# What Determines 401(k) Plan Fees? A Dynamic Model of Transaction Costs and Markups<sup>\*</sup>

Hanbin Yang<sup> $\dagger$ </sup>

March 30, 2025

#### Abstract

I show that most reductions in 401(k) plan fees over the past decade come from updating plan menus to incorporate newly introduced funds with lower fees. Because employer sponsors face transaction costs when selecting and switching plan providers, providers can delay menu updates, preventing participants from accessing these lower-fee funds. To quantify the impact of transaction costs, I develop a dynamic structural model of employers' provider choice and providers' fee competition. I estimate that transaction costs contribute 11 bps to plan fees on average or \$1 billion in total. However, mitigating transaction costs has limited effects once forward-looking providers revise their fee strategies in response. By contrast, consolidating plans of small employers can generate substantial fee savings.

<sup>&</sup>lt;sup>\*</sup>I am indebted to the guidance and support of John Campbell, Ariel Pakes, Mark Egan, Robin Lee, and Adi Sunderam. I benefited from discussions with Vivek Bhattacharya, João Cocco, Yufeng Huang, Samuel Jäger, Yizhou Jin, David Laibson, Alex MacKay, Elie Tamer, Pietro Tebaldi, and seminar participants at Boston College Carroll School of Management, Columbia Business School, European Winter Finance Conference, Harvard Finance and IO workshops, International Industrial Organization Conference, London Business School, London Financial Intermediation Theory Workshop, London School of Economics, University of North Carolina, and University of Washington Foster School of Business. All errors are my own.

<sup>&</sup>lt;sup>†</sup>London Business School. Email: hyang@london.edu

# 1 Introduction

Employer-sponsored defined contribution plans, primarily 401(k) plans, are a critical component of retirement savings in the U.S., with over \$7 trillion in assets and 60 million active participants as of 2021 (ICI, 2022a). These plans allow employee participants to accumulate retirement savings through tax-advantaged investments, but participants pay investment and services fees. Plan fees can have a significant impact on retirement savings. Even a 1 percentage point increase in fees can reduce a participant's 401(k) balance at retirement by 28 percent (DOL, 2019).

401(k) plan fees have come under increasing scrutiny in recent years. Employee participants have filed a growing number of lawsuits against their employers for offering expensive investment options or permitting high service fees. Notably, *Hughes v. Northwestern* reached the U.S. Supreme Court in 2022. Using data on 401(k) plans, I show that the average 401(k) plan fee across employers is 65 basis points (bps) as of 2019, with 10% of employers having fees above 100 bps. These levels of plan fees may seem high. As a comparison, individuals can invest in an S&P 500 index fund by paying less than 10 bps. Federal employees also pay less than 10 bps in fees in Federal Thrift Saving Plans (TSP)<sup>1</sup>.

While plan fees are paid by participants, employers negotiate these fees on behalf of participants with plan providers. Often referred to as recordkeepers, plan providers are financial firms such as Vanguard and Fidelity, who provide customized menus of investment options as well as administrative or advisory services. While these lawsuits and much of the existing literature focus on agency costs between employers and participants, there has been limited attention on frictions between employers and providers. As in many business-to-business interactions, employers may face transaction costs, such as the costs of selecting and switching providers. As a result, providers could exploit these transaction costs and charge high plan fees.

In this paper, I study how employers' transaction costs of selecting and switching providers affect 401(k) plan fees. Using comprehensive data on 401(k) plans between 2009 and 2019, I document an important and novel pattern of plan fees: menu updates to incorporate low-cost funds explain the majority of the secular fee decline during this sample period. Fund investment fees in 401(k) plans decrease by on average 30 bps from 2009 to 2019. Only 30 percent of this decline comes from fee reductions among existing funds included in 2009 plan menus, whereas the remaining 70 percent results from menu updates to incorporate funds that are similar to existing ones but have lower fees. Since employers need to work with providers to update plan menus, any transaction costs employers incur when interacting with providers can lead to delays in menu updates. As a result, participants cannot access funds that are cheaper than the existing ones on their menus.

The nature of fee decline means that providers could exploit employers' transaction costs by delaying menu updates. Employers face two main sources of transaction costs when dealing with providers. First, employers incur costs when choosing providers. Employers use a Request for Proposal (RFP) process to solicit customized bids from plan providers. Second, employers face costs when switching providers. I find that employers reduce fees by over 10 bps when they switch providers. However, only 3% to 5% of employers switch annually, consistent with transaction costs. Transaction costs can include monetary, time, and psychological costs. I do not differentiate between these sources in my analysis.

<sup>&</sup>lt;sup>1</sup>https://www.tsp.gov/tsp-basics/expenses-and-fees/

Building on these findings, I develop a dynamic structural model of transaction costs in a negotiated-price market. Since providers offer customized services and negotiate fees separately with each employer, I treat each employer as an independent market and model how multiple providers compete to serve a single employer. The model comprises three key components. First, the employer determines whether to conduct a costly RFP for potential fee reductions, as an optimal stopping problem (Rust, 1987). Second, the employer chooses a provider at an RFP, where the employer pays switching costs if choosing a non-incumbent provider. I model provider choice using a random-coefficient discrete choice model (Berry et al., 1995). Third, providers compete in a dynamic fee-setting game (Dubé et al., 2009; Cabral, 2016; Sweeting et al., 2022), where forward-looking providers may initially set low fees, anticipating smaller future declines due to transaction costs. I jointly solve these three components and compute outcomes under the Markov perfect equilibrium.

I overcome two challenges when identifying my model. First, I do not observe RFPs in my data, and I use the autocorrelation of fees to infer the probability of RFPs. I can group employer  $\times$  year observations into three cases: switching providers (assuming always after an RFP), no menu turnover (assuming no RFP), and lastly no switching with some menu turnover. The fee autocorrelation is close to zero when employers switch providers and reaches its upper bound without menu turnover. Then, I can infer the probability of RFPs in the last case by comparing its fee autocorrelation versus zero or the upper bound. Second, how employers trade off transaction costs versus fees depends on a fee sensitivity parameter, which reflects the employer's sophistication or effort. Identifying this parameter is difficult because I do not observe fees from non-chosen providers, a common issue in negotiated-price markets. In my setting, fee reductions when employers switch providers can approximate the differences in proposed fees between incumbent and non-incumbent providers. A small fee reduction from switching indicates that the employer is highly fee-sensitive, preventing the incumbent provider from exploiting switching costs to impose high fees.<sup>2</sup>

I estimate my model separately for employers in each quartile of plan assets and recover reasonable estimates of transaction costs. Transaction costs measured in basis points over plan assets decrease in employer size. From small employers in the first quartile to large employers in the fourth quartile, RFP costs range from 10 to 1 bps, while switching costs range from 46 to 2 bps. Scaled by total plan assets, RFP costs range from \$7,000 to \$18,000, and switching costs from \$30,000 to \$57,000. Based on estimated transaction costs, the difference between observed plan fees and fees employers could obtain should they switch providers is 11 bps for the average employer or \$1 billion per year in total.<sup>3</sup>

Based on my model estimates, I recover markups on plan fees by inverting providers' optimal fee-setting conditions. I find that markups account for 18% over plan fees on average. Larger employers tend to have lower markups. Specifically, markups for small employers are 18 bps or 23% over plan fees, decreasing to 6 bps or 13% for large employers. Large employers have lower markups because I estimate that they are more fee-sensitive when trading off transaction costs

 $<sup>^{2}</sup>$ I find that employers with more plan assets, more participants, higher employer contribution, higher participation rates, and higher ESG scores have lower fee reductions from switching providers, suggesting this variation is informative of employers' sophistication or effort captured by their fee sensitivities.

<sup>&</sup>lt;sup>3</sup>In Section 3.2, I discuss several restrictions when constructing my estimation sample and report aggregate fees in my sample only. The overall impact of these aggregated counterfactual outcomes can be three to five times larger when applied to the full universe of 401(k) plans with comparable plan assets.

versus fees. My estimates also point to economies of scale, where providers incur lower costs (plan fees minus markups) of providing 401(k) services to larger employers. I refer to providers' costs as production costs to distinguish them from employers' transaction costs. Production costs decline from 61 bps to 42 bps from small to large employers, explaining two-thirds of the difference in plan fees across employer sizes.

Having estimated the model, I consider counterfactual policies aimed at mitigating transaction costs. As highlighted in previous literature (Klemperer, 1995; Dubé et al., 2009; Cabral, 2016), transaction costs can influence fees through two opposing mechanisms. On one hand, providers exploit the transaction costs of their existing employers by delaying menu updates and maintaining high fees. On the other hand, anticipating high fees in the future, providers initially set lower fees when competing for employers at RFPs. I find that when both RFP and switching cost parameters are set to zero, fees remain roughly the same as the status quo, as providers respond by setting higher fees at RFPs in the new equilibrium. Interestingly, moderate transaction costs incentivize providers to compete more aggressively than the status quo, reducing plan fees by as much as 4 bps on average or \$500 million per year in aggregate. One novel mechanism in my model is that providers have an incentive to set lower fees to reduce the probability that employers run future RFPs. I show that if employers' RFP decisions become more fee sensitive, plan fees can be lower by 3 bps while employers incur similar levels of transaction costs as the status quo.

I also consider plan consolidation since I estimate that providers incur lower production costs and charge lower markups for larger employers. Specifically, I consolidate plans of smaller employers to match the average plan assets of large employers in the fourth size quartile. Consolidation can reduce fees by on average 16 bps or \$400 million annually, through reductions in both markups and production costs. Additionally, employers save \$110 million in transaction costs by eliminating duplicated RFPs and provider switches. Even without modifying employer preferences that lead to low markups, plan consolidation still reduces fees by on average 11 bps or \$300 million annually, thanks to lower production costs alone. This counterfactual suggests that recent policies facilitating multiple employer plans (MEPs) under the SECURE Act of 2019 could potentially benefit both participants and employers.

I make four contributions with this paper. First, I document a new channel of how providers exploit employers' transaction costs by delaying the inclusion of newly introduced funds with lower fees, contributing to the literature on friction in retirement plan designs. Among existing papers in this literature<sup>4</sup>, Pool et al. (2016, 2021) show how providers steer plan menus toward their proprietary funds and funds with revenue-sharing. Doellman and Sardarli (2016) and Badoer et al. (2020) show how providers trade off direct compensations from sponsors with indirect compensations through revenue sharing. Chalmers and Reuter (2020) and Reuter and Richardson (2022) study the demand for financial advice and the impact of advice on asset allocation in 401(k) plans. Loseto (2023) and Gropper (2023) study how plan menu design affects participants' welfare. Bhattacharya and Illanes (2022) study misaligned incentives of

<sup>&</sup>lt;sup>4</sup>More broadly, a large literature in household and behavioral finance focuses on employees' participation and contribution. See Benartzi and Thaler (2007) and Choi (2015) for a discussion of this literature. Papers in this literature have studied the effects of plan design such as automatic enrollment (Madrian and Shea, 2001; Beshears et al., 2009; Choi et al., 2007; Carroll et al., 2009) and firm matching (Choi et al., 2002; Duflo et al., 2006; Dworak-Fisher, 2011). Yogo et al. (2025) study access to retirement plans and participation across US households using administrative data.

employers and imperfect competition among providers with a static bargaining model between employers and providers. These papers tend to study providers' incentives in static settings, whereas my focus is on the dynamic implications of transaction costs.

Second, my counterfactual simulations provide insights into recent regulations and lawsuits involving 401(k) plans. My estimates suggest that fees are lower at larger employers due to both lower markups and lower production costs. Hence, consolidating plans of small employers could potentially lead to substantial fee savings. Bhattacharya and Illanes (2022) also study plan consolidation counterfactual but find that consolidation only reduces fees through modifying employer preference. In contrast, I estimate stronger economies of scale. In addition, my decomposition suggests that differences in employer preferences do not explain the majority of the dispersion in 401(k) plan fees. Much of this dispersion can be natural outcomes in a negotiated price market with heterogeneous services. Therefore, 401(k) plan fees can appear high in a cross-sectional comparison, even though employers plausibly fulfill their fiduciary duties. My results are also relevant for discussions of whether employers should periodically run RFPs in recent lawsuits. I show that instead of more frequent RFPs, more fee-sensitive RFPs can reduce fees without incurring additional burdens on employers.

Third, this paper provides an empirical application of the mostly theoretical literature on dynamic competition with switching costs. Klemperer (1995) and Farrell and Klemperer (2007) provide early summaries of this literature. More recently, Dubé et al. (2009) and Cabral (2016) show that the equilibrium impact of switching costs on prices is ambiguous due to the opposite effects of "invest" and "harvest" incentives, and that prices are generally the lowest with modest levels of switching costs. While existing papers in this literature assume demand-side agents are myopic, I allow both providers and employers to be forward-looking by endogenizing employers' RFP decisions. In my model, providers have an additional incentive to "invest" to reduce the probability that employers run another RFP.

Lastly, I develop novel estimation and identification strategies, contributing to the empirical industrial organization literature, especially studies of financial markets with negotiated prices.<sup>5</sup> A common challenge in negotiated-price markets is that fees set by non-chosen providers are unobserved. Previous papers (Allen et al., 2019; Allen and Li, 2025; Cuesta and Sepúlveda, 2021) use results from second-price auctions to infer the distribution of unobserved fees. I cannot directly adopt this approach because my model is dynamic. My paper is also related to studies of fee dispersion in financial markets. Prior papers primarily focus on fee dispersion within the same or homogeneous products, such as index funds (Hortaçsu and Syverson, 2004) and private equities (Begenau and Siriwardane, 2024). In contrast, 401(k) plans are potentially heterogeneous, and I develop an estimation strategy to allow for both observed and unobserved heterogeneity.

The rest of the paper is structured as follows. I discuss institution details in Section 2 and data in Section 3. In Section 4, I present motivating evidence of patterns of fee decline and transaction costs. I introduce the structural model in Section 5, explain estimation and identifi-

<sup>&</sup>lt;sup>5</sup>See Clark et al. (2021) for an overview of the broader literature on the industrial organization of financial markets. Previous works in this literature have studied car loans (Einav et al., 2012; Grunewald et al., 2020), credit cards (Nelson, 2018), insurance (Koijen and Yogo, 2016), mortgages (Allen et al., 2014, 2019; Robles-Garcia, 2019), municipal bonds (Brancaccio et al., 2017), pensions (Luco, 2019; Illanes, 2016; Illanes and Padi, 2019), financial advice (Bhattacharya et al., 2019; Egan, 2019; Guiso et al., 2022), and consumer and student loans (Bachas, 2018; Cuesta and Sepúlveda, 2021). Egan et al. (2022, 2023) estimate investor demand to recover beliefs and risk preferences.

cation in Section 6, and present estimation results in Section 7. I discuss policy counterfactuals in Section 8 and conclude in Section 9.

# 2 Institutional Background

This paper examines the interaction between employer sponsors of 401(k) plans and plan providers. More broadly, there are four main types of agents in the 401(k) sector: employers, employee participants, plan providers, and mutual fund companies. In the following, I introduce these agents and discuss the rationale for focusing on employer-provider interactions. Then, I discuss employers' transaction costs when selecting and switching providers. Lastly, I summarize allegations in excessive fee lawsuits.

# 2.1 Employer Sponsors of 401(k) Plans

Private sector employers sponsor 401(k) plans, allowing employees to make tax-deferred contributions from their salaries into these accounts. Employees then invest these contributions in various investment vehicles to save for retirement.<sup>6</sup>

The Department of Labor regulates employer sponsors under the Employee Retirement Income Security Act of 1974 (ERISA). ERISA imposes fiduciary responsibilities on employer sponsors, requiring them to act solely in the best interest of their plan participants. Employers must exercise prudence in selecting and reviewing investment options in their plan menus, including due diligence concerning risk, performance, fees, and diversification. Fiduciary duties also extend to the selection of plan providers and the assessment of their performance.

### 2.2 Plan Providers

Often referred to as recordkeepers or trustees, financial service providers administer 401(k) plans and design plan menus. They process employee contributions and distributions, set up online portals for participants to access account details and execute transactions, and may also offer educational resources to help participants with investment decisions. In Figure 1a, I display market shares based on the number of plans and assets for providers as of 2019. Some of these providers are asset managers, such as Fidelity and Vanguard.<sup>7</sup> Empower, ADP, and Principal Financial Group are third-party providers. They do not have proprietary investment products but may offer other bundled services. For example, ADP offers HR payroll services, while Principal Financial Group offers insurance products. Many banks, insurance companies, and asset managers also provide recordkeeping services. Due to their relatively small market shares, I group them collectively into the "other" category. In Figure 1b, I show that large providers have maintained stable market shares throughout my sample period.<sup>8</sup>

 $<sup>^{6}</sup>$ A 401(k) is a type of defined contribution (DC) plan. Other employer-sponsored DC plans include 403(b) plans for employees in education and nonprofit organizations, 457 plans for certain governmental employees, and Thrift Savings Plan (TSP) for federal government employees. My data and analysis include 401(k) and 403(b) plans.

<sup>&</sup>lt;sup>7</sup>When designing plan menus, recordkeepers who are asset managers typically favor their proprietary funds but may also offer funds from other mutual fund companies. The fraction of proprietary funds among Vanguard and Fidelity plans are around 90% and 50% on average, according to Appendix Table A3.

<sup>&</sup>lt;sup>8</sup>There are several consolidations among large recordkeepers after the end of my sample. For example, Empower acquired MassMutual's retirement businesses in 2010 and Prudential's in 2022. Consolidations in my sample are between smaller providers and are adjusted accordingly, as detailed in Appendix A.

### Figure 1: Provider Market Shares



Panel (a) plots market shares of large providers both in number of plans and total assets as of 2019, where smaller providers are grouped into the other category. Panel (b) plots the market share of providers over time, where I restrict to a balanced panel.

#### 2.3 401(k) Plan Fees, Participants, and Fund Managers

401(k) plan fees consist of investment fees and recordkeeping fees. First, employee participants pay investment fees on their selected funds, such as mutual fund expense ratios, to the asset managers who operate those funds.<sup>9</sup> Second, employee participants (and occasionally employers) pay recordkeeping fees to plan providers. Recordkeeping fees can be structured as either a fixed dollar amount per participant or a percentage of invested assets. Deloitte (2013) shows that investment fees make up 82% of total plan fees, with participants covering 87%.

In my analysis, I abstract away from the interaction between employers and their participants. I assume employers internalize fees participants pay and benefits participants receive from 401(k) plans when they choose plan providers and negotiate plan fees. Although individual participants may pay different fees based on their fund allocations, their choices are limited by the plan menus. As a result, employers' decisions play a significant role in determining the fees that participants pay.

