceed as if $\delta = 1$. Then, in Section 3.4, we estimate the implied elasticity from bunching at a regular tax (i.e., $\delta = 0$ and thus $d\tilde{\tau} = d\tau$) and compare the two elasticities to find the \tilde{R} that would allow them to be consistent with the same *structural* labor supply elasticity.

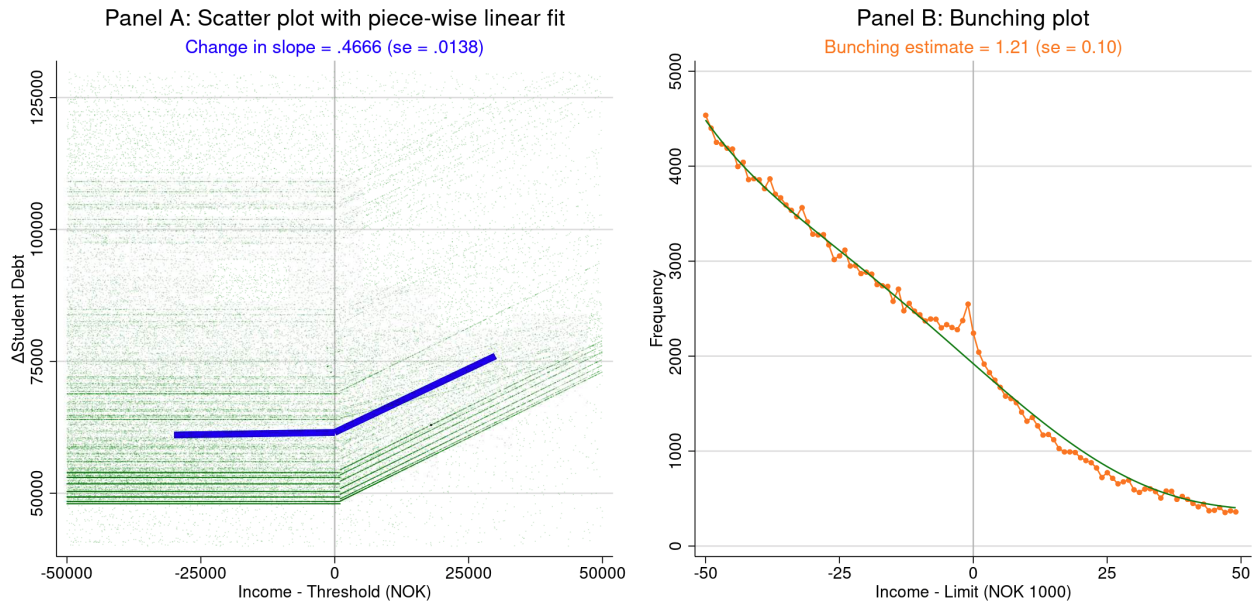
⁵ b is the relative excess mass near the threshold. This is the estimate that is typically reported in bunching plots. The standard approach is to assume that (i) the entire excess mass is located at the income bin directly to the left of the threshold and that (ii) absent the presence of the threshold, all earnings would be immediately to the right of the threshold. This means that we can multiply b by the width of the threshold to get a measure of how much earnings is shifted down due to the presence of the threshold. By further scaling this by the value of the threshold, we obtain a measure, B , of the relative reduction in earnings caused by the (effective) increase in the marginal tax rate.

Figure 1 shows some key details of the empirical analysis. Panel A verifies that earnings above the threshold lead to an increase in next period debt. Most students are on the expected kinked trajectory where each additional NOK of earnings increases debt by 0.50 NOK. The blue fitted line illustrates how we obtain our first-stage measure of the effect of excess earnings on debt accumulation. We find that the slope increases by 0.47. This is close to the nominal increase of 0.50 due to very few non-compliers.⁶

The second plot estimates the slope of changes in debt on earnings to be statistically indistinguishable from 0.5 above the threshold, consistent with the institutional setting. In other words, the marginal effective tax rate increases by 50 percentage points at the threshold. Cast in terms of the previous notation, this implies that $d\tau = 0.5$.

FIGURE 1: VERIFYING THE EFFECT OF EXCESS EARNINGS ON FUTURE DEBT AND EXAMINING BUNCHING RESPONSES

Panel (A) provides a scatter plot, in green, of the relationship between debt accumulation and student earnings around the debt-conversion threshold. The fitted blue line illustrates the estimation of the effect of earnings in excess of the threshold and accumulated debt. Panel (B) provides a graphical illustration of how the bunching estimate. The orange connected line shows the actual distribution of students around the conversion threshold. The fitted green line shows the estimated counterfactual distribution. The bunching estimate provides the relative excess mass (actual versus counterfactual) of students near the threshold. This is done using the Stata .ado file provided by [Chetty et al. \(2011\)](#). This program calculates the excess bunching between NOK -10,000 and NOK 6,000. Standard errors are computed from bootstrapping (N=1,000). All plots represent statistics from the pooled sample years 2004–2011.