On the other hand, since investment fees are paid to fund managers, there is a vertical relationship between recordkeepers and fund managers. I abstract away from analyzing this vertical relationship. Integrated providers that are asset managers themselves (e.g., Fidelity and Vanguard) receive investment fees directly. Third-party providers indirectly obtain a portion of investment fees through "revenue sharing" with fund managers. I assume plan providers maximize joint profit and can use "revenue sharing" or other transfers to align incentives with fund managers.<sup>10</sup>

#### 2.4 How Employers Choose and Monitor Providers

Employers incur transaction costs when choosing and switching providers. First, employers choose providers by running a Request for Proposal (RFP), a standard procedure in private-sector procurement. The RFP process can be complex and time-consuming. Employers may

 $<sup>^{9}</sup>$ A mutual fund can have multiple share classes with different expense ratios. Investors across different 401(k) plans pay the same expense ratio for the same share class.

<sup>&</sup>lt;sup>10</sup>Pool et al. (2016, 2021) and Bhattacharya and Illanes (2022) study how vertical integration and revenue sharing create misaligned incentives in menu design.

hire external consultants to navigate the RFP process, leveraging their industry knowledge and legal expertise. Potentially due to the time costs of handling RFPs and the monetary costs of hiring consultants, employers typically conduct an RFP every three to five years.<sup>11</sup>

Second, switching providers can create administrative burdens and disrupt participants' experience. Employers also value their overall relationships with their current providers, especially since their providers may offer other services.<sup>12</sup> According to Deloitte (2015, 2019), 45-50% of employers have maintained the same recordkeepers for over 10 years, with the primary reason being "overall relationships". In addition, there might be contractual restrictions or penalties associated with switching providers. In *Hughes v. Northwestern*, a significant portion of Northwestern University's plan asset is invested in TIAA annuity, with TIAA also serving as the recordkeeper. The annuity has specialized administrative requirements and substantial early withdrawal penalties. As a result, it would be highly costly for Northwestern University to switch away from TIAA.

After selecting a provider, employers need to monitor the provider's performance. Employers typically schedule annual or quarterly review meetings with providers, where they may discuss potential changes to plan menus or other services. Employers may also hire external consultants to benchmark their fees against a comparable peer group or conduct Requests for Information (RFIs) to solicit fee quotes from other potential providers. Fee comparison results from benchmarking or RFIs can help employers negotiate fees with their current providers. However, since providers offer customized services to each employer, benchmarking or RFIs may not fully capture the fees employers could obtain if they were to conduct a formal RFP.

#### 2.5 Excessive Fee Lawsuits

In Figure 2, I show the number of lawsuits bought against employer sponsors of 401(k) plans each year. The number of lawsuits peaked during the Great Recession when most cases focused on inappropriate investment choices that led to losses in participants' retirement savings. Recently, "excessive fees" lawsuits have surged again and focus primarily on plan fees.

I discuss three main allegations in recent lawsuits and how these legal arguments motivate key assumptions in my empirical approach. First, because providers design menus and set recordkeeping fees jointly, I assume employers and providers negotiate over total plan fees combining investment and recordkeeping fees together. Second, I do not control for index versus active funds on plan menus. There can be heterogeneous preferences across participants in different employers regarding active management. Alternatively, providers can charge high fees by offering expensive active funds. Third, I allow for unobserved heterogeneity in services and fund preferences across employers that affect plan fees.

Offering retail share classes with higher expense ratios instead of institutional share classes of the same mutual funds. According to the plaintiffs, offering a more expensive share class of the same mutual fund suggests employers are negligent in monitoring their plan menus. However,

<sup>&</sup>lt;sup>11</sup>In two polls conducted by the National Association of Plan Advisors (Adams, 2020, 2022), 50% and 61% plan advisors recommend that employers conduct an RFP every 3-5 years. See Appendix Figure A7 for a sample RFP timeline. There is no specific legal requirement to conduct RFPs. Although some plaintiffs argue that employers breach their fiduciary duties because they do not regularly solicit competitive bids with RFPs, courts tend to dismiss such claims, including in *Sacerdote v. New York University* and *Albert v. Oshkosh Corporation*.

 $<sup>^{12}</sup>$ According to Deloitte (2013), only 56% of employers in their survey reported having no other relationships outside their 401(k) plans.

Figure 2: Number of Cases Against Employer Sponsors of 401(k) Plans



Statistics based on Mellman and Sanzenbacher (2018); Aronowitz (2023)

although retail share classes have higher expense ratios, they may include higher revenue-sharing to offset recordkeeping fees paid by participants. On a net basis, a retail share class could result in lower overall fees compared to the institutional share class of the same fund. In *Forman v. TriHealth*, the court dismissed claims that the employer violated fiduciary duties by selecting a more expensive share class of the same mutual fund. In addition, to satisfy higher minimum investment thresholds of institutional share classes, employers may need to limit the number of funds offered, thereby restricting investment options for participants. <sup>13</sup>

Offering actively managed funds instead of the index funds with lower expense ratios. Similar to the previous claim, plaintiffs argue that selecting more expensive active funds suggests negligence. While active funds may underperform their index counterparts (Gruber, 1996), no regulatory guidance requires employer sponsors to offer index funds. In *Smith v. Common-Spirit Health*, the court ruled that active funds are a common component of retirement plans and employers should offer active funds as an option for risk-tolerant investors among plan participants.

Allowing plan providers to charge higher recordkeeping fees than comparable employers. Similar to claims regarding share classes, plaintiffs argue that higher recordkeeping fees indicate employer negligence in monitoring and negotiating plan fees. Conversely, employer sponsors tend to argue that they offer a different scope or quality of services than the employers used as benchmarks in participants' comparisons. Plaintiffs typically compare their recordkeeping fees with those of five to ten employers of similar sizes and rarely provide direct evidence that another provider would offer a lower fee for identical services. In *Matousek v. MidAmerican Energy Company, Smith v. CommonSpirit Health*, and *Albert v. Oshkosh Corp.*, courts ruled that plaintiffs failed to present sufficient evidence that employer sponsors allow higher recordkeeping fees for comparable services.

I do not have direct evidence on the extent of service heterogeneity across providers or employers. Industry practitioners tend to argue that such heterogeneity could exist (Aronowitz, 2022b). As indirect evidence, survey results from Deloitte (2019) indicate that when asked to rank top improvements from recordkeepers, 18% of employers cited participant readiness for retirement, 15% mentioned plan sponsor websites or tools, and 11% prioritized participant

<sup>&</sup>lt;sup>13</sup>See detailed argument in the Investment Company Institute's Amicus Brief https://www.supremecourt.gov/ DocketPDF/19/19-1401/198014/20211028135053555\_19-1401%20ICI%20Amicus%20Brief.pdf

experience. In contrast, only 10% and 7% identified direct recordkeeping fees and investment fees, respectively, as top concerns. The fact that more employers focus on service quality over fees suggests that they may not view recordkeeping as a homogeneous service.

# 3 Data and Summary Statistics

### 3.1 Data

The primary dataset for this study comes from BrightScope Beacon. BrightScope Beacon provides plan and fund level information for ERISA defined contribution plans, covering 84% of total plan assets. BrightScope collects data from plan sponsors and publicly available sources, including the Department of Labor's Form 5500 and the Securities and Exchange Commission (SEC). The dataset covers 70,000 401(k) plans from 2009 to 2019. For each 401(k) plan, BrightScope reports annual data on investment menus and the total assets allocated by participants to each investment option. The data also contain fund asset class classifications (e.g. large cap equities, bonds) and investment styles based on Morningstar categories.

BrightScope collects the identities of service providers, the service codes that allow me to identify the primary recordkeeping provider, and the total dollar amount of direct compensation. However, I do not observe whether recordkeeping fees are paid by employers or employees, or whether recordkeeping fees are structured as a dollar amount per participant or as a percentage of plan assets. To be consistent with investment fees, I measure direct compensation as a fraction of total plan assets. Third-party recordkeepers may also receive indirect compensation from fund managers in exchange for including their mutual funds on plan menus. To avoid double-counting expense ratios, I exclude indirect compensation from my analysis. Occasionally, some plans have multiple recordkeepers. In such cases, I designate the provider with the highest compensation or the longest tenure as the primary recordkeeper. See Appendix A for further discussion. It is worth noting that direct compensation data is noisy and may overestimate recordkeeping fees (Aronowitz, 2022a).

While I observe providers, I do not have data on when employers conduct RFPs to choose providers.

#### 3.2 Sample Construction

I merge investment menu data from BrightScope with historical expense ratios from CRSP. To construct my sample, I focus on employers that have at least five mutual funds on their menus and on average over 70% of assets allocated to mutual funds. This step eliminates some large corporations that typically request providers to customize funds only for their plan participants. While I observe these customized funds in my data, their historical expense ratios are not available. Additionally, I exclude employers with fewer than 100 participants. These employers are more likely to outsource the management of their 401(k) services and also delegate part of their fiduciary duties to providers.

Appendix Table A1 shows that my estimation sample represents 17% of 401(k) plan assets. The sample is less representative of the smallest employers due to missing menu information and of the largest employers because they are more likely to offer customized funds. My sample

is more representative of employers in the 25th to 75th percentiles of the size distribution, with plan assets ranging between \$10 million and \$40 million.

My dataset is an unbalanced panel. Of the 19,037 employers in my estimation sample as of 2019, 2,476 or 13% have data for each year. More employers began reporting Form 5500 to the Department of Labor over time. Among employers who consistently report Form 5500, plan menu and recordkeeping compensation are sometimes missing. I use the full estimation sample for structural estimation and restrict to the balanced panel for certain analyses to avoid selection biases.

#### **3.3 Summary Statistics**

In Table 1, I present summary statistics for employers in my estimation sample as of 2019. The median employer has approximately 300 participants with \$20 million in total plan assets. After excluding large corporations that offer customized funds, my estimation sample includes \$1.1 trillion assets. On average, investment fees are 47 bps. This average is computed by calculating asset-weighted mutual fund expense ratios for each employer and then taking a simple average across employers. Direct recordkeeping fees average 17 bps relative to total plan assets. Combining investment and direct recordkeeping fees, total plan fees are 65 bps on average. The dispersion in plan fees is large, with a standard deviation of 27 bps. Plan fees range from 30 bps at the 10th to 100 bps at the 90th percentile.

	Num Obs	Mean	St.Dev.	Pct10	Pct50	Pct90
Plan participants	19,037	875	2,990	130	288	$1,\!600$
Total plan assets (million)	19,037	58	236	5	18	109
Average account balance (thousand)	19,037	86	91	17	59	181
Fee						
Plan average expense ratio (bps)	19,037	47	24	17	45	79
Direct recordkeeping fees (bps)	19,037	17	19	1	11	44
Total plan fees (bps)	19,037	65	27	31	62	101
Plan menu						
Number of funds	19,037	26	9	16	26	35
Menu turnover (%)	19,037	13	20	0	5	39
Provider switch						
Switched provider $(\%)$	19,037	4	20	0	0	0

Table 1: Summary Statistics of 401(k) Plans in 2019

Summary statistics across plans as of 2019.

On average, employers offer 26 funds in their plan menus. To measure menu updates, I use menu turnover, defined as the fraction of funds added or replaced in a given year. On average, 13% of funds in plan menus are updated annually. Additionally, 4% of employers switch their providers in 2019.

In Appendix Table A3, I present summary statistics by providers. Fidelity and Vanguard typically serve larger employers, whereas ADP primarily targets smaller employers. Vanguard charges the lowest expense ratios, partially because it offers a higher proportion of index funds. Despite similar index fund allocations, third-party providers such as ADP and Principal Financial Group offer funds with higher expense ratios than Fidelity. Vertical integration allows

Fidelity and Vanguard to offer more proprietary funds and charge lower recordkeeping fees.<sup>14</sup>

# 4 Motivating Evidence

I first show that most of the decline in plan fees is driven by menu updates to incorporate lowcost funds, instead of reductions in fees of existing funds. Because updating menus is important, employers' transaction costs when dealing with providers can lead to delays in menu updates that result in high fees. To motivate the impact of transaction costs, I show that employers do not switch providers frequently but achieve large fee reductions when they switch.

### 4.1 Fee Decline Due to Menu Update

According to Morningstar (2021), the asset-weighted average mutual fund expense ratio has declined from 87 to 40 bps from 2001 to 2021. I discuss factors that contribute to this trend in Appendix B. For example, mutual fund managers benefit from economies of scale. Their fixed costs of operating mutual funds, measured in basis points, decrease over time as these funds accumulate more assets under management.

I show that this secular decline in fees is primarily by menu updates that incorporate funds with lower expense ratios rather than reductions in expense ratios among existing funds. To illustrate this, I compare different time series of expense ratios in Figure 3 using a balanced panel of employers. Observed fees (solid black line) decreased by almost 30 bps on average from 2009 to 2019, reflecting changes in fund expense ratios, plan menus, and participant allocations. To construct this time series, I first compute asset-weighted expense ratios for each employer and then take the simple average across employers in a given year. On the other hand, the dashed black line with circles on the top of the graph only accounts for expense ratio changes among existing funds, using fund menus and allocation as of 2009. When menu updates are not considered, the fee reduction is less than 10 bps, representing less than 30% of the total observed decline. I discuss the construction of these time series in further detail in Appendix E.1.

While participants can rebalance toward lower-cost funds, they may be constrained by available funds on their plan menus. The gray line in Figure 3 allows rebalancing among funds in 2009 menus, which is close to the dashed line without rebalancing. In addition, fee reductions from menu updates could also be driven by changes in participant preference, such as a shift toward index funds. In Appendix Figure A3, I show similar patterns within active and index funds. Menu updates allow participants to access cheaper funds within the active or index universe.

Because employers need to work with providers to update plan menus, providers can exploit employers' transaction costs to delay menu updates and maintain high fees. I find that 70% of the funds in 2019 menus were already available in the marketplace as early as 2009. Had employers offered available funds from their 2019 menus to their participants in 2009, they could have reduced average expense ratios as of 2009 by 20 bps, as shown by the dashed line with

<sup>&</sup>lt;sup>14</sup>Based on my conversations with industry practitioners, although integrated providers charge lower plan fees, they tend to promote their brokerage accounts and other products to participants, which can negatively impact participant welfare.





Figure shows average expense ratios of funds in 401(k) plans. I compute asset-weighted expense ratios for each employer and take simple averages across employers. The black solid line uses contemporaneous menus and allocation. The dashed line with circle markers is constructed using funds in 2009 menus, capturing only changes in expense ratios. The gray line with circle markers allows participants to rebalance among funds in their menus as of 2009. The dashed line with triangle markers is constructed using funds in 2019 menus that are available during each calendar year.

triangle markers in Figure 3.<sup>15</sup>

Just as mutual fund fees have declined over time, recordkeeping fees have also experienced a secular decline, thanks to factors such as improvements in providers' IT infrastructure. See Appendix B for further details. Providers can also take advantage of employers' transaction costs by delaying the renegotiation of recordkeeping fees. I focus on mutual fund expense ratios in this analysis, because patterns of recordkeeping fees are less obvious due to measurement errors in recordkeeping fee data.

#### 4.2 Employers Face Transaction Costs

I demonstrate that employers switch providers infrequently but can substantially reduce fees when they switch, consistent with transaction costs associated with conducting RFPs and switching providers, as discussed in Section 2.4.

In Table 1, I show that around 4% of employers switch providers each year.<sup>16</sup> I then use a difference-in-differences specification to estimate fee reductions when employers switch providers. Using a balanced panel of employers, I define the treatment group as employers that switched providers once between 2012 and 2016, while the control group consists of employers that never switched providers. The specification is as follows:

$$fee_{ijt} = \sum_{s \neq -1} \theta_s 1\{t = s\} + \Gamma_i + \Gamma_j + \tau_{jt} + X_{it}\beta + \epsilon_{ijt}$$
(1)

<sup>&</sup>lt;sup>15</sup>As another piece of evidence supporting delay in menu updates, in Appendix E.2, I show that when mutual fund companies introduce cheaper share classes of the same mutual funds, providers do not immediately incorporate those into plan menus even though employers are eligible based on their investment size.

 $<sup>^{16}</sup>$ In Figure A5a, I show a time-series of the fraction of employers who switch providers each year, which is between 3% and 5%.

The dependent variable is the plan fee, which includes both investment and recordkeeping fees, of employer *i* at year *t* with provider *j*. *s* corresponds to the number of years relative to the year of switch, with s = -1 assigned to employers that never switch providers. The coefficients of interest are denoted by  $\theta_s$ , and I normalize  $\theta_s = 0$  for s = -1. I control for employer and provider fixed effects  $\Gamma_i, \Gamma_j$ , provider-specific time trend  $\tau_{jt}$ , and plan characteristics  $X_{it}$  which include menu composition, types of services, and number of plan participants.

Figure 4: Difference in Differences When Employers Switch Providers



Figure displays diff-in-diffs coefficients  $\theta_s$  in Equation (1). Vertical bars show 95% confidence intervals. I cluster standard errors at the employer level. Due to staggered provider switches, I follow the procedure in Sun and Abraham (2021) where I estimate  $\theta_s$  separately for employers who switch each year and take weighted averages. Standard errors are bootstrapped over 200 repetitions.

I plot the coefficients  $\theta_s$  in Figure 4. There is an immediate reduction of over 10 bps of plan fees when employers switch providers. This reduction is around 16% relative to average plan fees in 2019 shown in Table 1. The large fee reduction combined with a low switching probability suggests substantial transaction costs. After switching providers, fee reductions are persistent but attenuate by half over three years. Even without switching providers, employers in the control group may reduce their fees by inviting other providers to compete with their existing providers during RFPs. Since employers are unlikely to conduct another RFP immediately after switching providers, plan fees tend to increase relative to the control group.

One concern is that employers may change their plan characteristics when they switch providers. My controls for menu composition and services help mitigate this concern. In Figure A5c, I estimate the same difference-in-differences specification where the outcome variable is computed using average expense ratios across all funds of the same asset class  $\times$  Morningstar category. This measure reflects changes in menu composition, which has a minor and almost statistically insignificant change after employers switch providers. Additionally in Figure A5b, the outcome variable corresponds to fees in excess of the cheapest eligible share class of the same mutual funds. Different share classes of the same mutual funds are identical products but charge different expense ratios. There is a statistically significant reduction of over 4 bps when employers switch providers. Thus, fee reductions resulting from provider switches are unlikely to be explained away by changes in plan characteristics.

# 5 Model

In this section, I introduce a structural model based on institutional details and motivating evidence discussed previously. Since providers customize services and negotiate fees with each employer, each employer represents her own market. I model how multiple providers compete to serve one single employer.<sup>17</sup> I assume the employer determines menu composition and types of services. Providers cater to the employer's demand and compete on fees.<sup>18</sup> I allow both providers and the employer to be forward-looking.

The employer (she) maximizes the utility her employee participants derive from the 401(k) plan, net of her transaction costs. The employer first decides whether to conduct an RFP, trading off the present value of potential fee reductions against RFP transaction costs  $\kappa_{rfp}$ . I model this decision as an optimal stopping problem (Rust, 1987), as potential fee reductions increase over time relative to the declining trend.

Second, the employer chooses a provider at an RFP. When making provider choices, the employer considers fees proposed by providers and provider quality. The employer also incurs a switching cost  $\kappa_{sw}$  when choosing a non-incumbent provider. I allow for heterogeneous fee sensitivities among employers, capturing how they trade off transaction costs against plan fees. I model the employer's provider choice at an RFP using a random-coefficient discrete choice model (Berry et al., 1995).

The Provider (he) competes with other providers to offer services to the employer. Since providers cater to the employer's preference for menus and services, they compete on fees. While fund expense ratios are set by fund managers, providers can effectively determine plan fees by selecting funds for plan menus and setting recordkeeping fees. In my model, providers compete in a dynamic fee-setting game at each RFP. Dynamic incentives can push providers to set fees below static optimal levels, if they anticipate earning higher fees in the future by exploiting employers' transaction costs.

Providers may also update menus outside of RFPs. For simplicity, I assume providers make binary decisions of whether to update menus, with a constant and exogenous probability.

**Stationarity**: My model is stationary for analytical tractability. Although plan fees decline over time, I attribute this decline to reductions in providers' costs of offering funds and services (production costs), without affecting the nature of provider competition. As shown in Figure 1b, the market shares of large providers remain stable over time, supporting the stationary assumption.

In Figure 5, I provide a graphical illustration of how observed plan fees remain roughly constant and decrease following menu updates in an RFP. If we control for the secular trend, the residualized fees increase over time until the next RFP, shown by the right hand side panel. In Section 6.1, I discuss how I residualize observed fees with a hedonic regression before structural estimation using my stationary model.

State variables: The model is dynamic in discrete time with two state variables: i) the identity of the incumbent provider s, ii) the residualized fee f. For the following of this model section, I will refer to state variable f as just fees instead of residualized fees to simplify

<sup>&</sup>lt;sup>17</sup>By focusing on a single employer, I abstract away from cross-employer considerations, such as cross-subsidization across employers or using certain employers to signal provider quality.

<sup>&</sup>lt;sup>18</sup>I abstract away from how providers steer menu composition, because most providers have access to a wide range of funds and can be relatively indifferent to employers' preferences for funds.

exposition.





Timing: Within each period (year), the following happens sequentially.

- 0. The employer and all J providers observe the two state variables: the incumbent provider s and fee f.
- 1. The employer decides whether to incur RFP cost  $\kappa_{rfp}$  to run an RFP.
- 2.a. At the RFP, each provider sets his optimal fee  $b_i(s)$  in a dynamic fee-setting game.
- 2.b. The employer realizes her idiosyncratic preference for each provider and chooses the provider j who generates the highest utility. States transition from (s, f) to  $(j, b_j(s))$ 
  - 3. The current provider (the winner of the RFP j or the incumbent from the beginning of the period s) receives fee f from participants of the employer. Fee f transitions to the next period depending on whether the provider updates the menu.

I introduce the employer's provider choice at stage 2 in Section 5.1 and her decision to run an RFP in Section 5.2. Then, I specify providers' fee competition in Section 5.3. Sections 5.4 and 5.5 explain fee transitions and the equilibrium concept.

#### 5.1 Stage 2.b: Employer's Provider Choice at an RFP

I model the employer's provider choice at an RFP as a discrete choice (McFadden, 1973; Berry et al., 1995). The employer's choice depends on the identity of the incumbent provider s due to switching cost  $\kappa_{sw}$  but not on fee f before the RFP.

The employer *i* derives utility from each provider *j* according to Equation (2). Each provider proposes a fee  $b_{ij}(s)$ , which I discuss later in Section 5.3. I allow the fee sensitivity parameter  $\alpha_i$  to vary across employers *i* following a lognormal distribution  $\ln \alpha_i \sim \mathcal{N}(\alpha, \sigma)$ . Provider quality is captured by fixed effect  $\nu_j$ . The employer incurs a switching cost  $\kappa_{sw}$  when choosing a non-incumbent provider  $j \neq s$ .  $\mathbb{E}V^E(j, b_{ij}(s))$  represents the employer's value function from choosing *j*, where the expectation is taken over fee transitions to the next period.  $\delta$  is the discount factor, which I assume to be the same for the employer and providers. The employer draws an idiosyncratic preference shock  $\epsilon_{ij}$  for each provider, which follows the Type I extreme value distribution with the scale normalized to 1.

$$\mathcal{U}_{ij}(s) + \epsilon_{ij} = -\alpha_i b_{ij}(s) + \nu_j - \kappa_{sw} \mathbb{1}\{j \neq s\} + \delta \mathbb{E} V^E(j, b_{ij}(s)) + \epsilon_{ij}$$
(2)

The probability that the employer chooses provider j is given by the following.<sup>19</sup>

$$q_{ij}(s) = \frac{e^{\mathcal{U}_{ij}(s)}}{\sum_{k=1}^{J} e^{\mathcal{U}_{ik}(s)}}$$
(3)

The employer's expected utility from the RFP follows the standard inclusive value expression.

$$\mathbb{E}\mathcal{U}_i(s) = \ln \sum_{k=1}^J e^{\mathcal{U}_{ik}(s)} \tag{4}$$

#### 5.2 Stage 1: Employer's Decision to Run an RFP

In an optimal stopping problem (Rust, 1987), the employer runs an RFP when the expected benefit exceeds the RFP cost  $\kappa_{rfp}$ . Since future fee transitions depend on the provider s, I condition on s in addition to fee f. If the employer runs an RFP, she receives utility following Equation (5), which accounts for the expected utility from an RFP  $\mathbb{E}\mathcal{U}_i(s)$  net of the RFP cost  $\kappa_{rfp}$ . If she does not run an RFP, she receives utility following Equation (6), and she remains with her incumbent provider at the current fee f.

I introduce an "attention" parameter  $\rho$ , which governs how the employer trades off her RFP cost against other utility components. This parameter can capture multiple drivers in a somewhat reduced form fashion. The attention can be either behavior or rational. The employer may have imprecise information about fees before receiving fee proposals from providers or follow a predetermined schedule for conducting RFPs every few years. Empirically, allowing a potentially lower level of fee sensitivity when the employer is choosing whether to run an RFP than when the employer is selecting providers at an RFP allows me to better rationalize the patterns in the data. Finally,  $\epsilon_{i,1}$ ,  $\epsilon_{i,0}$  are Type I extreme value shocks, again with the scale normalized to 1. <sup>20</sup>

$$\mathcal{U}_{i1}(s) + \epsilon_{i1} = \rho \mathbb{E} \mathcal{U}_i(s) - \kappa_{rfp} + \epsilon_{i1} \tag{5}$$

$$\mathcal{U}_{i0}(s,f) + \epsilon_{i0} = \rho \Big( \nu_s - \alpha_i f + \delta \mathbb{E} V^E(s,f) \Big) + \epsilon_{i0} \tag{6}$$

The probability of running an RFP is the following:

$$\lambda_i(s, f) = \frac{\exp\left(\mathcal{U}_{i1}(s)\right)}{\exp\left(\mathcal{U}_{i1}(s)\right) + \exp\left(\mathcal{U}_{i0}(s, f)\right)}$$
(7)

 $V^E_i(s,f)$  below is the employer's start of period value function, which again follows the standard

<sup>&</sup>lt;sup>19</sup>Since there is no outside option of not offering 401(k) plans, the employer-specific shifter of quality or costs cancels out when comparing the differences in  $\mathcal{U}_{ij}(s)$  across providers. Employers with arbitrarily different quality or costs but the same  $\alpha_i$  have the same provider choice probabilities.

 $<sup>^{20}</sup>$ I have also considered a scale parameter for the Type I extreme value shocks instead of  $\rho$  and obtained qualitatively similar estimates. I prefer the attention parameter  $\rho$  because it does not attribute a higher degree of randomness in RFP decisions to expected utility.

inclusive value expression.

$$V_i^E(s,f) = \ln\left(\exp\left(\mathcal{U}_{i1}(s)\right) + \exp\left(\mathcal{U}_{i0}(s,f)\right)\right)$$
(8)

#### 5.3 Stage 2.a: Providers' Dynamic Fee Competition

The provider's ex-ante value function at an RFP depends on the probability of winning the RFP and the difference in continuation values between winning and losing. Provider j's production cost is captured by cost fixed effect  $c_j$ .  $V_i^W(j, b_{ij}(s))$  denotes the winner's continuation value, and  $V_{ij}^L(k, b_{ik}(s))$  denotes j's continuation value when k wins instead. Expectations  $\mathbb{E}V_i^W$  and  $\mathbb{E}V_{ij}^L$  are taken over fee transitions.

$$V_{ij}^{P,rfp}(s) = q_{ij}(s) \left( b_{ij}(s) - c_j + \underbrace{\delta \mathbb{E}V_i^W(j, b_{ij}(s))}_{\text{Winner} \text{ continuation value}} \right) + \sum_{k \neq j} q_{ik}(s) \underbrace{\delta \mathbb{E}V_{ij}^L(k, b_{ik}(s))}_{\text{Loser continuation value}}$$
(9)

Provider continuation values as the winner or loser of the RFP depend on the probability  $\lambda_i(s, f)$  that the employer runs RFPs.

$$V_i^W(s,f) = \lambda_i(s,f)V_{is}^{P,rfp}(s) + (1 - \lambda_i(s,f))(f - c_s + \delta \mathbb{E}V_i^W(s,f))$$
$$V_{ij}^L(s,f) = \lambda_i(s,f)V_{ij}^{P,rfp}(s) + (1 - \lambda_i(s,f))\delta \mathbb{E}V_{ij}^L(s,f)$$

If the employer runs an RFP, the previous winner participates as the incumbent  $V_{is}^{P,rfp}(s)$  with a higher probability of winning again due to the switching cost, whereas previous losers participate as non-incumbent providers  $V_{ij}^{P,rfp}(s)$  with  $j \neq s$ . If the employer does not run an RFP, the winner earns fee f net of production costs  $c_s$  and his continuation value, whereas losers only internalize their continuation values.

To gain further intuition, I can solve these continuation values recursively and rewrite Equation (9) as the following

$$V_{ij}^{P,rfp}(s) = q_{ij}(s) \left( \begin{array}{c} \overbrace{\gamma_{ij}(s)}^{\text{expected future}} (b_{ij}(s) - c_j) + \overbrace{\eta_{ij}^{W}(s)}^{\text{Future gains}} \right) + (1 - q_{ij}(s)) \overbrace{\eta_{ij}^{L}(s)}^{\text{Future gains}}$$

 $\gamma_{ij}$  corresponds the expected periods that j provides services to employer i. By setting fees at  $b_{ij}(s)$ , provider j expects to earn  $\gamma_{ij}(s)(b_{ij}(s) - c_j)$  without updating menus. In addition, provider j also expects future gains  $\eta_{ij}^W(s)$ , including the secular decline in j's production costs and the present value of participating in the next RFP as the incumbent. On the other hand, if provider j loses the RFP, he obtains  $\eta_{ij}^L(s)$  which only includes the present value of participating in the next RFP as a non-incumbent provider. Note that  $\gamma_{ij}(s), \eta_{ij}^W(s), \eta_{ij}^L(s)$  depend on the probabilities of RFPs, which in turns are affected by fees set at the current RFP.

Each provider sets his optimal fee  $b_{ij}(s)$  by solving the FOC of  $V_{ij}^{P,rfp}(s)$  with respect to  $b_{ij}(s)$ , as best responses to fees set by other providers  $b_{i,-j}(s)$  and before the employer realizes

her idiosyncratic preferences  $\epsilon_{ij}$  for each provider.

$$b_{ij}(s) = \underbrace{c_j + \frac{1}{\alpha_i(1 - q_{ij}(s))}}_{\text{static cost + markups}} - \underbrace{\frac{\eta_{ij}^W(s) - \eta_{ij}^L(s)}{\gamma_{ij}(s)}}_{\text{future gains per period} > 0} + \underbrace{\frac{\frac{\partial \gamma_{ij}(s)}{\partial b_{ij}}(b_{ij}(s) - c_j) + \frac{\partial \eta_{ij}^W(s)}{\partial b_{ij}}}{\alpha_i(1 - q_{ij}(s))\gamma_{ij}(s)}}_{\text{higher fees raise prob of RFPs } < 0$$
(10)

The first term is the standard cost plus markups in static optimal fees. Because  $q_{is}(s)$  is higher for the incumbent provider s due to switching cost, the incumbent s charges higher markups. Employers with low fee sensitivities  $\alpha_i$  have higher markups and larger markup difference between incumbent and non-incumbent providers.

The second and third terms capture dynamic incentives. The second term is typical in dynamic models with switching costs. It corresponds to higher future gains per period after winning the current RFP. This future gain is competed away at the RFP, prompting providers to reduce their fees  $b_{ij}(s)$ . The third term captures an additional incentive due to endogenous RFP decisions. Since the employer is more likely to run another RFP if her current fee is high, setting a lower fee will reduce the probabilities of RFPs and raise expected profits. I discuss how I solve the FOC in Appendix F.1.

#### 5.4 Stage 3: Fee Transitions and Menu Updates

The transition of the fee state variable between periods depends on the rate of decline in providers' production costs and whether providers update menus outside of RFPs. Without menu updates, providers maintain high fees rather than reducing fees along with the secular decline in costs of providing funds and services. Hence, the fee state variable f increases between periods without menu updates. I allow the rate of increase  $g_s$  to be provider-specific, capturing heterogeneity in production cost decline.<sup>21</sup> I consider  $g_s$  as an exogenous provider attribute.

Whether providers update menus and how much they update outside RFPs depends on the negotiation between the employer and her incumbent provider during annual or quarterly review meetings. For simplicity, I abstract away from directly modeling this negotiation step. I assume all providers follow a constant and exogenous menu update probability  $\phi$ . In practice, the employer can use the threat of RFP to push her provider to update menus. In my model, this mechanism is front-loaded at an RFP, where the provider has an incentive to set a lower fee to reduce the probability that the employer runs RFPs in the future. In addition, I focus on the extensive margin of whether providers update plan menus and abstract away from modeling the magnitude of menu turnover. Whenever providers update menus, I assume that fee growth is set to  $g_s = 0$  for all providers.<sup>22</sup> While my approach is potentially equivalent to a model without menu updates, where fees grow at some average rate between  $g_s$  and zero, allowing for menu updates enables me to use menu turnover data to recover the unobserved probabilities of RFPs, which I discuss later in Appendix G.3.

The expected provider continuation values integrate over the probability of menu update  $\phi$ .

<sup>&</sup>lt;sup>21</sup>This heterogeneity can also include different mutual fund fee strategies. For example, compared to other fund managers, Vanguard manages a smaller set of funds, most of which are fee-competitive. The fee difference between existing funds on plan menus and newly introduced funds is smaller for Vanguard. As a result, Vanguard's  $g_s$  is lower.

 $<sup>^{22}</sup>$ I have estimated a parameter for fee growth when providers update menus and obtained an estimate very close to zero. I have also estimated a version of the model with endogenous menu updates and obtained similar results.

Again, fee f increases at provider-specific rates  $g_s$  without menu updates and stays constant otherwise.

$$\mathbb{E}V_{i}^{W}(s, f) = \phi V_{i}^{W}(s, f) + (1 - \phi)V_{i}^{W}(s, f + g_{s})$$
$$\mathbb{E}V_{ii}^{L}(s, f) = \phi V_{ii}^{L}(s, f) + (1 - \phi)V_{ii}^{L}(s, f + g_{s})$$

The employer's expected value function also integrates over menu update probability  $\phi$ .

$$\mathbb{E}V_i^E(s,f) = \phi V_i^E(s,f) + (1-\phi)V_i^E(s,f+g_s)$$

#### 5.5 Equilibrium

I solve for the Markov Perfect equilibrium, where

- 1. The employer decides whether to run RFPs given her current incumbent provider and fee (s, f), according to the optimal stopping problem in Section 5.2
- 2. At RFPs, the employer chooses a provider to maximize expected utility, according to Section 5.1
- 3. Providers set optimal fees at RFPs as best responses to fees set by other providers following the FOC in Equation (10) of Section 5.3
- 4. The employer and providers have the same beliefs regarding transitions of fees given each other's optimal policies

The equilibrium corresponds to a fixed point of the employer's provider choices, the employer's RFP decisions, and providers' optimal fees conditional on state variables (s, f). In my model, fees set by providers at RFPs enter the employer's and providers' value functions. Additionally, the employer's probabilities of RFPs affect provider value functions. As a result, I solve the employer's RFP decisions and providers' optimal fees jointly.

# 6 Estimation and Identification

My estimation strategy is motivated by the following decomposition of observed plan fees  $fee_{ijt}$  of employer *i* at time *t* whose provider is *j*. Plan fees include asset-weighted mutual fund expense ratios and direct recordkeeping fees as a fraction of total plan assets.

$$fee_{ijt} = \underbrace{markup_{ijt} + \overbrace{c_j}}_{\text{fee state variable } f_{ij}} \underbrace{\tau_{jt}}_{\text{Provider cost trend }} \underbrace{\tau_{jt}}_{\text{Provider cost trend }} \underbrace{X_{it}\beta}_{\text{Observed characteristics }} \underbrace{\Gamma_i + \gamma_{ijt}}_{\text{Unobserved characteristics }}$$
(11)

I decompose observed plan fees into markups and various components of providers' production costs. Provider cost shifter corresponds to the  $c_j$  in my structural model. I assume time trend  $\tau_{jt}$  capture provider-specific reductions in production costs, driving the secular decline in fees. Provider cost shifters correspond to  $c_j$  in my structural model. Lastly, there are heterogeneous characteristics of funds or services across employers. I assume no provider has any competitive advantage of offering specific funds or services, so these characteristics only affect providers' production costs. Some characteristics  $X_{it}$  are observed, such as the asset class composition of

plan menus. Other characteristics are unobserved, including unobserved features of services and some fund selections that could reflect either markups or production costs (e.g. index vs active funds). I use employer shifters  $\Gamma_i$  and a time-varying  $\gamma_{ijt}$  to capture unobserved characteristics.  $\gamma_{ijt}$  can also reflect measurement errors in fees.

My estimation follows two steps. First, I use a hedonic regression to estimate observed plan characteristics  $X_{it}$  and time trend  $\tau_{jt}$ . Second, I estimate structural parameters in the model to recover markups. The fee state variable f in my structural model captures markups and cost shifters  $c_j$ , whereas the residualized fees from the hedonic regression in the first step also include unobserved characteristics. Since I do not have a direct empirical counterpart of the fee state variable f, I estimate my model using indirect inference (Gourieroux and Monfort, 1996).<sup>23</sup>

I also estimate several parameters offline when they directly correspond to certain empirical moments without having to solve the structural models. In Table 2, I summarize which parameters I estimate in each step and the moments in the data I use for identification. I set the discount factor  $\delta = 0.95$  for both employers and providers.