In Panel (B), the yellow dotted line shows the distribution of students around the earnings threshold. The green line is the counter-factual distribution, which is a 5th-order polynomial

⁶Some non-compliers exist, for example, because they may have moved in with their parents during the fall semester, which would exclude them from receiving any conversion for fall semester loans. Such moves must be reported to the educational loan fund, but not to the tax authorities from which we receive address data.

fitted to the non-bunching region. We obtain a measure of the excess mass of individuals near the threshold by comparing the actual and counter-factual distributions. This offers a bunching estimate of 1.21, which means that there are 121% more individuals around the threshold than the counter-factual distribution implies. To translate this bunching estimate into a labor earnings elasticity, we multiply it by the size of the bins (1,000), and divide by the average threshold amount (120,162) to get the percent change in labor earnings caused by the presence of a threshold. This provides an estimate of B in equation 1. We then further divide by 47%, which itself is divided by the average marginal (net-of-tax) keep rate at this earnings level of 75%. Per equation 1, this produces an elasticity of labor earnings to the net-of-tax (or net-of-debt-increase) wage of 0.0162.⁷ Standard errors are similarly found to be 0.0013.⁸

3.2 Determinants of non-bunching

We now investigate potential determinants of this non-bunching behavior. Our main approach is to plot student characteristics against their position relative to the conversion threshold.⁹ This is a visual exercise where we attempt to draw conclusions from visual breaks in the relationship between a given characteristic and students' earnings occurring around the conversion threshold.

In Figure 2, Panel (A), we find that the amount of ex-ante liquid assets drops sharply right above the threshold. This suggests that non-bunchers have less liquid wealth, consistent with these students being financially constrained. Panel (B) of Figure 2 shows how future labor earnings vary with the student's position relative to the threshold. This reveals no sharp rise or decrease in realized future incomes above the threshold, which suggests that non-bunchers do not differ significantly in terms of near-future earnings prospects.

Taken together, these findings emphasize financial frictions as a key channel in driving the insensitivity to the conversion threshold. Those who earn above the threshold have similar future earnings prospects, but have significantly less liquid assets. Holding less assets may both causally affect the extent to which the agents are constrained and be a proxy for financial frictions as it indicates a preference towards smoothing consumption toward the present.

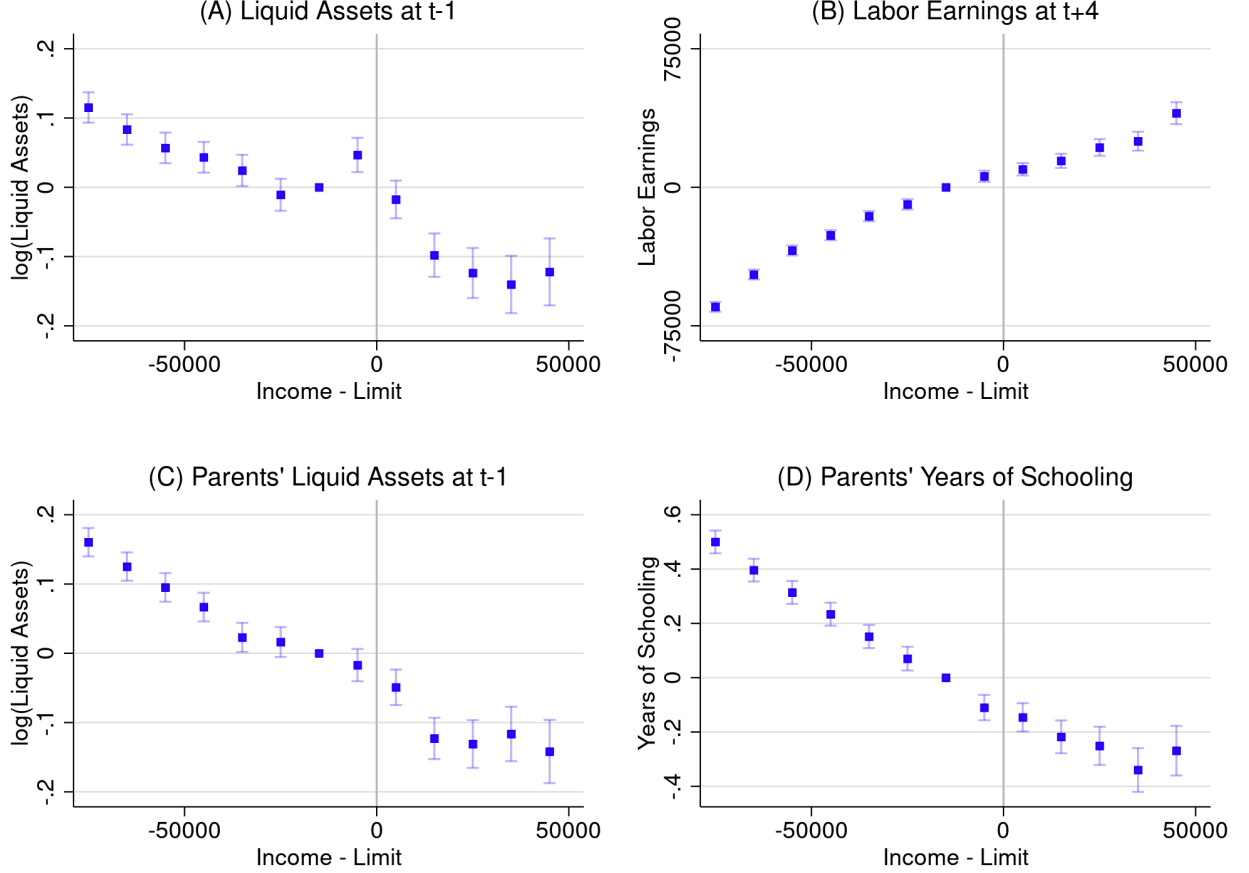
⁷These calculations do not adjust for the fact that any accumulated debt is interest-free while in school. Adjusting for a 3-year 3%-interest discount would increase the elasticity by around 9%.