To allow for heterogeneity across employers of different sizes, I break employers into quartiles based on their total plan assets.<sup>24</sup> I control for size quartiles in the hedonic regression and estimate structural parameters separately for each quartile, allowing employers' transaction costs, fee sensitivities, preferences for providers, and providers' production costs to vary across sizes of employers.

Parameters	Moments
Structural Estimation	
Transaction costs: $\kappa_{sw}, \kappa_{rfp}$	Probability of switching and fee autocorrelation
Employer fee sensitivity and attention: $\alpha, \sigma_{\alpha}, \rho$	Fee reductions from switching, sensitivity of switching wrt. fees
Provider propensity to update menu: $\phi$	Probability of menu updates
Small providers net quality: $\nu_{other} - \alpha c_{other}$	Probability of switching from small providers
Hedonic Regression & Offline Large provider net quality: $\nu_j - \alpha c_j$ Employer specific net quality: $\nu_i - \alpha c_i$	Employer FE & choice probability conditional on switching
Offline Fee state variable growth rates: $g_j$	Increase in residualized fees without menu updates

Table 2: Overview of Parameter Estimation

### 6.1 Step 1: Residualizing Fees with Hedonic Regression

I residualize the observed characteristics that I assume correspond to exogenous variation in production costs and do not affect markups. Specifically, I use the following hedonic regression. The dependent variable  $fee_{ijt}$  represents the observed plan fee of

<sup>&</sup>lt;sup>23</sup>I control for provider fixed effects in both steps to account for  $c_i$ .

<sup>&</sup>lt;sup>24</sup>I compute average total plan assets over time for each employer, so employers do not move across size quartiles. I have also tried grouping employers based on the number of plan participants and obtained similar results. I prefer grouping by total assets because plan fees are measured as fractions of total assets.

employer i at time t whose provider is j.

$$fee_{ijt} = \underbrace{X_{it}\beta}_{\text{Observed}} + \underbrace{\Gamma_j + \tau_{jt} + \tau_{s(i)t}}_{\text{Provider FE}} + \underbrace{\Gamma_i + \epsilon_{ijt}}_{\tilde{f}_{ijt}: \text{Residualized fee}}$$
(12)

 $X_{it}$  captures observed employer and plan characteristics, including menu composition and types of services that provider j offers. I also use time trends  $\tau_{jt}, \tau_{s(i)t}$  to residualize provider-specific and employer-size-quartile-specific decline in providers' production costs. Lastly, I use provider fixed effects  $\Gamma_j$  to control for persistent differences in the quality or production costs across providers.  $\Gamma_j$  captures both provider cost shifters  $c_j$ and also difference in markups between providers. I will later include another set of provider fixed effects in the structural estimation step. In Appendix H, I discuss the hedonic regression in greater detail.

Residualized fees  $f_{ijt}$  captures both markups and unobserved characteristics. While I include employer fixed effects  $\Gamma_i$  in the hedonic regression, I also consider  $\Gamma_i$  as part of residualized fees because I cannot reliably estimate many employer fixed effects with a short panel. I then use the variation in residualized fees, provider switches, and menu turnovers to estimate structural parameters.

### 6.2 Step 2: Structural Estimation with Indirect Inference

The indirect inference estimation procedure relies on simple auxiliary models to capture the variation in the data that allows me to identify structural parameters.<sup>25</sup>

I use indirect inference because I do not directly observe the fee state variable f. The closest empirical counterpart is residualized fees  $\tilde{f}_{ijt}$  from the hedonic regression. However, residualized fees  $\tilde{f}_{ijt}$  still include unobserved characteristics, causing an omitted variable bias. For example, an employer may not switch her provider despite high fees because these fees compensate for more valuable services. As a result, I would underestimate the sensitivity of her switching decision with respect to fees.

Estimating a large number of employer fixed effects or allowing for a mixture of latent types is difficult in non-linear structural models with short panels. However, in linear auxiliary models, employer-specific unobserved characteristics can be absorbed by employer fixed effects. These employer fixed effects correspond to employer-specific quality or cost shifters that cancel out when comparing utilities or fees across providers in my structural model. Therefore, I can estimate comparable auxiliary model coefficients using observed data and data simulated by my model.

<sup>&</sup>lt;sup>25</sup>Because I do not observe fees proposed by non-chosen providers at RFPs, I cannot use standard demand inversion as in Berry et al. (1995). I also cannot use analytical solutions of a second-price auction (Allen et al., 2019; Cuesta and Sepúlveda, 2021) due to the dynamic feature of the model where fees enter the continuation values. As a result, I adopt this full-solution approach, where I simulate equilibrium outcomes from my structural model and compare them with their empirical counterparts.

I estimate the coefficients of these auxiliary models using observed data  $\mathcal{B}^{data}$  and also estimate these coefficients using data simulated by my model  $\mathcal{B}^{model}(\Theta)$  at candidate sets of parameters  $\Theta$ . The estimation procedure searches for structural parameters  $\Theta$  to minimize the distance between the two sets of auxiliary model coefficients  $\hat{\Theta} = \arg \min_{\Theta} ||\mathcal{B}^{data} - \mathcal{B}^{model}(\Theta)||_2$ . I introduce seven auxiliary models to identify seven structural parameters: RFP costs  $\kappa_{rfp}$ , switching costs  $\kappa_{sw}$ , average and standard deviation of fee sensitivity  $\alpha, \sigma$ , attention in RFP decision  $\rho$ , providers' menu update probability  $\phi$ , and lastly the net quality for providers in the other category  $\xi_{other} = \nu_{other} - \alpha c_{other}$ . I do not separate quality  $\nu$  from cost c, since only their difference matters after plugging in the expression of optimal fees Equation (10) into employers' utility in Equation (2). Details of the auxiliary models are included in Appendix G.

#### 6.3 Step 2: Identification of Structural Parameters

I focus on the identification of transaction costs  $\kappa_{rfp}$ ,  $\kappa_{sw}$ , fee sensitivity  $\alpha$ , and RFP attention  $\rho$ . Assuming I have perfect data, the level and slope of RFP probabilities identify RFP costs  $\kappa_{rfp}$  and the mixture of fee sensitivity and attention  $\alpha\rho$ , similar to Rust (1987). Differences in choice probabilities between incumbent and non-incumbent providers at RFPs combined with differences in equilibrium fees set by incumbent and non-incumbent providers identify switching costs  $\kappa_{sw}$  and fee sensitivity  $\alpha$ , as standard discrete choice models (Berry et al., 1995). However, I do not observe RFPs or fees set by non-chosen providers. To make progress, I use fee autocorrelation and reduction in fees from switching to approximate the level of RFP probabilities and differences in fees between incumbent and non-incumbent providers at RFPs. I provide the intuition of identification and specify how each auxiliary model identifies corresponding parameters in Appendix G.

**Transaction costs**: The probability that employers switch providers identifies the sum of the two transaction costs  $\kappa_{rfp} + \kappa_{sw}$ . Since switching requires both running an RFP and choosing a non-incumbent provider, higher switching costs and/or RFP costs lead to a lower probability of switching.

Although RFPs are not directly observed in the data, I can infer the probability of RFPs using the autocorrelation of fees. This fee autocorrelation then allows me to separate RFP costs  $\kappa_{rfp}$  from switching costs  $\kappa_{sw}$ . The lower bound of the fee autocorrelation is zero, assuming that providers do not condition on previous fees when setting fees at RFPs. The upper bound can be directly estimated when there is no menu turnover.<sup>26</sup> For employers who do not switch providers, an autocorrelation close to zero suggests frequent RFPs and low RFP costs, whereas an autocorrelation near

 $<sup>^{26}</sup>$ The upper bound is not necessarily equal to one in the data due to changes in fund expense ratios, reallocation across different funds, and measurement errors in recordkeeping fees.

the upper bound suggests infrequent RFPs and high RFP costs. See Appendix G.3 for more details.

Employer fee sensitivity: The sensitivity of switching providers with respect to fees identifies the mixture of average fee sensitivity  $\alpha$  and the attention parameter  $\rho$ . If higher fees strongly predict switching, employers should be attentive in their RFP decisions and/or put more weight on fees rather than transaction costs.

To separately identify  $\alpha$ , I use fee reductions when employers switch providers. In my model, where providers only compete on fees, employers with higher fee sensitivities have smaller fee reductions from switching. Suppose an employer is extremely fee-sensitive, such that providers engage in perfect competition. In this case, her fee reduction from switching providers only captures the accumulated decline in production costs (or equivalently growth in markups) from the previous RFP. Any additional reductions reflect markup differences between incumbent and non-incumbent providers, which depend on the employer's fee sensitivity. When considering switching providers at an RFP, the employer trades off lower markups against switching costs. If the employer puts less emphasis on fees, non-incumbent providers have to set lower markups to incentivize switching. As a result, the employer achieves larger fee reductions when choosing non-incumbent providers.<sup>27</sup>

As robustness checks, I examine how fee reductions from switching vary across employer characteristics. Appendix G.6 shows that fee reductions from switching are smaller for larger employers, employers that contribute more to their plans, employers with higher participation rates, and employers with higher ESG scores. These employers are likely more sophisticated and value employee welfare more. It seems reasonable that they place greater emphasis on plan fees paid by their participants rather than their own transaction costs.

**Random coefficient in employer fee sensitivity**: Since fee reductions from switching identify the average employer fee sensitivity, the variance of fee reductions from switching should capture the variance of employer fee sensitivity  $\sigma^2$ . However, using the variance of fee reductions directly can be problematic, because it may include the variance of changes in measurement errors and unobserved production costs.

To make progress, I assume that, after controlling for provider fixed effects, fee changes from switching in the data consist of two components: strictly negative markup changes (positive markup reductions) related to employer fee sensitivity and zero-mean noise. Any increase in fees from switching providers has to come from noise rather than markups. Thus, the truncated mean of positive fee changes from switching can be used to estimate the variance of noises, allowing me to recover the variance of employer fee sensitivity  $\sigma^2$ .

 $<sup>^{27}</sup>$ Appendix Figure A6 shows that fee reductions from switching under estimated fee sensitivities are much larger compared to fee reductions assuming perfect competition.

**Limitations:** Using the average and variance of fee reductions from switching providers to recover the distribution of fee sensitivity  $\alpha_i$  imposes certain assumptions. Employers with the lowest fee sensitivities rarely conduct RFPs or switch providers. Provider switches in my data primarily reflect behaviors of relatively fee-sensitive employers, and I rely on the lognormal functional form to extrapolate.

In addition, I find variation in the data that corresponds to employer fee sensitivities and transaction costs, which determine markups in my model. However, there may be other sources of markups that I do not capture in my structural estimation. Consequently, I may underestimate both the level and dispersion of markups. As an alternative strategy, I directly use the variance of fees to identify the variance of employer fee sensitivity  $\sigma^2$ , assuming no variation in unobserved characteristics about production costs within each size quartile. Results from this approach likely capture an upper bound of the level and dispersion of markups. Using this approach, I estimate slightly higher levels of markups, and markups explain a greater fraction of the fee dispersion somewhat mechanically. Estimates of the effects of transaction costs and counterfactual results are qualitatively similar.

### 6.4 Parameters Estimated Offline

Other parameters in the model are directly observed or can be estimated offline. First, I assume there are six providers at each RFP, including the five largest providers (Fidelity, Vanguard, Empower, ADP, and Principal Financial Group) and one representative from the other category of smaller providers.<sup>28</sup> Choice probabilities conditional on switching can recover quality net of cost  $\xi_j$  for the five large providers, which I explain in Appendix F.2. Second, I measure fee state variable growth rates  $g_j$  directly in the data, using median changes in residualized fees  $\tilde{f}_{ijt}$  without menu updates to capture how much expense ratios increase relative to the secular trend. In Appendix Table A10, I report conditional choice probabilities and fee state variable growth rates.

# 7 Estimation Results

I present parameter estimates for each employer size quartile in Table 3. To compare magnitudes, I divide transaction cost parameters by average fee sensitivities  $\bar{\alpha}$  so they have the same unit as fees. Both RFP and switching costs in terms of basis points over plan assets decrease with employer size. A large component of these costs consists of fixed costs, which naturally become smaller when divided by larger plan assets. From small employers in the first quartile to large employers in the fourth quartile, RFP costs decrease from 10 to 1 bps, and switching costs decrease from 46 to 2 bps. Measured

<sup>&</sup>lt;sup>28</sup>If the incumbent prior to the RFP is a provider from the other category, I assume there are seven providers to allow for switching between providers in the other category.

in dollars, both transaction costs increase with employer size except between the third and fourth quartile. There can also be some variable components, and these transaction costs become higher for employers with more participants and potentially more complex plan structures. It is reasonable that switching costs are much larger than RFP costs, since changing providers involves handling disruptions to participants' experiences and potentially changes to bundled non-401(k) services.

Employer fee sensitivity  $\alpha_i$  follows a lognormal distribution, and I report the average  $e^{\alpha + \frac{1}{2}\sigma^2}$  and standard deviation  $((e^{\sigma^2} - 1)e^{2\mu + \sigma^2})^{1/2}$ . I estimate that larger employers are more fee-sensitive, especially those in the fourth size quartile. In the first quartile, employers' fee elasticity when choosing providers at RFPs is only 1.6, suggesting a relatively low level of fee sensitivity and potentially high markups. On the other hand, employers in the fourth quartile have an elasticity of 50, suggesting that providers engage in near-perfect competition for large employers. The standard deviation of  $\alpha_i$  increases along with the mean, while the coefficient of variation (mean divided by standard deviation) is around one for all four quartiles, suggesting similar dispersion in fee sensitivities across employer size.

When deciding whether to conduct RFPs, employers tend to be less attentive with  $\rho < 1$ . The attention parameter is lower for larger employers, but the multiple of  $\rho \times \alpha$  is still higher for larger employers. Outside of RFPs, the probability that providers update menus ranges from 50% for small employers to almost 70% for large employers. Larger employers are in a stronger position to demand more frequent menu updates.

The net equality  $\xi$  should be interpreted on a relative basis across providers.<sup>29</sup> Among large providers, Fidelity has the highest net quality, consistent with its large market share shown in Figure 1a. Vanguard has relatively higher  $\xi$  among large employers and relatively lower  $\xi$  among small employers. Figure 1a shows that Vanguard has a larger market share in plan assets than in number of plans, consistent with higher net quality among large employers. ADP has a higher  $\xi$  on a relative basis for small employers in the first quartile. ADP bundles recordkeeping with HR payroll, which is more appealing for small employers.

### 7.1 Goodness of Fit

In Appendix Table A11, I assess the goodness of fit by comparing auxiliary model coefficients estimated using data simulated by the model and coefficients estimated using observed data. The standard errors of coefficients using observed data are shown in parentheses and indicate whether the differences in coefficients are statistically signif-

<sup>&</sup>lt;sup>29</sup>Technically, there are two  $\xi_{other}$  for other providers in my model. In Table 3,  $\xi_{other}$  corresponds to the net quality of one incumbent provider from the other category. At an RFP, the employer may invite multiple small providers from the other category to participate at an RFP. I only model one representative provider of the other category, and the  $\xi_{Other}$  of this representative also includes a variety effect. I normalize  $\xi_{Other} = 0$ . Without the variety effect,  $\xi_{other}$  is negative.

Employer quartiles		1st		2	2nd		3rd		4th	
		Coef	SE	Coef	SE	Coef	SE	Coef	SE	
RFP cost bps \$ thousand	$\kappa_{rfp}/\bar{\alpha}$	$10.2 \\ 7$	(0.2)	$\begin{array}{c} 6.4 \\ 10 \end{array}$	(0.4)	4.8 17	(0.1)	0.8 18	(0.5)	
Switching cost bps \$ thousand	$\kappa_{sw}/\bar{\alpha}$	$\begin{array}{c} 46.0\\ 30 \end{array}$	(0.8)	$\begin{array}{c} 24.2\\ 38 \end{array}$	(1.3)	$15.9 \\ 57$	(0.4)	$\begin{array}{c} 1.7\\ 39 \end{array}$	(1.0)	
Fee sensitivity: mean Coefficient Elasticity	$\bar{\alpha}$	$\begin{array}{c} 0.11 \\ 1.6 \end{array}$	(0.00)	$\begin{array}{c} 0.24\\ 3.4 \end{array}$	(0.01)	$0.41 \\ 5.2$	(0.01)	$\begin{array}{c} 4.66\\ 64.4\end{array}$	(2.46)	
Fee sensitivity: st.dev. Coefficient Elasticity	$\sigma_{lpha}$	$\begin{array}{c} 0.14 \\ 1.6 \end{array}$	(0.00)	$0.32 \\ 3.5$	(0.02)	$0.53 \\ 5.1$	(0.02)	$15.00 \\ 58.9$	(10.15)	
RFP attention Coefficient	ρ	0.78	(0.01)	0.59	(0.03)	0.47	(0.01)	0.26	(0.15)	
Pr menu update	$\phi$	0.50	(0.00)	0.62	(0.00)	0.68	(0.00)	0.63	(0.00)	
Net quality of other providers Relative to other providers	$\xi/\bar{\alpha}$	-14.6	(0.3)	-13.4	(0.7)	-10.5	(0.2)	-1.1	(0.6)	
ADP		-0.8		2.0		1.0		0.2		
Principal Fin.		-2.4		0.9		1.7		0.2		
Empower		1.8		3.7		3.3		0.4		
Fidelity		1.6		3.4		3.9		0.4		
Vanguard		-1.7		2.0		2.9		0.4		

 Table 3: Structural Parameter Estimates

Table reports parameter estimates. Standard errors are computed using the asymptotic variance formula, where the variance-covariance matrix of auxiliary model coefficients is computed using the bootstrap method over 250 iterations. In each iteration, I sample employers with replacement and estimate auxiliary model coefficients.

icant. Coefficients are generally not distinguishable at the 5% significance level (two-tailed). The only coefficient that I fit slightly poorly is the probability of switching from large providers.

In Appendix Table A10, I show that the probabilities that employers run RFPs range from 15% to 30% based on model simulation. This is broadly consistent with the industry common practice of conducting an RFP every three to five years. When some employers claim they run an RFP every three years, they could refer to formal RFPs, informal Requests for Information (RFIs), or benchmarking. RFP decisions in my model should correspond to formal RFPs, and it seems reasonable that the RFP frequency from my model is close to the lower bound of the industry practice. Large firms run RFPs less frequently potentially because they adopt more formal RFP procedures that require coordination across multiple departments.

### 7.2 Markups

Having estimated my model, I can recover markups. At RFPs, markups can be computed from providers' FOC conditions Equation (10). After RFPs, markups evolve according to growth rates  $g_j$  until the next RFPs. I report markups separately for each size quartile in Figure 6 and show aggregated results in Table 4.

Average level of markups: In the first column of Table 4, I show that markups are on average 18% of fees, among all employers in my sample. Around 12 bps out of the 65 bps average plan fees as of 2019 are markups, while the rest are providers' production costs.

In Figure 6a, I separate plan fees into production costs and markups for each size quartile. Markups are 18 bps or 23% for small employers, decreasing to 6 bps or 13% for large employers. Lower markups for larger employers are consistent with larger employers having higher fee sensitivity. Between the first and fourth size quartiles, the difference in plan fees is 20 bps. Roughly one-third of this difference is due to markups, and two-thirds comes from production costs, which account for economies of scale of large employers.

Sources of markups: I decompose sources of markups in Figure 6b. First, nonincumbent providers set lower markups at RFPs to incentivize switching. Providers charge 4 bps in markups to small employers in the first quartile and negative markups to large employers. Forward-looking providers are willing to charge negative markups temporarily, anticipating higher future markups taking advantage of employers' transaction costs. Second, incumbent providers set relatively higher markups at RFPs to exploit switching costs. Dark gray bars show that markups from switching costs range from over 10 bps for small employers to 5 bps for large employers. Lastly, because employers do not run RFPs every year, markups increase over time relative to the declining production costs between RFPs. Light gray bars show that the last source of markups between RFPs is between 2 and 4 bps across employers with different sizes.

In Table 4, I show that employers (across all four size quartiles) can reduce fees by 11 bps on average if they were to run RFPs and switch providers incurring both transaction costs. In aggregate, employers can save participants \$1 billion in annual plan fees. I compute this difference by comparing average markups from the model with average markups non-incumbent providers charge at RFPs. This difference may not equal the observed fee reductions in the data due to the selection of employers who endogenously run RFPs and switch.

These results highlight that transaction costs play significant roles in preventing 401(k) participants from accessing low-cost investment options and services. However, whether fee reductions of such magnitude are feasible also depend on providers' equilibrium response, which I discuss later with counterfactual exercises in Section 8.

Distribution of markups: In Figure 6c, I plot the distribution of markups simu-

lated by my model across different states. Within each size quartile, markups vary across employers due to differences in fee sensitivity and within employers across state (s, f)because employers incur transaction costs at different points in time. The solid black line indicates that markups center around zero for large employers. Since the market for large employers is nearly perfectly competitive, providers shuffle markups across time but earn a small amount of markups on average. As a comparison, markups are higher on average for smaller employers and also more dispersed. Due to the wide dispersion, there are substantial overlaps between the four distributions, suggesting that some large employers can have higher markups than small employers at least temporarily.





Figure illustrates model generated markups for employers in each size quartile. Panel (a) separates plan fees into markups and production costs. Panel (b) decompose markups into components set by providers at RFPs and growth in markups between RFPs. Panel (c) plots distribution of markups across states. Panel (d) shows the decomposition of fee dispersion.

**Decomposition of fee dispersion**: I decompose the dispersion in 401(k) plan fees into different sources of markups and production costs. Using the law of total variance, I separate the variance of markups into the between-employer variation due to different fee sensitivity and the within-employer variation because employers incur transaction costs at different points in time.

$$Var(f_i) = \underbrace{Var(\bar{f_i})}_{\text{between: fee sensitivity}} + \underbrace{E_i[Var(f_i|i)]}_{\text{within: transaction costs}}$$

Further separating markups into these two components is useful. The variance driven by fee sensitivity is more likely to capture how well employers manage their plans. On the other hand, some variance of markups is natural in this type of negotiated-price market with transaction costs. Even homogeneous employers may have different fees in the cross-section, and this component of variance in markups is not informative of employers' behaviors. <sup>30</sup>

I attribute the variation in plan fees not explained by markups to production costs. The variance of production costs includes the component explained by observed plan characteristics in the hedonic regression and the remaining fee variation not explained by either the hedonic regression or the structural model.

In the first column of Table 4, I report the variance decomposition by pooling employers across size quartiles. The variation in markups driven by fee sensitivity explains 38% of the overall fee dispersion. Within-employer variation in markups due to transaction costs explains 13%. The remaining 48% is attributed to production costs.

Across different size quartiles, fee sensitivity explains a greater fraction of the fee dispersion for small employers in the first quartile and large employers in the fourth quartile. There may be greater heterogeneity in sophistication across employers within these two quartiles at the boundaries of the overall size distribution. In Appendix Figure A2, I show a wider dispersion of employer sizes within the first and fourth quartiles.

This decomposition can shed light on recent excessive fee lawsuits. While differences in fee sensitivity play an important role, a substantial fraction of the fee dispersion is attributed to employers' transaction costs when selecting and switching providers as well as providers' production costs when providing heterogeneous funds and services. As a result, 401(k) plan fees can appear high in a cross-sectional comparison, even though employers plausibly fulfill their fiduciary duties.

# 7.3 Robustness

I estimate my model under two alternative assumptions. First, my overall estimation strategy is to find the variation in the data that plausibly captures markups, leaving the residual as production costs. Using this approach, I may underestimate the level and dispersion of markups. To address this concern, I make an alternative assumption that there are no unobserved production costs. Under this assumption, I estimate the random coefficient for fee sensitivity by matching the variance of markups simulated from my structural model to the variance of residualized fees, instead of using the

 $<sup>^{30}</sup>$ I am able to make this decomposition because I allow employers within a size quartile to have different fee sensitivities but face the same transaction costs. In practice, both can vary across employers. In my model, my transaction cost parameters can be viewed as some common costs of doing business across employers. On top of this common friction, an employer who internalizes higher transaction costs behaves similarly to an employer who is less fee-sensitive.

	Baseline	No unobserved production costs	Myopic
Markups (bps)			
Average	11	16	14
When employers switch providers	0	5	3
Fee reductions if employers switch providers			
Average (bps)	11.3	10.7	10.5
Aggregate (\$ billion)	1.04	1.03	0.71
% level of fee			
Markups	18	24	21
Production costs	82	76	79
% variance of fee			
Markups: fee sensitivity	38	63	28
Markups: transaction costs	13	13	14
Production costs	48	23	58

Table 4: Decomposition of Plan Fees and Fee Dispersion

Table reports the decomposition of the level and the dispersion of plan fees into markups and production costs. The first column corresponds to estimates from the baseline specification. The second column corresponds to the alternative assumption without unobserved production costs. The last column assumes providers and employers are myopic.

variance of fee changes from switching providers. See Appendix G.7 and Appendix G.8 for a comparison of auxiliary models under these two assumptions.

In Appendix Table A12, I report estimates under the assumption of no unobserved production costs. I estimate similar parameters for large employers in the fourth quartile. Even based on my baseline assumptions, I estimate a very small variation in unobserved production costs. It seems plausible that providers incur similar production costs once employers achieve certain economies of scale. For example, providers may always need to designate a dedicated staff member to work with a large employer, despite different plan features or services offered.

For smaller employers in the bottom three quartiles, I estimate that fee sensitivities are around 50% lower. Because the level of markups is a decreasing and convex function of fee sensitivity based on Equation (10), a lower fee sensitivity leads to a higher level of markups and also a larger variance of markups.

In the second column of Table 4, I show that markups make up a slightly higher percentage of plan fees (24% instead of 18% under the baseline assumption). The variation in markups due to fee sensitivity also explains a greater fraction of fee dispersion (63% instead of 38%). Even with potential upper bounds of markups, there is still a meaningful amount of variation in plan fees not explained by employers' sophistication or negotiation effort.

Second, I assume employers and providers are myopic with discount factors equal to zero. Myopic providers do not consider future profits when they set markups at current RFPs. Myopic employers trade off transaction costs against fee savings in the current period only. In Appendix Table A14, I report estimates under the myopic assumption. While I can also rationalize the data (see Appendix Table A15 for a comparison of goodness of fit), I obtain different parameter estimates. Since myopic employers only focus on fee differences in the current period, rationalizing infrequent provider switches requires lower transaction costs. I estimate that employers in the first quartile face 2 bps of RFP costs and 13 bps of switching costs. In dollar terms, the costs of running RFPs and switching providers are only \$1,000 and \$9,000, which seem unreasonably small. For large employers in the fourth quartile, the costs of running RFPs are only \$5,000, which seems negligible for employers with nearly 3,000 employees on average. When comparing the magnitudes of transaction costs, it seems more sensible to assume that employers and providers are forward-looking.

With the myopic assumption, I still estimate a similar level of markups as a percentage of plan fees, according to the last column of Table 4. The variation in markups driven by fee sensitivity explains an even smaller fraction of the fee dispersion (28% instead of 38%).

# 8 Counterfactual

Using model estimates, I simulate counterfactual policies aimed at mitigating various drivers of plan fees. To evaluate counterfactual outcomes, I focus on plan fees paid by participants and transaction costs incurred by employers. If employers have to take on substantially higher transaction costs, these policies are somewhat unpractical even if they can reduce plan fees for participants.<sup>31</sup> For fees and transaction costs, I report both averages across employers and also the aggregated dollar amount which weight each employer by total plan assets. I report outcomes under the stationary distribution across states and abstract away from the across-time variation. Since my estimation sample is a subset of all plans, the overall impact of these counterfactual policies could three to five times larger according to Appendix Table A1.

I report outcomes under the status quo in the first row of Table 5. Pooling employers across size quartiles, the average plan fee is 65 bps or \$5.8 billion annually. Employers face 6 bps or \$180 million transaction costs per year, where transaction costs are calculated by summing the probability of RFP times RFP costs and the probability of switching times switching costs. In Figure 7, the black square marker captures plan fees (y-axis) and transaction costs (x-axis) under the status quo.

In Appendix Figure A9 and Appendix Table A16, I also report counterfactual outcomes under the alternative model assumptions of no unobserved characteristics related

 $<sup>^{31}</sup>$ I do not compute consumer surplus, which can be hard to interpret in my setting. First, without an outside option of not offering 401(k) plans, the scale of consumer surplus is not identified. Next, I assume employer preferences for provider quality represent their participants' preferences. In practice, employers may request providers to offer certain types of funds or services that their participants do not value.

to production costs within each employer size quartile. Results are qualitatively similar to those using baseline estimates.

	Plan fees		Transa	action costs
Counterfactual	(bps)	(\$ billion)	(bps)	(\$ billion)
Status quo	64.5	5.80	5.9	0.18
Employers ignore transaction costs	64.3	5.41	58.7	2.89
RFP/switch that minimize fees	60.6	5.32	18.7	1.21
High RFP attention $(\rho = 0.95)$	61.1	5.36	5.5	0.20
Consolidate plans, production cost only	53.8	5.52	5.9	0.18
Consolidate plans	48.7	5.40	0.6	0.07
High fee sensitivity $(90 \text{th pct})$	54.3	5.02	11.0	0.34

Table 5: Plan Fees and Transaction Costs under Counterfactual Policies

Table reports outcomes under the status quo and counterfactual policies. I report both average plan fees or transaction costs measured in basis points across employers and also the aggregated dollar amount of annual plan fees or transaction costs.

Figure 7: Plan Fees and Transaction Costs under Counterfactual Policies



Figure plots counterfactual plan fees across different transaction costs. The y-axis corresponds to average plan fees, and the x-axis corresponds to employers' transaction costs.

### 8.1 More Frequent RFPs and Provider Switches

I consider counterfactual scenarios where employers behave as if transaction costs are lower, such that they run RFPs or switch providers more often. I simulate model equilibrium by reducing RFP and switching cost parameters but measure outcomes using estimated transaction costs under the status quo. Hence, I am effectively subsidizing RFPs and provider switches, rather than considering IT advancements that streamline the RFP and switching process.

Following the literature on switching costs since Klemperer (1987a,b), transaction costs have two opposing effects on fees when providers engage in an "invest-harvest"

pricing strategy. On one hand, providers can exploit transaction costs to "harvest", by charging higher markups at RFPs and delaying menu updates until the next RFPs. My estimates suggest that providers earn 11 bps higher fees on average relative to fees set when incentivizing employers to switch at RFPs. On the other hand, forward-looking providers anticipate higher future markups and have incentives to "invest" by setting low or even negative markups initially. Once employers ignore transaction costs, while employers can no longer exploit transaction costs to "harvest" markups, they also have no incentives to "invest". Combining these two effects, the overall impact on average plan fees is theoretically ambiguous.

Gray scatters in Figure 7 correspond to counterfactuals where I reduce RFP and switching cost parameters by different magnitudes. Moving from left to right, employers run RFP and switch more frequently and incur higher transaction costs as a result. The U-shape curve illustrates the two opposing effects on fees. More RFP and switching from the status quo initially decrease plan fees. But as employers run RFP or switch much more frequently, plan fees can actually become higher.

The second row of Table 5 corresponds to outcomes when employers behave as if there were no RFP or switching costs. In this "Ignore transaction costs" counterfactual, average plan fees remain largely unchanged from the status quo, while aggregate fees decrease by \$400 million. This discrepancy is because the "invest" and "harvest" incentives have different relative weights between small and large employers.<sup>32</sup> In fact, I show in Appendix Table A16 that the average plan fees can be even higher than the status quo when employers ignore transaction costs under the alternative modeling assumption of no unobserved production costs.

As a comparison, the "RFP/switch that minimize fees" counterfactual corresponds to the bottom of the U-shape curve of gray scatters in Figure 7, where I find transaction cost parameters that lead to the lowest plan fees. In this counterfactual, employers run RFP and switch providers more frequently, incurring three times the transaction costs compared to the status quo. As a result, providers cannot "harvest" as much as under the status quo but still retain some incentives to "invest" that lead to overall reductions in fees. Plan fees decrease by 4 bps on average or around \$500 million in aggregate relative to the status quo.

My findings that fees are minimized at modest transaction costs are consistent with simulation and theoretical results in Dubé et al. (2009) and Cabral (2016). Intuitively, in a static problem, if providers reduce markups from the static optimal level, they increase their probability of being chosen by the employer at the RFP, but their markups are lower conditional on winning. However, with dynamic considerations, providers consider

 $<sup>^{32}</sup>$ One reason providers have lower incentives to "invest" in large employers is because I estimate that they run RFPs less frequently, shown in Appendix Table A10. The incentive to invest is stronger when providers can recoup their investment sooner. With less frequent RFPs, the present value of charging higher incumbent markups in the next RFP is smaller.

the present value of markups over multiple future periods and are willing to reduce markups below the static optimal level. As a result, some modest transaction costs can be beneficial to incentivize competition among providers.

In recent excessive fee lawsuits, some participants accuse their employers of not running RFPs regularly or staying with the same providers for a long period of time.<sup>33</sup> My results suggest that more frequent RFPs or switching may not necessarily achieve as large reductions in plan fees as one anticipates due to providers' dynamic fee-setting incentive.

# 8.2 **RFP** Attention

In my setting, providers consider the following three components of "future gains" when they set fees at RFPs: expected higher markups from competing in the next RFP as the incumbent, increasing markups relative to the declining production costs until the next RFP, and reducing the probability that employers run another RFP. While the first incumbent markup component is common in existing models of dynamic competition with switching costs and the second component can be addressed by defining each period of the model to be multiple years between RFPs, the last component is novel in my model. Because I endogenize employers' RFP decisions, providers are incentivized to set lower fees to reduce the probability of future RFPs and extend the expected duration of providing services to their existing employers.

To illustrate this incentive, I increase the RFP attention parameter  $\rho$  up to 0.95, so employers behave almost as fee-sensitive when determining whether to run RFPs as when they choose providers at RFPs. With higher RFP attention, providers would reduce fees further to prevent future RFPs. The dashed gray line with circle markers in Figure 7 illustrates that plan fees become lower, whereas employers incur similar transaction costs as the status quo. In the fourth row of Table 5, I show that when employers become more attentive in their RFP decisions, plan fees reduce by 3 bps on average, saving participants \$440 million annually.

This counterfactual suggests that it can be more efficient for employers to run RFPs conditional on the levels of plan fees rather than following certain fixed schedules. In practice, employers could use lower cost benchmarking or Request for Information to stay informed of whether their fees are reasonable, which allows them to use formal RFP as a threat when they meet with their providers during periodic review meetings.

While I do not specify the sources of RFP inattention, information asymmetry could be one potential explanation. Then, requiring more information disclosure can benefit employers and their employee participants. Indeed, similar arguments motivated DOL Rule 408(b)(2) implemented in 2012, which required providers to disclose estimated di-

 $<sup>^{33}</sup>$ These allegations are typically not central to the cases but serve as suggestive evidence of employer negligence. See Singh v. Deloitte LLP for a recent example.

rect and indirect compensation to employer sponsors. Such disclosure made it easier for employers to determine whether the compensation is reasonable and whether contracts with providers involve conflicts of interest that may affect the service providers' performance to employee participants (Badoer et al., 2020). In the discussion of an alternative regulation to direct mandate at employers, Federal Register (2010) acknowledged information asymmetry which can prevent employers from assessing the reasonableness of compensation and argued that "A mandate directed solely at fiduciaries...would merely create a brighter line of obligation for the fiduciary without empowering him to satisfy that obligation; perpetuate the information asymmetry, therefore not correcting the market failure".

One caveat is that I do not consider costly monitoring or other principal-agent frictions in this counterfactual. It can be less costly for employee participants to monitor their employers if they run RFPs every few years.

# 8.3 Consolidating Plans

Next, I consider consolidating 401(k) plans of small employers, motivated by regulations around multiple employer plans (MEPs) in the SECURE Act of 2019. My estimation suggests that having many small employers with limited scale and a lack of sophistication offer 401(k) plans to their respective employees is potentially inefficient. In addition to higher production costs and markups, it is also wasteful for many employers to incur their own transaction costs.

To quantify the benefits of scale, the solid gray line with triangle markers in Figure 7 illustrates plan fees and transaction costs as I incrementally group employers in smaller size quartiles such that they are as large as employers in the second, third, and fourth quartiles. To show these counterfactual outcomes, I compute different weighted averages of status quo results across size quartiles, without recomputing model equilibrium. Consolidation can save both plan fees for participants and transaction costs for employers.

The fifth and sixth rows of Table 5 display outcomes when I consolidate employers in the bottom three quartiles. In the third row, I only account for the effect of reducing production costs, assuming small employers still maintain their fee sensitivity and incur their own transaction costs. Production costs alone can reduce fees by over 10 bps and save participants \$300 million annually. In the fourth row, I also assume employers adopt the fee sensitivity and transaction costs of large employers. As small employers consolidate together their 401(k) plans, they potentially have the resources to hire professionals who are more sophisticated and put more effort into negotiating with providers. Plans fees decrease by 16 bps on average or \$400 million in aggregate. Bhattacharya and Illanes (2022) also find that plan consolidation has larger effects with modifications to employers' preferences. Compared to their results, I estimate stronger
economies of scale, and plan consolidation can achieve larger benefits through reducing production costs.