⁸We ignore the (very small) standard errors involved with estimating the change in debt per additional NOK earned above the threshold.

⁹Another application of this type of analysis can be found in concurrent work by [Bastani and Waldenström \(2020\)](#) who examine how ability covaries with taxpayers' position relative to a regular tax threshold to infer the ability gradient in tax responsiveness.

FIGURE 2: CHARACTERISTICS OF STUDENTS BELOW AND ABOVE THE
INCOME-CONTINGENT DEBT-CONVERSION THRESHOLD

The graphs below show the financial characteristics of students who are near the threshold. Panel A considers the liquid assets of students. These consist of deposits, stocks, bonds, and mutual fund holdings. Panel B shows future log labor earnings, measured 4 years later. Panel C shows the amount of liquid assets held by the student's parents. Panel D shows the educational attainment of the parents, measured as the maximum number of years of school among the set of parents. Standard errors used to provide 95% confidence intervals are clustered at the student level.



To investigate this liquidity channel further, we also show how *parents'* liquidity correlates with the students' earnings location in Panel (C) of Figure 2. This documents a noteworthy negative relationship between the parents' financial resources and the in-school labor earnings of the child. This suggests that parents play an important role in determining the amount of time students may dedicate to their studies. More relevant to the present study, is the finding that parents' assets drop shortly above the earnings threshold. This indicates that non-bunchers have access to fewer financial resources, which is consistent with financial frictions playing a key role in driving the observed non-responsiveness to the conversion threshold. However, wealth may proxy for human capital which influences tax responsiveness ([Bastani and Waldenström, 2020](#)). Therefore, we plot parental educational attainment on the y-axis in Panel (D). This shows that there is no break in the relationship between educational attainment, measured in

the maximum years of schooling among the parents and the child's position relative to the conversion threshold. This addresses the hypothesis that less resources, in a human capital, rather than financial, sense can explain the irresponsiveness to the threshold. If anything, extrapolating from the below-threshold relationship, non-bunchers may have higher-educated parents. To the extent that this is correlated with the students' life-time wealth, this may explain some of the desire of students to front-load consumption through the incurring higher student loans.

3.3 Bunching heterogeneity

We proceed by a supplementary, more standard approach of investigating heterogeneity in earnings sensitivity to the threshold in Figure 3. This approach splits the sample into subsets based on student and parental characteristic to compute heterogeneous bunching elasticities. We see that the largest contribution to the total excess mass in the preceeding Figure 1 is from students who themselves and their parents have above-median liquid assets. Figure 3 also suggests that the main driver of bunching responses is the parents' rather than the students' own liquid assets. Moving from the left to the right panels, which improves the parents' liquidity, more than doubles the bunching estimates.¹⁰

What can this heterogeneity tell us about how the severity of financial constraints vary with parents' liquidity? The source of variation in the implied labor supply elasticities across the liquidity subsamples are the bunching estimates provided in Figure 3. The ratio of bunching estimates therefore provide the relative implied elasticities. If we impose the same structural labor supply elasticity (e.g., same constant Frisch elasticity of labor supply) across the samples, differences can only be attributed to differences in the (gross) discount rate, \tilde{R} (See Proposition 1 in Section 4 for a theoretical example). This follows from substituting the expression for $\tilde{\tau}$ in equation 2 into the expression for the e in 1 and setting the fraction payable today, δ , to zero.