In addition to saving fees for participants, plan consolidation also reduces transaction costs by 5 bps on average or \$110 million. Since menu updates are important to reduce fees, it would be more efficient to incur necessary transaction costs once and spread the benefit of fee savings to as many participants as possible. I do not discuss the extensive margin of whether employers sponsor 401(k) plans. One reason that prevent employers from doing so can be administrative burden such as transaction costs. Hence, plan consolidation may also incentivize more employers to offer 401(k) plans to their employees.

This counterfactual suggests that consolidating small 401(k) plans could achieve substantial savings in both plan fees and transaction costs. The SECURE Act of 2019 made multiple employer plans (MEPs) less restrictive and more attractive for small employers. Previously, employers have to share a common bond (e.g. in the same industry) and the entire MEP could be penalized if one employer within the group violates its fiduciary duty (the "bad apple" rule).

One caveat is that a multiple employer plan could come with a lower fee sensitivity, higher production costs, or higher transaction costs than a large employer with comparable size due to the complexity of its organization structure (Chen and Munnell, 2024). To the extent that a third-party organization is hired to oversee multiple employer plans, there is an additional source of friction where this third party could extract rent from plan consolidation. Without considering these cases, I likely overestimate the benefits of consolidation. Furthermore, my model does not capture potential supply-side effects. Some providers may rely on small employers to stay profitable. If plan consolidations lead to provider exits, the remaining providers may have higher marker power and can charge higher markups than my estimates.

#### 8.4 Increasing Employer Fee Sensitivity

Lastly, to assess the impact of employer fee sensitivity, I simulate a counterfactual where employers become more fee-sensitive. In Figure 7, the dashed black line with square markers displays outcomes when I set the fee sensitivity parameter  $\alpha_i$  for all employers within each size quartile equal to various percentiles based on estimated distributions. From top left to bottom right, square markers correspond to  $\alpha_i$  at the 50th, 75th, 90th, and 95th percentiles. Although higher fee sensitivity can reduce markups, employers incur more transaction costs as they are more willing to conduct RFPs and switch providers to reduce fees. The last row of Table 5 reports outcomes where fee sensitivity is set at the 90th percentile. Plan fees decrease by 10 bps on average or almost \$800 million annually. Meanwhile, employers take on 5 bps higher transaction costs or \$160 million. If recent excessive fee lawsuits push employers to become more fee-sensitive, participants may benefit from reduced fees without substantially higher transaction costs for employers. However, I caution against over-interpreting the results of this counterfactual, as changing employer fee sensitivity can be difficult in practice.

### 8.5 Counterfactual Discussion

In terms of why the average 401(k) plan fee as of 2019 is as high as 65 bps, a common theme across these counterfactual analyses speaks to the decentralized nature of retirement plans sponsored by individual employers. In particular, small employers can be unsophisticated when negotiating fees with their providers and do not benefit from economies of scale. In a negotiated price market, employers also duplicate transaction costs and may be uninformative about fees when they decide whether to run RFPs.

In Appendix C, I compare 401(k) plan fees with government-sponsored defined contribution plans. Federal Thrift Saving Plans (TSP) for federal and civil services employees only charge participants 5.7 to 9 bps. TSP is the largest defined contribution plan in the world with over \$800 billion asset. Its fee structure and plan menus can also be easily found online. Hence, there can be substantial reduction in 401(k) plan fees through economies of scale and transparency.

Employer-sponsored plans do have certain advantages. Employers are better informed about the risk preferences of their participants and can be in a better position to offer investment advice. I cannot measure the utility participants derive from these benefits in this paper.

# 9 Conclusion

In this paper, I demonstrate that employers' transaction costs when choosing and switching providers can affect 401(k) plan fees. These transaction costs are also important to understand what determines the level and the dispersion of 401(k) plan fees.

Despite the introduction of low-fee mutual funds, I show that providers can exploit employers' transaction costs and delay the inclusion of these funds in 401(k) menus, effectively earning higher markups over time relative to the secular decline in providers' production costs. The mechanism by which providers maintain sticky prices relative to declining costs to sustain markups extends beyond 401(k) plans. A similar dynamic can be observed in fixed-rate mortgages during periods of declining interest rates. Hence, my empirical results and structural model can be applied more broadly.

The descriptive evidence motivates my dynamic model to capture how employers choose providers and how providers compete. Estimating my model with fees and provider choice data, I recover employers' transaction costs, the markup components of plan fees, and the production costs providers incur when providing 401(k) services. I find that transaction costs measured as basis points over plan assets decrease in the size of employers' plans, reflecting large fixed-cost components. Similarly, both markups and production costs decrease in employer size, suggesting that larger employers are more sophisticated and benefit from economies of scale.

I use model estimates to quantify whether counterfactual policies can effectively reduce plan fees. I show that consolidating plans could achieve substantial fee savings through both reducing production costs thanks to economies of scale and also higher fee sensitivity of larger employers. The effects of reducing transaction costs on fees are more complex due to providers' dynamic incentives. Subsidizing employers to run RFPs or switch providers more frequently does not necessarily reduce fees because providers have smaller incentives to use low markups to incentivize employers to switch. As a comparison, more fee-sensitive RFPs can lead to lower fees while employers do not incur additional transaction costs. Overall, my counterfactual speaks to the potential inefficiency of retirement plans sponsored by individual employers, especially small employers.

# References

- ABOWD, J. M., F. KRAMARZ, AND D. N. MARGOLIS (1999): "High Wage Workers and High Wage Firms," *Econometrica*, 67, 251–333.
- ADAMS, N. (2020): "Reader Poll: Has Litigation Led to Change?" National Association of Plan Advisors.
- (2022): "Reader Radar: Revisiting and Revising RFPs," National Association of Plan Advisors.
- ALLEN, J., R. CLARK, AND J.-F. HOUDE (2014): "The Effect of Mergers in Search Markets: Evidence From the Canadian Mortgage Industry," *American Economic Review*, 104, 3365–96.
- (2019): "Search Frictions and Market Power in Negotiated-Price Markets," Journal of Political Economy, 127, 1550–1598.
- ALLEN, J. AND S. LI (2025): "Dynamic Competition in Negotiated Price Markets," Journal of Finance, 80, 561–614.
- ARCIDIACONO, P. AND R. A. MILLER (2011): "Conditional Choice Probability Estimation of Dynamic Discrete Choice Models With Unobserved Heterogeneity," *Econometrica*, 79, 1823– 1867.
- ARONOWITZ, D. (2022a): "Debunking Recordkeeping Fee Theories in Excessive Fee Cases," Euclid Fiduciary Whitepaper.

— (2022b): "A Deep Dive into the Hy-Vee Excessive Fee Case," *Euclid Fiduciary*.

——— (2023): "The Key Fiduciary Liability Storylines of 2022," Euclid Fiduciary.

- BACHAS, N. (2018): "The Impact of Risk-Based Pricing and Refinancing on the Student Loan Market," Tech. rep.
- BADOER, D. C., C. P. COSTELLO, AND C. M. JAMES (2020): "I Can See Clearly Now: The Impact of Disclosure Requirements on 401(k) Fees," *Journal of Financial Economics*, 136, 471–489.
- BEGENAU, J. AND E. SIRIWARDANE (2024): "Fee Variation in Private Equity," Journal of Finance, 79, 1199–1247.
- BENARTZI, S. AND R. THALER (2007): "Heuristics and Biases in Retirement Savings Behavior," Journal Of Economic Perspectives, 21, 81–104.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): "Automobile Prices In Market Equilibrium," econometrica, 63, 841–90.
- BESHEARS, J., J. J. CHOI, D. LAIBSON, AND B. C. MADRIAN (2009): The Importance of Default Options for Retirement Saving Outcomes: Evidence from the United States, University of Chicago Press.
- BHATTACHARYA, V. AND G. ILLANES (2022): "The Design of Defined Contribution Plans," Tech. rep.
- BHATTACHARYA, V., G. ILLANES, AND M. PADI (2019): "Fiduciary Duty and the Market for Financial Advice," Tech. rep.
- BRANCACCIO, G., D. LI, AND N. SCH"URHOFF (2017): "Learning by Trading: The Case of the US Market for Municipal Bonds," *Unpublished Paper. Princeton University*.

- CABRAL, L. (2016): "Dynamic Pricing in Customer Markets with Switching Costs," *Review of Economic Dynamics*, 20, 43–62.
- CARROLL, G. D., J. J. CHOI, D. LAIBSON, B. C. MADRIAN, AND A. METRICK (2009): "Optimal Defaults and Active Decisions," *Quarterly Journal of Economics*, 124, 1639–1674.
- CHALMERS, J., O. S. MITCHELL, J. REUTER, AND M. ZHONG (2021): "Auto-Enrollment Retirement Plans for the People: Choices and Outcomes in Oregonsaves," Tech. rep.
- CHALMERS, J. AND J. REUTER (2020): "Is conflicted Investment Advice Better Than No Advice?" Journal of Financial Economics, 138, 366–387.
- CHEN, A. AND A. MUNNELL (2024): "Will Multiple Employer Plans Help Close the Coverage Gap?" Center for Retirement Research at Boston College.
- CHOI, J. J. (2015): "Contributions to Defined Contribution Pension Plans," Annual Review of Financial Economics, 7, 161–178.
- CHOI, J. J., D. LAIBSON, B. C. MADRIAN, AND A. METRICK (2002): "Defined Contribution Pensions: Plan Rules, Participant Choices, and the Path of Least Resistance," *Tax Policy* and the Economy, 16, 67–113.
- CHOI, J. J., D. LAIBSON, B. C. MADRIAN, A. METRICK, AND J. M. POTERBA (2007): For Better or for Worse: Default Effects and 401(k) Savings Behavior, University of Chicago Press, 81–126.
- CLARK, R., J.-F. HOUDE, AND J. KASTL (2021): "The Industrial Organization of Financial Markets," in *Handbook of Industrial Organization*, Elsevier, vol. 5, 427–520.
- CUESTA, J. I. AND A. SEPÚLVEDA (2021): "Price Regulation in Credit Markets: A Trade-Off Between Consumer Protection and Credit Access," *Working Paper*.

DELOITTE (2013): "Inside the Structure of Defined Contribution/401(k) Plan Fees,".

—— (2015): "Annual Defined Contribution Benchmarking Survey. Ese of Use Drives Engagement in Saving for Retirement," .

— (2019): "The Retirement Landscape Has Changed – Are Plan Sponsors Ready? 2019 Defined Contribution Benchmarking Survey Report," .

- DOELLMAN, T. W. AND S. H. SARDARLI (2016): "Investment Fees, Net Returns, and Conflict of Interest in 401(k) Plans," *Journal of Financial Research*, 39, 5–33.
- DOL (2019): "A Look At 401(K) Plan Fees," U.S. Department of Labor, Employee Benefits Security Administration (EBSA).
- DUBÉ, J.-P., G. J. HITSCH, AND P. E. ROSSI (2009): "Do Switching Costs Make Markets Less Competitive?" Journal of Marketing Research, 46, 435–445.
- DUFLO, E., W. GALE, J. LIEBMAN, P. ORSZAG, AND E. SAEZ (2006): "Saving Incentives for Low-and Middle-Income Families: Evidence From a Field Experiment With H&R Block," *Quarterly Journal of Economics*, 121, 1311–1346.
- DWORAK-FISHER, K. (2011): "Matching Matters in 401(k) Plan Participation," Industrial Relations: A Journal of Economy and Society, 50, 713–737.
- EGAN, M. (2019): "Brokers Versus Retail Investors: Conflicting Interests and Dominated Products," Journal of Finance, 74, 1217–1260.

- EGAN, M. L., A. MACKAY, AND H. YANG (2022): "Recovering Investor Expectations From Demand for Index Funds," *Review of Economic Studies*, 89, 2559–2599.
- (2023): "What Drives Variation in Investor Portfolios? Evidence from Retirement Plans," *Working Paper*.
- EINAV, L., M. JENKINS, AND J. LEVIN (2012): "Contract Pricing in Consumer Credit Markets," *Econometrica*, 80, 1387–1432.
- FARRELL, J. AND P. KLEMPERER (2007): "Coordination and Lock-in: Competition with Switching Costs and Network Effects," *Handbook of industrial organization*, 3, 1967–2072.
- FEDERAL REGISTER (2010): "Reasonable Contract or Arrangement Under Section 408(b)(2)—Fee Disclosure; Interim Final Rule," Journal of Financial Economics, 75, 41600–38.
- GOURIEROUX, C. AND A. MONFORT (1996): Simulation-Based Econometric Methods, Oxford university press.
- GROPPER, M. (2023): "Lawyers Setting the Menu: The Effects of Litigation Risk on Employer-Sponsored Retirement Plans," *Working Paper*.
- GRUBER, M. J. (1996): "Another Puzzle: The Growth in Actively Managed Mutual Funds," Journal of Finance, 51, 783–810.
- GRUNEWALD, A., J. A. LANNING, D. C. LOW, AND T. SALZ (2020): "Auto Dealer Loan Intermediation: Consumer Behavior and Competitive Effects," Tech. rep.
- GUISO, L., A. POZZI, A. TSOY, L. GAMBACORTA, AND P. E. MISTRULLI (2022): "The Cost of Steering in Financial Markets: Evidence From the Mortgage Market," *Journal of Financial Economics*, 143, 1209–1226.
- HORTAÇSU, A. AND C. SYVERSON (2004): "Product Differentiation, Search Costs, and Competition in the Mutual Fund Industry: A Case Study of SP 500 Index Funds," *Quarterly Journal of Economics*, 119, 403–456.
- HOTZ, V. J. AND R. A. MILLER (1993): "Conditional Choice Probabilities and the Estimation of Dynamic Models," *Review of Economic Studies*, 60, 497–529.
- ICI (2022a): "The BrightScope/ICI Defined Contribution Plan Profile: A Close Look at 401(k) Plans, 2019," .
  - (2022b): "Trends in the Expenses and Fees of Funds, 2021," Investment Company Institute Research Perspective, 28.
- ILLANES, G. (2016): "Switching Costs in Pension Plan Choice," Working Paper.
- ILLANES, G. AND M. PADI (2019): "Retirement Policy and Annuity Market Equilibria: Evidence From Chile," Tech. rep.
- KLEMPERER, P. (1987a): "The Competitiveness of Markets with Switching Costs," The RAND Journal of Economics, 138–150.
- (1987b): "Markets with Consumer Switching Costs," *Quarterly Journal of Economics*, 102, 375–394.
- (1995): "Competition When Consumers Have Switching Costs: An Overview With Applications to Industrial Organization, Macroeconomics, and International Trade," *Review of Economic Studies*, 62, 515–539.

KOIJEN, R. S. AND M. YOGO (2016): "Shadow Insurance," Econometrica, 84, 1265–1287.

- LOSETO, M. (2023): "Plan Menus, Retirement Portfolios, and Investors' Welfare," Working Paper.
- LUCO, F. (2019): "Switching Costs and Competition in Retirement Investment," American Economic Journal: Microeconomics, 11, 26–54.
- MADRIAN, B. C. AND D. F. SHEA (2001): "The Power of Suggestion: Inertia in 401(k) Participation and Savings Behavior," *Quarterly Journal of Economics*, 116, 1149–1187.
- MCFADDEN, D. (1973): "Conditional Logit Analysis of Qualitative Choice Behavior," .
- MELLMAN, G. S. AND G. T. SANZENBACHER (2018): "401 (k) Lawsuits: What Are the Causes and Consequences?" *Issue in Brief*, 18–8.
- MORNINGSTAR (2021): "2021 U.S. Fund Fee Study," Morningstar Manager Research.
- NELSON, S. (2018): "Private Information and Price Regulation in the US Credit Card Market," *Working Paper*.
- NEPC (2015): "NEPC 2015 Defined Contribution Plan & Fee Survey: What a Difference a Decade Makes," .
- PECHTER, K. (2022): "Let's Concentrate on Recordkeeping," Retirement Income Journal.
- POOL, V. K., C. SIALM, AND I. STEFANESCU (2016): "It Pays to Set the Menu: Mutual Fund Investment Options in 401(k) Plans," *Journal of Finance*, 71, 1779–1812.

(2021): "Mutual Fund Revenue Sharing in 401(k) Plans," Working Paper.

- REUTER, J. AND D. P. RICHARDSON (2022): "New Evidence on the Demand for Advice Within Retirement Plans," .
- ROBLES-GARCIA, C. (2019): "Competition and Incentives in Mortgage Markets: The Role of Brokers," *Working Paper*.
- RUST, J. (1987): "Optimal Replacement of GMC Bus Engines: An Empirical Model of Harold Zurcher," *Econometrica*, 999–1033.
- SUN, L. AND S. ABRAHAM (2021): "Estimating Dynamic Treatment Effects in Event Studies With Heterogeneous Treatment Effects," *Journal of Econometrics*, 225, 175–199.
- SWEETING, A., D. JIA, S. HUI, AND X. YAO (2022): "Dynamic Price Competition, Learning-By-Doing, and Strategic Buyers," *American Economic Review*, 112, 1311–33.
- YOGO, M., A. WHITTEN, AND N. COX (2025): "Financial Inclusion Across the United States," Journal of Financial Economics, 166, 104003.

# **Online Appendix**

# A Data

### A.1 Identifying Providers

According to Form 5500 Schedule C, an employer interacts with multiple providers such as recordkeepers, investment advisors, auditors, and legal consultants. I focus on the main provider with primarily recordkeeping functions. Employers report a list of service codes for each provider in Form 5500 Schedule C, and BrightScope consolidates those service codes to provider functions that help me identify the main recordkeeper.

I use a hierarchy of provider functions to designate the main recordkeeper. I first consider the provider that is "Recordkeeper". If there is no recordkeeper, I move on to "Trustee", "Custodian", and then "Administrator". The vast majority of employers have a recordkeeper following this sequence of functions.<sup>34</sup> My results are relatively insensitive to the precise sequence of services in this hierarchy.

Occasionally, an employer may hire multiple record keepers, and I use the following rules to determine the main record keeper. These rules err on the side of not introducing provider switches due to adding or removing one of the multiple record keepers. First, for each employer  $\times$  provider, I define a tenure covering the consecutive years this provider performs record keeping functions. I drop a record keeper whose tenure is completely covered within the tenure of another record keeper. When multiple providers have the same tenures, I select the provider with higher compensation. Second, two record keepers overlap during the year of provider transition. In this case, I choose the record keeper with higher compensation. Only a few employers still have multiple providers after these two steps, and I drop these employers.

### A.2 Identifying Provider Switches

After designating a unique recordkeeper for each provider  $\times$  year, I define a provider switch when the employer changed its recordkeeper from the previous year. The precise year of recordkeeper switch may not correspond to the year with menu changes. For example, an employer may hire a new recordkeeper in late 2010 but revise the menu of funds in 2011. I use menu turnover from BrightScope plan menu data to revise switch timing. Around a window of plus and minus one year from the year of switch, I revise the year of switch to be the year with the highest menu turnover.

There has been some consolidation among recordkeepers that can result in recordkeeper changes. To adjust for these cases, I look at pairs of old and new providers

<sup>&</sup>lt;sup>34</sup>Bhattacharya et al. (2019) used a similar procedure.

with high frequencies in the data and cross check with publicly available information about recordkeeper mergers or acquisitions. These include Great-West's mergers with the defined contribution business of Putnam Investments and JP Morgan Retirement Plan Services and rebranding as Empower in 2014. Other consolidations are between smaller providers in the other categories. When recordkeeper A bought B in year T, I reassign the recordkeeper from B to A prior to T and do not consider this as a provider switch.

#### A.3 Provider Compensation

BrightScope collects provider direct compensation data from Form 5500 Schedule C. Providers receive both direct compensation from employers and indirect compensation through revenue sharing from mutual fund companies. Because I abstract away from separately modeling integrated versus third party recordkeepers, I consider the total fees paid by employers and employee participants. The total fees include both mutual fund expense ratios and direct compensation for recordkeeping services.

Direct compensation can be a fixed dollar amount per participant or a percentage of invested assets, but Form 5500 only reports the total amount of compensation. I convert the amount of direction compensation to a fraction of total plan assets. If an employer switches to a new provider in the middle of a year, the amount of compensation.

The compensation from Form 5500 may include services other than recordkeeping (Aronowitz, 2022a). Recordkeepers may also offer consulting, advisory, or legal services, and these compensations are not separately reported. In my analysis, I control for other provider functions. But those information are coarse and the control is likely not sufficient. In addition, recordkeeping compensation also includes commissions for individual trades, loan-origination, Qualified Domestic Relations Order (QDRO) related fees that are charged to the individual participants, rather than on a plan-wide basis. For example, in *Matousek v MidAmerican Energy Company*, plaintiffs claim a recordkeeping fee of \$326 to \$525 per participants computed from Form 5500 Schedule C over five years. But the employer disputed that the recordkeeping fee was only \$37.

# **B** Drivers of Mutual Fund Fee Decline

The decline in mutual fund fees over the past 25 years has contributed to the reduction in 401(k) fees. From 1996 to 2021, the average expense ratio of equity mutual funds decreased from 104 to 47 bps. The average expense ratios of hybrid and bond mutual funds fell from 95 to 57 bps and 84 to 39 bps, respectively (ICI, 2022b). This decline can be attributed to several factors (Morningstar, 2021; ICI, 2022b).

First, mutual funds benefit from economies of scale and scope. Certain costs such as transfer agency fees, director fees, accounting and audit fees, are relatively fixed in dollar terms, and these costs become increasingly negligible as funds gather more assets. In addition, the fixed operating costs of an asset manager also spread across a greater number of mutual funds, leading to cost savings.

Second, investors have become more fee-conscious over the past two decades. There has been a growing preference for index funds, which replicate the returns of market indices. Index funds are generally cheaper to manage compared to actively managed funds since they do not require extensive research on underlying companies, and their portfolio turnovers are lower. Even within the respective active and index fund universe, investors tend to gravitate towards the lowest-cost funds, resulting in increased competition and fee reduction.

Third, the popularity of no-load funds has contributed to the decline in expense ratios. Investors increasingly compensate for financial advice separately and prefer funds without sales charges. This explanation is less relevant for 401(k) plans since most funds offered in 401(k) menus are no-load share classes.

In addition to mutual fund fees, recordkeeping fees in 401(k) plans have also declined over time. Average recordkeeping fees have decreased, driven by improvements in IT technology, the scale of recordkeepers, and some standardization of recordkeeping services (Pechter, 2022). Based on its survey, NEPC (2015) shows that the average recordkeeping fees declined from 57 to 46 bps over plan assets or from \$118 to \$64 per participant from 2006 to 2015.

# C Comparison of Plan Fees with Government Sponsored Defined Contribution plans

In this appendix section, I compare plan fees from employer sponsored 401(k) plans with government-sponsored defined contribution plans. Started in 2017, OregonSaves is one of the earliest state-sponsored defined contribution plans for self-employed and employees at small employers that do not offer 401(k) plans. This program charges participants around 50 bps each year. There is another \$16 annual fee, which covers higher administrative costs due to low account balances (the average account balance

of \$1324 as of 2022) as participants have yet accumulated large amounts of savings in this new program. On the other hand, Federal Thrift Saving Plans (TSP) for federal and civil services employees charge participants 5.7 to 9 bps. TSP is the largest defined contribution plan in the world with over \$800 billion assets, and the low fees can reflect its economies of scale.<sup>35</sup>

To compare government sponsored programs with 401(k) plans with comparable scale, I fit a linear relationship between log plan fee and the log of total plan assets using my 401(k) data as of 2019, In Appendix Figure A1, I show the estimated plan fee - total asset curve and indicate fees charged by OregonSaves and TSP. The comparison shows that plans with comparable total assets (with the curve extrapolated out of sample to compare with TSP) would charge similar plan fees as these government-sponsored programs. Assuming these government bodies ensured that providers do not charge an unrealistic amount of markup, employers are likely to offer competitive fees to plan participants on average.

Appendix Figure A1: Comparison with Government Sponsored DC plans



Figure plots curve of total plan fee versus log of total plan assets, estimated from  $\ln \text{fee} = \beta_0 + \beta_1 \times \ln \text{plan}$  asset. Plan fees and total assets of OregonSaves and TSP are displayed as comparisons.

<sup>&</sup>lt;sup>35</sup>OregonSaves program detail: https://www.oregonsaves.com/savers/program-details. OregonSaves account statistics are based on https://crr.bc.edu/wp-content/uploads/2022/07/0regon\_July-2022.pdf. Also see Chalmers et al. (2021) for an analysis of OregonSaves. Other state-sponsored defined contribution plans, for example, CalSavers, have similar cost structures. TSP program costs https://www.tsp.gov/tsp-basics/expenses-and-fees/ and account balance statistics from https://www.frtib.gov/pdf/reading-room/FinStmts/TSP-FS-Dec2021.pdf

# **D** Additional Summary Statistics

	All plans	1st quartile	2nd quartile	3rd quartile	4th quartile
Number of plans	630,724	473,448	18,055	$11,\!075$	9,715
Coverage by number of plans					
Has plan menu (%)	9	5	65	86	94
Estimation sample $(\%)$	3	1	28	41	38
Coverage by plan assets					
Has plan menu (%)	84	18	70	89	95
Estimation sample $(\%)$	17	4	29	42	18

Appendix Table A1: Sample Coverage as of 2019

Appendix Table A2: Summary Statistics Across Employer Size Quartile

	1st quartile	2nd quartile	3rd quartile	4th quartile
Plan participants	306	398	624	2,703
Total plan assets (million)	6	16	36	224
Average account balance (thousand)	33	67	110	164
Fee				
Plan average expense ratio (bps)	52	48	46	40
Direct recordkeeping fees (bps)	27	18	13	8
Total plan fees (bps)	79	66	58	49
Number of funds	26	27	26	25
Menu turnover (%)	14	14	13	11
Switched provider $(\%)$	5	5	4	3

## Appendix Figure A2: Employer Plan Size Distributions



Figure shows the distribution of log average plan assets across employers in each size quartile.

	Fidelity	Vanguard	Empower	Principal Fin	ADP	Others
Market share $(\%)$	27	6	10	2	7	48
Plan participants	966	1,182	702	763	433	884
Total plan assets (million)	69	138	36	37	17	53
Average account balance (thousand)	90	149	75	83	51	82
Fee						
Plan average expense ratio (bps)	53	18	43	47	57	47
Direct recordkeeping fees (bps)	9	10	18	28	16	23
Total plan fees (bps)	62	29	61	75	73	70
Number of funds	28	27	26	24	27	25
Plan menu						
Fraction of index funds (%)	27	75	38	30	31	33
Menu turnover (%)	13	8	14	15	17	13
Fraction of proprietary funds (%)	50	88	0	16	0	10

Appendix Table A3: Provider Specific Summary Statistics

# E Additional Motivating Evidence

### E.1 Decline in Mutual Fund Fees

I explain the mutual fund fee decomposition details in Figure 3. Let  $p_{ft}$  denote the expense ratio of fund f in year t. Let  $w_{ift}$  denote the fraction of assets of employer i's plan allocated to fund f in year t. Let  $\iota_{ft}$  denote whether fund f is available in the marketplace as of year t. The fund is not available if its fund manager has not yet introduced this fund or if the fund has been delisted. I keep a balanced panel of employers for this analysis. The observed fee decline (black solid line) and the decline corresponding to changes in expense ratios (dashed line with circle markers) are constructed as follows

$$p_t^{\text{observed}} = \frac{\sum_{if} w_{ift} \times p_{ft}}{\sum_{if} w_{ift}}$$
$$p_t^{\Delta \text{fee in 2009 menu}} = \frac{\sum_{if} w_{if,2009} \times \iota_{ft} \times p_{ft}}{\sum_{if} w_{if,2009} \times \iota_{ft}}$$
$$p_t^{\Delta \text{fee in 2019 menu}} = \frac{\sum_{if} w_{if,2019} \times \iota_{ft} \times p_{ft}}{\sum_{if} w_{if,2019} \times \iota_{ft}}$$

The gray line in Figure 3 allows participants to rebalance among funds in 2009 menus. I consider a simple calculation to approximate rebalancing. I first hold allocation to each asset class fixed. These are nine broad asset classes: large/mid/small cap equities, international equities, bonds, alternatives, allocation funds, and target date funds. I assume, for example, a participant close to retirement will not shift from lower-risk bond funds to higher-risk equity funds because of a few basis point differences in expense ratios. Within each asset class, I allocate assets proportionally based on observed allocation but within funds included in 2009 menus. When all funds in 2009 menus are removed in an asset class, I extrapolate based on the average rate of fee decline within each asset class.

Let c(f) denote the asset class of fund f. The asset class average fee after rebalancing  $\hat{p}_{ict}$  is

$$\hat{p}_{ict} = \begin{cases} \frac{\sum_{f:c(f)=c} p_{ft} \times w_{ift} \times \iota_{ft}}{\sum_{f:c(f)=c} w_{ift} \times \iota_{ft}} & \text{if } \sum_{f:c(f)=c} w_{ift} \times \iota_{ft} > 0\\ \hat{p}_{ic,t-1} + \Delta p_c & \text{if } \sum_{f:c(f)=c} w_{ift} \times \iota_{ft} = 0 \end{cases}$$

Let  $w_{ict}$  denote the fraction of total balance that participants in employer *i* allocate to asset class *c* at time *t*. The gray line in Appendix Figure A3 is constructed as

$$p_t^{\Delta \text{fee in 2009 menu \& rebalance}} = \frac{\sum_{ic} \hat{w}_{ict} \times \hat{p}_{ict}}{\sum_{ic} \hat{w}_{ict}}$$

I have also estimated a demand system allowing time-varying fee sensitivities and

used this demand system to predict counterfactual allocation using the menu as of 2009. I obtained similar results to the gray line displayed in Figure 3.

Appendix Figure A3 displays the decomposition separately for active and index funds. Employers can reduce fees from updating menus toward cheaper and comparable funds within the active and index universe.



Appendix Figure A3: Decline in Mutual Fund Fees

Figure shows the decomposition of mutual fund fee changes separately for active and index funds.

### E.2 Providers Do Not Immediately Update Plan Menus

In this section, I show that providers do not immediately include cheaper funds when they become available.

Appendix Figure A4: Slow Adoption of Cheapest Share Class



Panel (a) shows the fraction of employers who offer Vanguard Target Date Funds (TDFs) that include cheaper institutional share class since its introduction in 2015. Panel (b) displays the fraction of employers who offer the cheapest share class over the age of the cheapest share class (initial offer year minus the current year).

I start with a case study of when Vanguard introduced institutional share classes for its target date fund (TDF) series in 2015. The institutional share class had an expense ratio of 10 bps, but the pre-existing investor share class had expense ratios ranging from 14 to 16 bps, depending on the retirement target date. Figure A4a shows that employers who offer Vanguard TDFs adopt the lower-cost institutional share class gradually, reaching 20% in 4 years. The institutional share class comes with higher investment limit. I restrict to employers whose participants invest over \$5 million into Vanguard TDFs according to the investment limit<sup>36</sup>

I extend the analysis more broadly in Figure A4b. For each fund  $\times$  year with multiple share class, I compute the fraction of employers who offer the cheapest share class, by the age of the cheapest share class (current year minus year when fund was first offered). To approximate investment limit, I look at the bottom 10th percentile of investment among plans who offer the cheapest share class in a given year. The gray bars show density for the age of cheapest share class across fund  $\times$  year.

Less then 10% of employers adopt the cheapest share class when they are first introduced. Adoption increases gradually reaching around 50% after 15 years. Due to revenue sharing discussed in Section 2, employers and providers may deliberately offer share classes that have higher expense ratios but lower fees on a net basis, which explains why the fraction may never reaches 100%. The density shown by the gray bars illustrates that it is rather frequent for asset managers to issue a new share class with the lowest expense ratio. But these cheaper share classes are not included to 401(k)plan menus with a timely fashion.

<sup>&</sup>lt;sup>36</sup>https://investor.vanguard.com/investor-resources-education/mutual-funds/ share-classes-of-vanguard-mutual-funds

## E.3 Additional Difference-in-Differences Results

Appendix Figure A5: Difference in Differences around Provider Switch



(a) Probability of Switching Providers

Panels (a) and (b) follow the same specification as Figure 4, except that controls  $X_{it}\beta$  only include plan number of participants. The outcome variable in (a) is fee over cheapest share class of the same mutual fund. To determine eligibility for the cheapest share class, I compare the asset invested in this mutual fund with the 10th percentile of asset invested across all employers who offer the cheapest share class in the same year. The outcome is unconditional: when the mutual fund does not have multiple share class or when the employer does not qualify, the outcome variable is zero for this employer × fund. The outcome variable in (b) is average fee across all mutual funds in the same asset class × Morningstar category. Panels (c) and (d) separately look at expense ratios and recordkeeping fees, as the two components of total plan fees.

# **F** Model Implementation

### F.1 Two-Stage RFP Simplification

I break the RFP into a two-stage process. In the first stage, her incumbent provider proposes a fee. The employer decides whether to accept her incumbent's offer. If she rejects, she incurs the switching cost and enters the second stage, where other providers compete with the incumbent. For simplicity, I assume the employer still incurs the switching costs by choosing the incumbent provider at the second stage. This captures cases when the previous incumbent is willing to match menus and services proposed by competing providers, and so there will be substantial changes even with the same provider.

This two-step process has two computational benefits. The first is dimension reduction. In the dynamic price setting game described in Section 5.3, each provider j sets an optimal fee conditional on the previous incumbent s identity. Suppose there are Jproviders in total. I need to solve for  $J^2$  optimal RFP fee. With the two-stage RFP simplification, all second-stages of RFPs are homogeneous and no longer depend on the previous incumbent s. Each incumbent provider still sets different fees in the first stage. So there are in total 2J optimal fees, instead of  $J^2$ . This simplification abstracts away from the variation in fees and transition probabilities across providers. However, since there are relatively few observations of provider switches, this variation may contain too much noise.

Second, all provider switches occur at the second stages of RFPs only. Since I observe provider choice probabilities conditional on switching, I can use them as CCP (Hotz and Miller, 1993) and do not need to estimate provider heterogeneity parameters  $\xi_j$  (except for providers in the other categories). I discuss this point later in Appendix F.2.

At the second stage of RFP, each provider sets a fee  $b_j^{2nd37}$ . These fees no longer depends on the identity of the previous incumbent. Employer *i* perceives the following utility from provider *j*:

$$\mathcal{U}_{j}^{2nd} + \epsilon_{ij} = -\alpha b_{j}^{2nd} + \nu_{j} + \delta \mathbb{E} V^{E}(j, b_{j}^{2nd}) + \epsilon_{ij}$$
(A1)

The probability that the employer chooses provider j is given by

$$q_j^{2nd} = \frac{e^{\mathcal{U}_j^{2nd}}}{\sum_k e^{\mathcal{U}_k^{2nd}}} \tag{A2}$$

Employer's expected utility from the second stage of RFP follows the standard in-

<sup>&</sup>lt;sup>37</sup>In the main text, my notation conditions on employer *i* to be more general, but it should become clear that the condition on *i* is redundant since employer specific quality  $\nu_i$  and cost  $c_i$  shifters are canceled out. So for simplicity in this appendix section, I drop the *i* subscript.

clusive value expression

$$\mathbb{E}\mathcal{U}^{2nd} = \ln\sum_{k} e^{\mathcal{U}_{k}^{2nd}}$$
(A3)

Provider's ex-ante value function at the second stage of RFP depends on the probability of winning and the continuation values as the winner providing services or the loser waiting for the next RFP.

$$V_j^{P,2nd} = q_j^{2nd} \left( b_j^{2nd} - c_j + \underbrace{\delta \mathbb{E} V_j^W(j, b_j^{2nd})}_{\text{Winner continuation value}} \right) + \sum_{k \neq j} q_k^{2nd} \underbrace{\delta \mathbb{E} V_j^L(k, b_k^{2nd})}_{\text{Loser continuation value}}$$
(A4)

Providers set an optimal fee  $b_j^{2nd}$  to maximize  $V_j^{P,2nd}$ . Solve FOC with respect to  $b_j^{2nd}$ :

$$b_{j}^{2nd} = c_{j} \underbrace{-\delta\left(\mathbb{E}V_{j}^{W}(j, b_{j}^{2nd}) - \frac{\sum_{k \neq j} q_{k}\mathbb{E}V_{j}^{L}(k, b_{k}^{2nd})}{1 - q_{j}^{2nd}}\right)}_{\text{Incremental continuation value}} + \left(\frac{1 + \delta \frac{\partial\mathbb{E}V_{j}^{W}(j, b_{j}^{2nd})}{\partial b_{j}^{2nd}}}{\alpha(1 - q_{j}^{2nd})}\right) \quad (A5)$$

At the first stage of RFP, the incumbent provider j proposes fee  $b_j^{1st}$ . The utility from rejecting j's offer and proceeding to the 2nd stage depends on the expected utility from Equation (A3) and switching cost  $\kappa_{st}$ .

$$\mathcal{U}_{j1}^{1st} + \epsilon_{i1} = -\alpha b_j^{1st} + \nu_j + \delta \mathbb{E} V^E(j, b_j^{1st}) + \epsilon_{i1}$$
(A6)  
$$\mathcal{U}_{j0}^{1st} + \epsilon_{i0} = \mathbb{E} \mathcal{U}^{2nd} - \kappa_{sw} + \epsilon_{i0}$$

The probability of accepting j's offer and the employer's expected utility at the 1st stage of RFP are, respectively:

$$q^{1st}(j) = \frac{e^{\mathcal{U}_{j1}^{1st}}}{e^{\mathcal{U}_{j1}^{1st}} + e^{\mathcal{U}_{j0}^{1st}}}$$
$$\mathbb{E}\mathcal{U}^{1st}(j) = \ln\left(e^{\mathcal{U}_{j1}^{1st}} + e^{\mathcal{U}_{j0}^{1st}}\right)$$

Ex-ante value function for provider k at the first stage depends on whether k = j is the incumbent.

$$V_{k}^{P,1st}(j) = \begin{cases} q^{1st}(j) \left( b_{j}^{1st} - c_{j} + \delta \mathbb{E} V_{j}^{W}(j, b_{j}^{1st}) \right) + (1 - q^{1st}(j)) V_{k}^{P,2nd} & \text{if } j = k \\ q^{1st}(j) \delta \mathbb{E} V_{k}^{L}(j, b_{j}^{1st}) + (1 - q^{1st}(j)) V_{k}^{P,2nd} & \text{if } j \neq k \end{cases}$$
(A7)

Similar to Equation (A5), I solve for the optimal fee at the 1st stage  $b_j^{1st}$  with FOC of Equation (A7) when j = k:

$$b_j^{1st} = c_j - \delta \left( \mathbb{E} V_j^W(j, b_j^{1st}) - V_k^{P,2nd} \right) + \left( \frac{1 + \delta \frac{\partial \mathbb{E} V_j^W(j, b_j^{1st})}{\partial b_j^{1st}}}{\alpha (1 - q^{1st}(j))} \right)$$
(A8)

Lastly, when solving employer's dynamic binary choice problems, value function at the 1st stage of RFP  $\mathbb{E}\mathcal{U}^{1st}(j)$  replace the value functions in Equation (5).