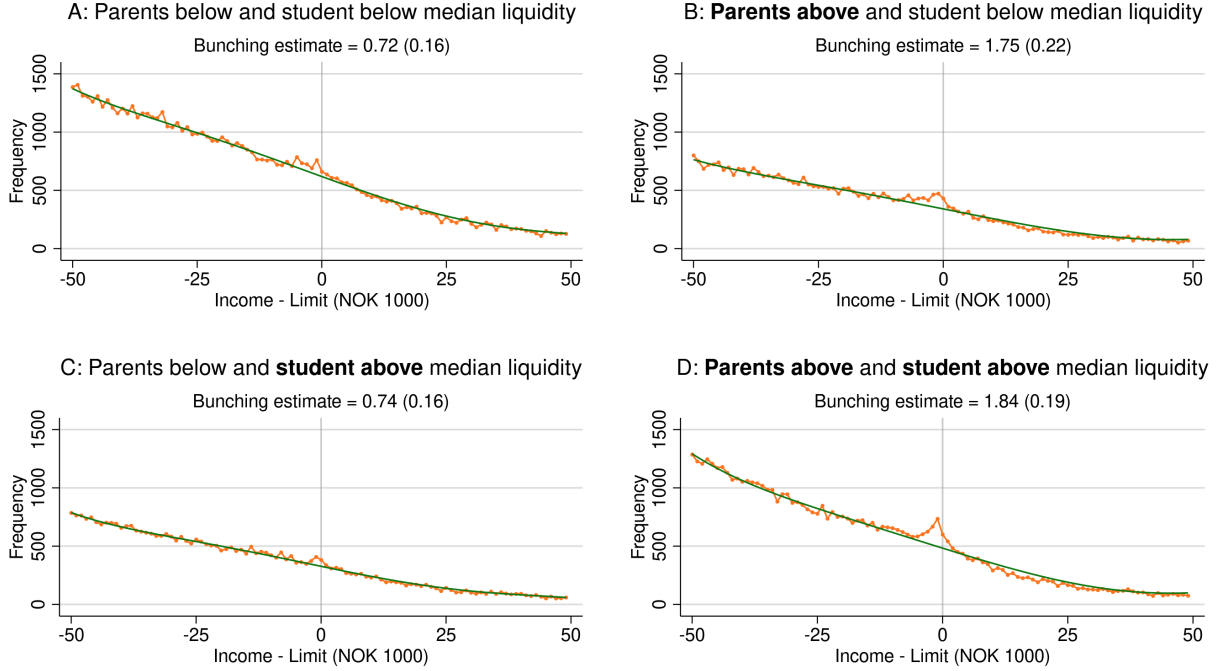
In Figure 3, we see that the elasticity increases from 0.72 to 1.84 when moving from below to above the median in terms of both students' and their parents' resources. This thus implies that the doubly-below median students have an average gross discount rate that is 2.56 times larger. Annualizing this, assuming a 10-year horizon, implies a $2.56^{1/10}=1.0986$ times higher annualized gross discount rate or approximately a 10 percentage point higher interest rate.¹¹

¹⁰In this case, it doesn't matter whether we compare excess mass in terms of students or earnings, since bin widths and thresholds are the same.

¹¹If the doubly-above median group has a baseline gross interest rate of 1.10, then the below-median group has a discount rate that is $(1.0986-1)*1.10 = 10.85$ pp. higher.

FIGURE 3: HETEROGENEITY IN BUNCHING BY AMOUNT OF LIQUID ASSETS

These plots calculate the bunching elasticity for different subsamples. Students are split into four subsamples based on whether their and their parents' $\text{LiquidAssets}_{t-1}$ are below or above median. These medians are calculated separately for each year in the sample.



3.4 Analysis of bunching at a regular tax threshold

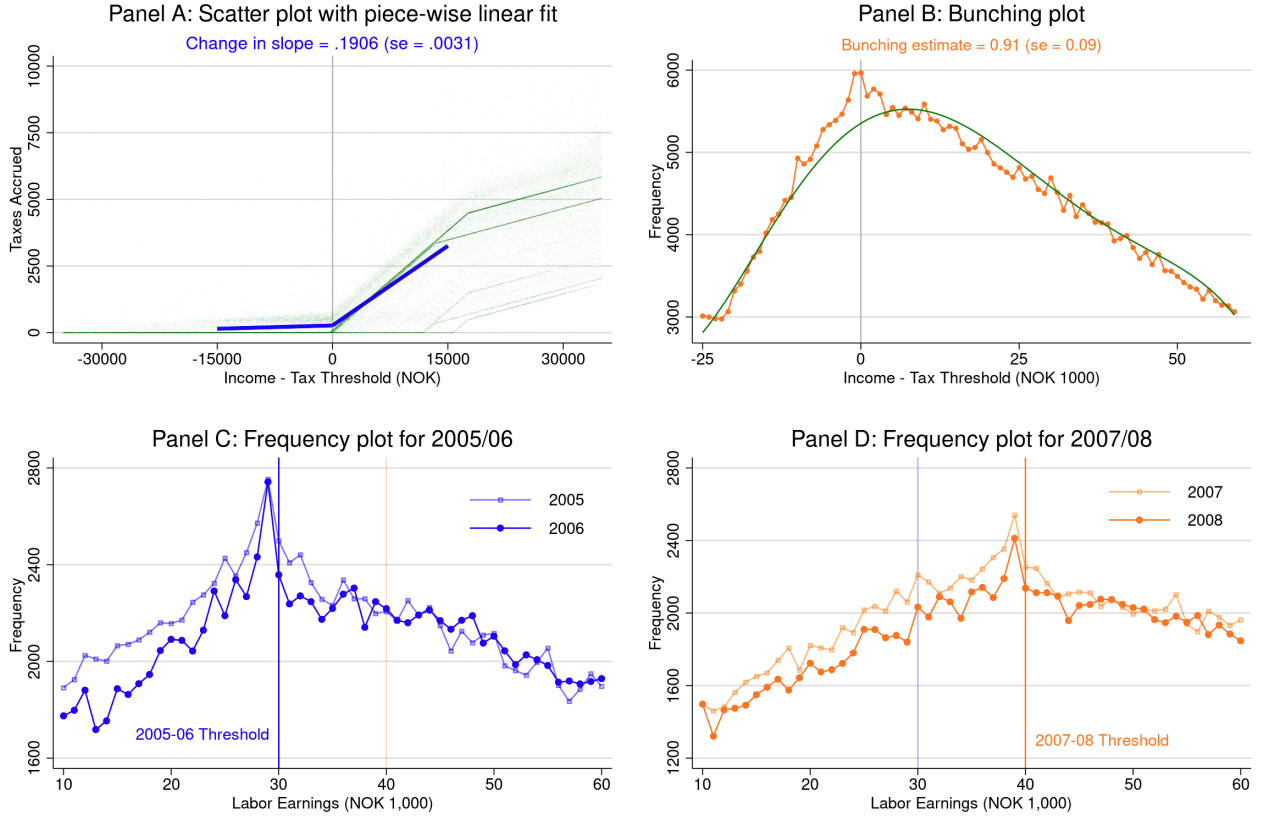
In this section, we repeat the introductory analyses done in Figure 1 using a *tax* threshold rather than the debt-conversion threshold. The purpose of this exercise is to obtain a reference estimate of the implied labor earnings elasticity at a tax threshold where marginally accrued taxes are not delayed. We focus on the first tax threshold in the progressive income taxation system. This threshold is located at NOK 30,000 during 2005–06 and NOK 40,000 during 2007–2011.¹² At this threshold, the marginal income tax increases from 0 to around 25 percent for most tax-payers.

In Figure 4, we investigate this complementary empirical setting. Panel (A) provides a scatter-plot which verifies the presence of a rise in the marginal income tax rate by plotting total taxes accrued that year against incomes. It also provides the fitted kinked line, from which we infer an average increase in the marginal tax rate of 19 percentage points at the threshold. The coefficient is lower than the nominal increase of 25 percentage points since some individuals may be eligible for higher standard deductions.

¹²We omit 2004. During this year the threshold was only NOK 23,000, which substantially reduces how much of the left tail we can use to estimate a counterfactual distribution.

FIGURE 4: BUNCHING AT A REGULAR TAX THRESHOLD

The first and second plots shows the relationship between labor incomes (“pensionable income”) and taxes accrued that year (payable same or next year) in the form of a scatter and binscatter plot, respectively. The third plot shows the distribution of students around the income tax threshold. The fourth plot calculates the bunching elasticity, in terms of the implied excess fraction of students located in the NOK 1,000 bin directly to the left of the threshold using the Stata .ado file provided by [Chetty et al. \(2011\)](#). This program calculates the excess bunching between NOK -10,000 and NOK 6,000. Standard errors are computed from bootstrapping (N=1,000). All plots represent statistics from the pooled sample years.



Panel (B) illustrates how the bunching estimate of $b = 0.091$ is calculated. While this bunching estimate is smaller than that found at the debt-conversion threshold, this comparison is uninformative for two reasons. First, the tax threshold is lower. This means that the implied excess mass, measured in relative income rather than relative amount of individuals, is higher. To obtain this number we multiply 0.091 with 1,000 (width of a bin) and divide by 36,706 (average threshold) to obtain an excess income mass of $B = 2.48\%$. This is higher than at the debt-conversion threshold. Second, we need to account for the fact that the reduction in the net-of-tax wage is considerably smaller at the tax threshold. Thus to obtain the implied elasticity, we divide by $19.6\%/100\%$ rather than the much higher $46.7\%/75\%$ at the debt-conversion threshold. This produces an implied elasticity of labor earnings to the net-of-tax wage of $e = 0.13$.

In Panel (B), we see that the bunching mass occurs at the mode of the distribution. If the location of the mode is not driven by students’ responses to the tax threshold, then the

co-location of the mode and threshold could lead to an upward bias in the bunching estimate. To address this concern, we show in Panels (C) and (D) that the location of the mode is driven by the location of the tax threshold. From 2005 to 2006 and from 2007 to 2008 there was no changes to the mode of the distribution. However, when the tax threshold rose from 2006 to 2007, the mode precisely followed. This reassures us that there is indeed substantial responsiveness to the tax threshold not driven by happenstance co-location of the mode and threshold.

The elasticity of 0.13 is 8 times larger than the elasticity of 0.0162 found in when analyzing responsiveness to the debt-conversion threshold. For these to be consistent with the same structural labor supply elasticity, the relevant average cumulative gross discount rate would have to be 8. Given that these loans have an average maturity of 10 years, this necessitates an annualized net discount rate of $23\% = 8^{1/10} - 1$. This number is comparable to average credit-card rates that lie slightly above 20%.¹³ This implies that students are willing to borrow from the educational loan fund at a rate exceeding that offered by financial institutions. This may be driven in part by credit rationing, but likely primarily from the fact that the loan fund does not require payments while students are still in school and generally have a long maturity with the additional opportunity to delay payments for up to three years.

We can use this implied elasticity to get an idea of how much bunching would be caused by the debt-conversion threshold in the absence of financial frictions. In other words, how much bunching would there be in Figure 1 if students responded to the debt-conversion threshold as if it were a regular income tax threshold? To find this, we reverse the calculation used to infer labor supply elasticities from bunching estimates. This offers a counter-factual bunching estimate of 23.43.¹⁴ This is considerably larger than the empirical bunching estimate of 1.21.

4 A simple model of labor supply under delayed taxation

The purpose of this section is to introduce a simple model that can guide the comparison of implied labor supply elasticities. The central take-away is Proposition 1, which emphasizes the role of discount rates in differential responses to regular and delayed taxation. However, the simple model allows us to go one step further to examine the effect of delaying tax payments on *future* labor supply as well, which results in Proposition 3.

¹³Source: Statistics Norway's Statistics on Interest Rates in Banks and Credit Institutions, source table 12844, 2019Q4: 21.6%

¹⁴ $= 0.13 * (120162 / 1000) * (75 / 50)$

Model environment. The agent works and consumes for two periods. A fraction, δ , of period 1 taxes are paid in the first period. The remainder, $1 - \delta$, is paid in the second. The period 2 wage, w_2 , is net-of-tax, and payable in period 2. The agent faces the following maximization problem.

$$\max_{c_1, c_2, l_1, l_2, s} u(c_1, l_1) + \beta u(c_2, l_2), \quad (3)$$

$$\text{s.t. } c_1 + s = y_1 + l_1 w_1 (1 - \tau \delta) \quad (4)$$

$$\text{and } c_2 = \bar{R}(s) + y_2 + l_2 w_2 - l_1 w_1 \tau (1 - \delta). \quad (5)$$

Where w_1 is the first-period gross wage, c_t is consumption, l_t is labor supply, l_t is exogenous income, s , and is savings. τ is the nominal tax rate for period 1 income.

The financial friction is the following. When agents save an amount s greater than $\bar{s} < 0$, they face a gross interest rate of 1, which is the same gross interest rate that the tax authorities charge on delaying tax payments. When they save less than \bar{s} (in general, borrow), they face a gross interest rate $R > 1$.

$$\bar{R}(s) = s + (R - 1)\mathbb{1}[s < \bar{s}](s - \bar{s}) \quad (6)$$

$$\equiv s\tilde{R} - (R - 1)\mathbb{1}[s < \bar{s}]\bar{s} \quad (7)$$

Assume that we have additively separable (dis)preferences for consumption and labor supply, and that the per-period utility takes the following form:

$$u(c, l) = \frac{1}{1 - \gamma} c^{1 - \gamma} - \psi \frac{l^{1 + \nu}}{1 + \nu}. \quad (8)$$

We will focus on the cases where the first order conditions (FOCs) bind. In other words, we consider cases where the optimally chosen s is not at the kink point, $s = \bar{s}$, in the agent's budget constraint. We can think of this as focusing on either the unconstrained agent, where $s > \bar{s}$, or the highly constrained agent who chooses $s < \bar{s}$ even if this entails borrowing at $\tilde{R} = R > 1$. This could mimic a setting in which the only source of loans available are high-interest credit cards. We conjecture that the responses of any agent who optimally chooses the kink point, $s = \bar{s}$, would be consistent with behavior “in the middle” of these two types of agents that we

analyze.¹⁵

FOCs for these cases where the optimal s is different from \bar{s} , are:

$$s : \quad c_1^{-\gamma} - \beta \tilde{R} c_2^{-\gamma} = 0 \quad (9)$$

$$l_1 : \quad w_1(1 - \tau\delta)c_1^{-\gamma} - \psi\nu l_1^\nu - w_1\tau(1 - \delta)\beta c_2^{-\gamma} = 0 \quad (10)$$

$$l_2 : \quad w_2 c_2^{-\gamma} - \psi\nu l_2^\nu = 0 \quad (11)$$

We see that the delayed tax scheme alters the standard optimization problem by introducing the third term in equation 10. Effectively, it adds an intertemporal component to the standard intratemporal trade-off between leisure and consumption.

To simplify the main proposition below, it is useful to define a somewhat stricter notion of not being at the kink point $s = \bar{s}$. This will allow us to not worry about agents not hitting the kink point if we make changes to the tax environment.

Definition of IHS: The intertemporal first-order condition (9) holds strongly (**IHS**) whenever the agent could increase saving by the delayed portion of period 1 income taxes without changing the marginal interest rate. More formally, this condition says that

$$s + l_1 w_1 \tau (1 - \delta) / R < \bar{s} \quad \text{or} \quad s > \bar{s}. \quad (12)$$

This condition is weak in the sense that it can always be satisfied by considering a small enough fraction, $1 - \delta$, of period 1 taxes that are delayed.

Proposition 1. It does not matter whether the agent faces a tax where a fraction, $1 - \delta$, is paid in period 2 or whether the agent faces a tax rate in period 1 of $\tilde{\tau} = \tau(\delta + (1 - \delta)/\tilde{R})$, i.e., a discounted tax, whenever the IHS condition holds.

Proof: First, note that δ only enters in to the period 1 intratemporal first-order condition (10). Use equation 9 to substitute for $\beta c_2^{-\gamma}$ in equation 10. Then we see below that a tax of $\tilde{\tau}$, payable in period 1, is equivalent, in terms of the FOCs, to the case where δ of the tax is paid in period 2.

¹⁵Think of the agents at the kink points as “moderately constrained” agents who would choose to borrow if the interest rate they faced were at least slightly lower than R .

$$w_1 \left(1 - \tau \{ \delta + (1 - \delta) / \tilde{R} \} \right) c_1^{-\gamma} - \psi \nu l_1^\nu = 0 \quad (13)$$

We must also increase savings, s , by $\tilde{s} = l_1 w_1 \tau (1 - \delta) / \tilde{R}$ in the budget constraints, having the agent effectively pre-pay the tax by saving (or borrowing less to pay) for it. We need \tilde{R} to remain constant for this to work, which is satisfied by the IHS. In other words, we can replace the delayed tax by a period 1 tax with the same NPV from the perspective of the agent, while leaving first-order conditions unaffected and still satisfy the budget constraints. This holds true for any δ , but the IHS is obviously easier to satisfy if δ is close to 1.

Corollary 1. If the elasticity of labor supply to the net of tax wage is positive (negative) and $\tau > 0$, then marginally decreasing δ – allowing more to be paid later – will strictly increase (decrease) labor supply for constrained agents ($s < \bar{s}$) whenever the IHS holds.

Proof: From Proposition 1, we know that the agent optimizes as if facing a regular income tax, payable during period 1, of $\tilde{\tau} = \tau(\delta + (1 - \delta)/\tilde{R})$. A marginal decrease in δ will thus strictly decrease the effective net-of-tax wage, since the $s < \bar{s}$ implies $\tilde{R} > 1$, and thereby increase (decrease) labor supply whenever the elasticity of labor supply to the net-of-tax wage is positive (negative).

Proposition 2: Intertemporal labor supply. Unconstrained agents who face $\tilde{R} = 1$ see no effect on labor supply growth from changing the fraction of period 1 taxes paid in period 1, δ . Constrained agents with binding FOCs, with $\tilde{R} = R > 1$, see an increase in labor supply growth if δ is reduced.

Proof: First insert the FOC for period 2 labor supply into the FOC for period 1 labor supply and divide through by w_1/w_2 .

$$\frac{w_2}{w_1} \left(\frac{l_2}{l_1} \right)^\nu = (1 - \tau \delta) \frac{c_1^{-\gamma}}{c_2^{-\gamma}} - \tau(1 - \delta)\beta \quad (14)$$

The LHS is strictly increasing in l_2/l_1 . Thus the sign of the derivative of the RHS with respect to δ provides the sign of the derivative $d(l_1/l_2)/d\delta$.

$$\frac{d}{d\delta} \frac{w_2}{w_1} \left(\frac{l_2}{l_1} \right)^\nu = -\tau \frac{c_1^{-\gamma}}{c_2^{-\gamma}} + (1 - \tau\delta) \frac{d}{d\delta} \frac{c_1^{-\gamma}}{c_2^{-\gamma}} + \tau\beta \quad (15)$$

$$= \tau \left(\beta - \frac{c_1^{-\gamma}}{c_2^{-\gamma}} \right) + (1 - \tau\delta) \frac{d}{d\delta} \frac{c_1^{-\gamma}}{c_2^{-\gamma}} \quad (16)$$

If the intertemporal FOC holds, then the RHS would be $\tau\beta(1 - \tilde{R}) \leq 0$. Thus a decrease in the fraction paid today would strictly increase l_2/l_1 if $s < \bar{s}$ and thus $\tilde{R} = R > 1$.

Proposition 3: Life-time labor supply. If the IHS holds and the elasticity of labor supply to the net of tax wage is positive, then highly constrained agents respond to a shift towards (more) delayed taxation by supplying more labor in both time periods. More specifically, if the IHS holds, then a reduction in δ increases both l_1 and l_2 .

Proof: From Corollary 1, we know that if the elasticity of labor supply to the net-of-tax wage is positive, then postponing the payment of period 1 taxes will increase labor supply in period 1 for constrained agents ($s < \bar{s}$) for whom the IHS holds. From Proposition 2, we learn that this would also increase labor supply growth between period 1 and period 2, which implies that period 2 labor supply increases as well.

5 Discussion

This paper introduces the hypothesis that delaying labor income tax payments may reduce their distortionary effects in the presence of financially constrained agents. We exploit a unique setting in Norway that allow us to test this hypothesis empirically. Our results indicate that delaying the payment of taxes, while keeping time of accrual constant, materially reduces the distortionary effects of income taxation when agents are credit constrained. These findings highlight delayed taxation as a promising new tool in optimal taxation and a fertile ground for more theoretical and empirical research.

The most feasible implementation of a system of delayed taxation is likely in connection with raising a given marginal income tax rate. Imagine an economy with a flat tax rate of 30%. Policy-makers are considering increasing the marginal rate to 40% at some threshold, say \$100,000. Allowing marginal taxes (10% on the amount above \$100,000) to be delayed is likely to reduce the distortionary effects among constrained agents. By construction, tax liabilities would only accrue to high-income individuals, which may limit the potential severity of issues such as adverse selection.

In a life-cycle model calibrated to U.S. workers, [Scott et al. \(2021\)](#) find that workers aged 25 would require a match rate above 1800% to participate in employer-sponsored retirement saving, which is driven by borrowing constraints and an upward-sloping earnings profile. This match rate decreases with age and approaches zero around age 45. This provides a useful statistic to explore the potential effects of delayed taxation: It indicates that the average 25 year old may be *eighteen times* less responsive to a labor income tax that can be paid at retirement and that incentives to delay accrued taxes would start vanishing at age 45. This suggests that a reasonable implementation of delayed taxation might involve age limits if policy-makers trade off tax-revenue effects and potential costs from mortality-induced non-payment.

There may be important costs associated with non-payment or debt-overhang (see, e.g., [Donaldson et al. 2019](#) and [Cespedes et al. 2020](#)) induced by such a scheme, particularly if, e.g., hyperbolic discounting plays an important role in determining the extent to which the agent is constrained. These costs should be weighed against the potential benefits from reduced distortions while bearing in mind that the effect on over-all indebtedness may be limited due to substitution away from other sources of credit.

References

- Andreoni, J. (1992). IRS as loan shark tax compliance with borrowing constraints. *Journal of Public Economics*, 49(1):35–46.
- Bastani, S. and Waldenström, D. (2020). The ability gradient in tax responsiveness.
- Bernstein, A. (2016). Household debt overhang and labor supply. *Unpublished Working Paper*.
- Brown, J. and Matsa, D. A. (2020). Locked in by leverage: Job search during the housing crisis. *Journal of Financial Economics*, 136(3):623–648.
- Carvalho, L. S., Meier, S., and Wang, S. W. (2016). Poverty and economic decision-making: Evidence from changes in financial resources at payday. *American Economic Review*, 106(2):260–84.
- Cespedes, J., Parra, C., and Sialm, C. (2020). The effect of principal reduction on household distress: Evidence from mortgage cramdown. *Available at SSRN*.
- Chetty, R., Friedman, J. N., Olsen, T., and Pistaferri, L. (2011). Adjustment costs, firm responses, and micro vs. macro labor supply elasticities: Evidence from danish tax records. *The Quarterly Journal of Economics*, 126(2):749–804.
- Dávila, E. and Hébert, B. M. (2019). Optimal corporate taxation under financial frictions. *NBER*

Working Paper.

- Donaldson, J. R., Piacentino, G., and Thakor, A. (2019). Household debt overhang and unemployment. *The Journal of Finance*, 74(3):1473–1502.
- Epper, T. (2017). Income expectations, limited liquidity, and anomalies in intertemporal choice. *Working Paper.*
- Golosov, M. and Tsyvinski, A. (2015). Policy implications of dynamic public finance. *Annual Review of Economics*, 7(1):147–171.
- Keane, M. P. (2011). Labor supply and taxes: A survey. *Journal of Economic Literature*, 49(4):961–1075.
- Kleven, H. J. (2016). Bunching. *Annual Review of Economics*, 8:435–464.
- Lockwood, B. B. (2020). Optimal income taxation with present bias. *American Economic Journal: Economic Policy*, 12(4):298–327.
- Lozachmeur, J.-M. (2006). Optimal age-specific income taxation. *Journal of Public Economic Theory*, 8(4):697–711.
- Ndiaye, A. (2020). Flexible retirement and optimal taxation. *Working Paper.*
- Ong, P. (2020). The effect of child support on labor supply: An estimate of the frisch elasticity. *Working paper.*
- Saez, E. (2010). Do taxpayers bunch at kink points? *American Economic Journal: Economic Policy*, 2(3):180–212.
- Scott, J., Shoven, J. B., Slavov, S., and Watson, J. G. (2021). Is automatic enrollment consistent with a life cycle model? *NBER Working Paper.*
- Seim, D. (2017). Behavioral responses to wealth taxes: Evidence from sweden. *American Economic Journal: Economic Policy*, 9(4):395–421.
- Søgaard, J. E. (2019). Labor supply and optimization frictions: Evidence from the danish student labor market. *Journal of Public Economics*, 173:125–138.
- Stantcheva, S. (2020). Dynamic taxation. *Annual Review of Economics*, 12:801–831.
- Zator, M. (2019). Working more to pay the mortgage: Household debt, consumption commitments, and labor supply. *Working paper.*