#### F.2 Recover Provider Fixed Effects from CCP

I only estimate one provider quality net of cost  $\xi_{other} = \nu_{other} - \alpha c_{other}$  for providers in the other category, and use conditional choice probabilities to recover the rest  $\xi_j$ .

First, I clarify why I estimate  $\xi$  instead of  $\nu$  and c separately. To solve providers' problem, keeping track of markup  $mk_j = b_j - c_j$  is sufficient. Plugging optimal fee, employer's utility in Equation (2) depends on the quality net of cost subtracted by markup. The employer's problem depends on  $\xi_i$  rather than  $\nu_i$  and  $c_j$  separately.

$$\underbrace{\nu_j - \alpha c_j}_{\xi_j} - \alpha m k_j$$

Next, I only need to solve for  $\xi_{other}$  for providers in the other category. I use choice probabilities  $q_j^{2nd}$  in Equation (A2) at the second stage of RFP to recover the rest of  $\xi_j$ . I observe choice probabilities conditional on switching in the data, which is close but not exactly  $q_j^{2nd}$ , because second stage choice shares also include some non-switch choices.

A simple adjustment step can bridge the gap. Let  $s_j(k)$  denote probabilities of choosing provider j conditional on switching from previous provider k. I have J(J-1)of the following equations, with only J unknown  $q_j^{2nd}$  for  $j \in \{1, ..., J\}$ . A minimum distance estimation allows me to recover  $q_j^{2nd}$ .

$$s_j(k) = \frac{q_j^{2nd}}{1 - q_k^{2nd}}$$

Then, I can use RFP choice probabilities at the second stage to recover  $\xi_j = \nu_j - \alpha c_j$ . Employer's expected utility from second stage RFP Equation (A3) can be rewritten as the following using the probabilities of choosing any provider in the other category (Arcidiacono and Miller, 2011). I use "Other" with the capitalized initial letter to indicate any provider in the other category with variety effect, where  $\xi_{Other} > \xi_{other}$ due to  $\nu_{Other} > \nu_{other}$ . Cost, markup, and continuation value are the same since I assume all providers in the other category are homogeneous.

$$\ln\left(\sum_{k} \exp(\mathcal{U}_{k}^{2nd})\right) = \underbrace{\xi_{Other} - \alpha m k_{other}^{2nd} + \delta \mathbb{E} V^{E}(other, b_{other}^{2nd})}_{\mathcal{U}_{Other}^{2nd}} - \ln(q_{Other}^{2nd})$$
(A9)

I normalize  $\xi_{Other} = 0$ .  $mk_{other}^{2nd}$  and  $V^E(other, b_{other}^{2nd})$  are both endogenized in the model equilibrium. Hence observed  $q_{Other}^{2nd}$  is sufficient to pin down the expected utility from second stage RFP for the employer, which I need to solve the rest of the model.

Next, compare the difference in the log of second stage RFP choice probabilities:

$$\ln(q_j^{2nd}) - \ln(q_{Other}^{2nd}) = \mathcal{U}_j^{2nd} - \mathcal{U}_{Other}^{2nd}$$
$$= \left(\xi_j - \alpha m k_j^{2nd} + \delta \mathbb{E} V^E(j, b_j^{2nd})\right) - \left(-\alpha m k_{other}^{2nd} + \delta \mathbb{E} V^E(other, b_{other}^{2nd})\right)$$

Then, I can recover  $\xi_j$  for all providers except those in the other category, which I estimate in the model.

Second stage choice probability  $\ln(q_{other}^{2nd})$  is not observed. So I use the probability of switching away from a provider in the other category to identify  $\xi_{other}$ , as discussed in Equation (G2).

#### F.3 State Variables

As discussed in Section 5, my model has two state variables: the identity of the current provider s and fee f. In my implementation, I separate fee f into optimal fee set at RFP b and an excess amount  $\tilde{f}$  due to higher markup due to outdated menus.

$$f = b + \tilde{f}$$

Provider s sets different fees at RFP depending on whether s wins the RFP from the first or second stage. I introduce a third state variable  $h \in \{0, 1\}$  corresponding to the stage of RFP. Then conditioning on (s, f) is the same as conditioning on  $(s, h, \tilde{f})$ .

The benefit of conditioning on  $(s, h, \tilde{f})$  is computational. Since  $\tilde{f} \ge 0$ , it is easier to project  $\tilde{f}$  on a grid. As a comparison, due to dynamic incentives, markup f can be negative. I would need to have a much larger grid. In practice, I discretize  $\tilde{f}$  over an evenly spaced grid with 121 points from 0 to 60 basis points.

### F.4 How Optimal Fees Affect Provider Continuation Values

I discuss how I solve for  $\frac{\partial \mathbb{E}V^W(j,b_j(s))}{\partial b_j}$  in optimal fee Equations (A5) and (A8). First, as explained in Appendix F.3, I solve for  $V^W(j,h,0)$  in practice where  $\tilde{f} = 0$  at RFP by

definition. I can rewrite the derivative of provider continuation values as

$$\frac{\partial \mathbb{E}V^W(j, b_j(s))}{\partial b_j} = \frac{\partial \mathbb{E}V^W(j, h, \tilde{f} = 0)}{\partial \tilde{f}}$$

Let  $\Delta$  denote the grid size for the grid of  $\tilde{f}$ . I can then approximate the derivative using the finite difference.

$$\frac{\partial \mathbb{E} V^W(j,h,0)}{\partial \tilde{f}} \approx \frac{\mathbb{E} V^W(j,h,\Delta) - \mathbb{E} V^W(j,h,0)}{\Delta}$$

where  $V^W(j, h, 0)$  and  $V^W(j, h, \Delta)$  are provider continuation values evaluated at  $\tilde{f} = \{0, \Delta\}$ , which come out of value function iterations.

#### F.5 Solving for Equilibrium

Conditional on structural parameters, I follow a nested algorithm to solve for the model equilibrium. In the outer loop, I solve for equilibrium markup at 1st and 2nd stage RFP, and provider value function at the 2nd stage of the RFP. This requires solving for 3J values, where J is the number of providers. The dimension is relatively low, and I use the MatLab function **fsolve**. In the inner loop, given markup and provider value function at RFP, I solve for employer value functions  $V^E(s, h, f)$  using value function iterations.

I rewrite Equations (A5) and (A8) as follows. Note that provider decisions only depend on markup.

$$mk_{j}^{2nd} = b_{j}^{2nd} - c_{j} = -\delta \left( \mathbb{E}V_{j}^{W}(j, 2, b_{j}^{2nd}) - \frac{\sum_{k \neq j} q_{k} \mathbb{E}V_{j}^{L}(k, 2, b_{k}^{2nd})}{1 - q_{j}^{2nd}} \right) + \frac{1 + \delta \frac{\partial \mathbb{E}V_{j}^{W}(j, 2, b_{j}^{2nd})}{\partial b_{j}^{2nd}}}{\alpha (1 - q_{j}^{2nd})}$$
$$mk_{j}^{1st} = b_{j}^{1st} - c_{j} = -\delta \left( \mathbb{E}V_{j}^{W}(j, 1, b_{j}^{1st}) - V_{k}^{P, 2nd} \right) + \frac{1 + \delta \frac{\partial \mathbb{E}V_{j}^{W}(j, 1, b_{j}^{1st})}{\partial b_{j}^{1st}}}{\alpha (1 - q^{1st}(j))}$$

The nested algorithm is as follows

- 1. Start with an initial guess for  $mk_i^{1st,0}, mk_i^{2nd,0}, V_i^{P,2nd,0}$
- 2. Inner loop value function iterations.
  - (a) Start with an initial guess for employer value function  $V^{E,0}(s,h,f)$
  - (b) Solve for employer's expected utility at second stage RFP  $\mathbb{E}\mathcal{U}^{2nd}$  using Equation (A9)
  - (c) Solve for first stage RFP provider choice probabilities  $q_j^{1st}$  and employer's expected utility  $\mathbb{E}\mathcal{U}^{1st}$
  - (d) Solve for the probability that the employer runs RFP  $\lambda(s, h, f)$  and update employer value function  $V^{E,1}(s, h, f)$ .

(e) Iteration step (a)-(d) until

$$||V^{E,t}(s,h,f)-V^{E,t+1}(s,h,f)||_{\infty}<\zeta$$

- 3. Solve recursively for provider value functions  $V_s^W(s, h, f)$  and  $V_j^L(s, h, f)$
- 4. Update provider value function at second stage of RFP  $V_j^{P,2nd,1}$  using Equation (A4)
- 5. Using the finite difference approximation for derivative in markup  $\frac{\partial V^W(j,s,b_j)}{\partial b_j}$  discussed in Appendix F.4
- 6. Update markup  $mk_j^{2nd,1}$  and  $mk_j^{1st,1}$
- 7. Iterate step 1 to 7 until

$$\max\left\{||mk_{j}^{2nd,1} - mk_{j}^{2nd,0}||_{\infty}, \ ||mk_{j}^{1st,1} - mk_{j}^{1st,0}||_{\infty}, \ ||V_{j}^{P,2nd,1} - V_{j}^{P,2nd,0}||_{\infty}\right\} < \zeta$$

### F.6 Random Coefficient

I use Monte Carlo to numerically integrate over employers with different fee sensitivity parameters  $\alpha_i$  with 50 Sobol draws. I have tried 100 Sobol draws and obtained quantitatively similar results. I have also tried evaluating integration using quadrature and obtained similar results. While quadrature method is more computationally efficient when evaluating integration, Monte Carlo simulation generates outcomes of employers at variance percentiles of fee sensitivity  $\alpha_i$ .

# G Auxiliary Models

The indirect inference estimation minimizes the distance between the vector of auxiliary model coefficients  $\mathcal{B}^{data}$  estimated using observed data, and  $\mathcal{B}^{model}(\Theta)$  estimated using data generated by my model at structural parameters  $\Theta$ . When computing  $\mathcal{B}^{model}(\Theta)$ , I use one observation per state (s, f), weighted according to the stationary distribution across states. To compute the stationary distribution, I start from a given distribution and iteratively apply the transition matrix across all states (s, f) until convergence. Transition probabilities between states are based on optimal decisions of the employer and providers at the model equilibrium.

The distance between  $\mathcal{B}^{model}(\Theta)$  and  $\mathcal{B}^{data}$  is the sum of squared differences between the two vectors of coefficients, where coefficients are weighted equally. I have seven auxiliary models for seven parameters, so the weighting matrix should not matter and I use the identity matrix for simplicity.

I estimate the following seven structural parameters: RFP cost  $\kappa_{rfp}$ , switching cost  $\kappa_{sw}$ , the average  $\alpha$  and standard deviation  $\sigma$  of fee sensitivity, attention in RFP decision  $\rho$ , providers' menu update probability  $\phi$ , and lastly the quality net of cost for providers in the other category  $\xi_{other}$ .

$$\Theta = \{\kappa_{rfp}, \kappa_{sw}, \alpha, \sigma, \rho, \phi, \xi_{other}\}$$

I use the following seven auxiliary models to identify these structural parameters. Let  $(y_{it}, \tilde{f}_{ijt}, t_{ijt})$  collect observed provider choice  $y_{it}$ , residualized plan fee  $\tilde{f}_{ijt}$ , and menu turnover  $t_{ijt}$  for each employer i in year t with provider  $y_{it} = j$ .

### G.1 Probability of Switching Providers: $\kappa_{rfp} + \kappa_{sw}$

The probability that employers switch providers identifies the sum of two transaction costs  $\kappa_{rfp} + \kappa_{sw}$ . Employers are less likely to switch providers if either RFP cost or switching cost is high. I do not observe RFPs and use Equation (G3) to separate  $\kappa_{rfp}$  and  $\kappa_{sw}$ . I also impose a condition that the incumbent provider prior to switching is not one from the other category, and later use the probability of switching away from a provider from the other category to recover  $\xi_{other}$ .

$$1\{y_{it} \neq y_{i,t-1}\} = \beta_1 + \epsilon_{ijt} \quad \text{if } y_{i,t-1} \neq other \tag{G1}$$

# G.2 Probability of Switching From Smaller Providers: $\xi_{other}$

In Equation (G2), I estimate the probability that employers switch away from a provider in the other category. Compared with  $\beta_1$  above, this probability identifies quality net of costs for providers in the other category:  $\xi_{other} = \nu_{other} - \alpha c_{other}$ . Providers in the other category have a lower net quality if  $\beta_2 > \beta_1$ .

$$1\{y_{it} \neq y_{i,t-1}\} = \beta_2 + \epsilon_{ijt} \quad \text{if } y_{i,t-1} = other \tag{G2}$$

I do not separate quality  $\nu$  from cost c, since only their difference matters after substituting the expression of optimal fees Equation (10) into employers' utility in Equation (2). I recover  $\xi_j$  for large providers (e.g. Fidelity and Vanguard) offline from conditional choice probabilities when employers switch providers, discussed in Section 6.4.

I cannot use the conditional probabilities of choosing any small providers in the data because employers may invite multiple providers from the other category to participate in their RFPs. Thus, I use Equation (G2) to recover  $\xi_{other}$  and leave the conditional probabilities of choosing small providers as a free moment to account for the variety from multiple small providers.<sup>38</sup>

### G.3 Fee Autocorrelation Separates $\kappa_{rfp}$ vs. $\kappa_{sw}$

Both RFP and switching costs affect the probability of switching providers in Equation (G1). To separate these two transaction costs, I use the autocorrelation of fees to infer the probability of RFPs. In my model, providers do not condition on previous fees when they set fees at RFPs. Under this assumption, the autocorrelation of fees at RFPs should be zero, which is the lower bound. Then, I estimate an upper bound of the autocorrelation when there is no menu turnover, which I can observe in the data. The upper bound is less than one in the data due to changes in expense ratios, reallocation across different funds, and measurement errors in recordkeeping fees.

When employers do not switch providers, an autocorrelation close to zero suggests frequent RFPs, low RFP costs, and high switching costs. Alternatively, an autocorrelation close to the upper bound suggests infrequent RFPs, high RFP costs, and low switching costs.

In my auxiliary model, I regress residualized fees on lag residualized fees interacted with an indicator for any menu update. I condition on observations without provider switches. The ratio of  $\beta_{3t}/\beta_{3n}$  corresponds to the relative frequency of RFPs.

$$\tilde{f}_{ijt} = \beta_{3t} 1\{t_{ijt} > 0\} \tilde{f}_{ij,t-1} + \beta_{3n} 1\{t_{ijt} = 0\} \tilde{f}_{ij,t-1} + \gamma 1\{t_{ijt} = 0\} + \Gamma_i + \epsilon_{ijt} \quad \text{if } y_{it} = y_{i,t-1}$$
(G3)

In Appendix Table A4, I report estimates of autocorrelation in residualized fees. In the first column, I show that the autocorrelation in residualized fees when employers switch providers is statistically indistinguishable from zero, consistent with the assumption in my model. In the second column, the upper bound of fee autocorrelation when employers do not switch providers and have no menu updates is 0.42. The ratio of

 $<sup>^{38}\</sup>xi_{other} = \nu_{other} - \alpha c_{other}$  depends on employers' fee sensitivities, which vary across employers. I assume  $\nu_{other}$  and  $c_{other}$  also vary across employers such that  $\xi_{other}$  is constant across employers.

the autocorrelation with menu updates to the autocorrelation without menu updates is 0.76, which is the empirical counterpart to the autocorrelation in the fee state variable of my structural model.

	(1)	(2)
VARIABLES	Full sample	Not switch providers
Switch providera	0.04***	
Switch providers	-0.04	
	(0.00)	
Lag residualized fee $\times$ switch providers	0.02	
	(0.02)	
Lag residualized fee $\times$ not switch providers	$0.33^{***}$	
	(0.01)	
Lag residualized fee $\times$ menu update		0.32***
		(0.01)
Lag residualized fee $\times$ no menu update		0.42***
		(0.01)
Menu update		-0.01***
-		(0.00)
Observations	$95,\!538$	91,550
Employer FE	X	Х

Appendix Table A4: Fee Autocorrelation

Table shows regression of residualized fees on lag residualized fees, by whether employers switch providers, not switch providers and have menu updates. I control for employer fixed effects and cluster standard errors at the employer level.

## G.4 Probability that Providers Do Not Update Menus: $\phi$

In Equation (G4), I estimate the probability that providers do not update menus, which identifies the constant and exogenous probability that providers update menus  $\phi$  outside RFPs.

$$1\{t_{ijt} = 0\} = \beta_4 + \epsilon_{ijt} \tag{G4}$$

In the model, I assume providers always update menus at RFPs. So  $\beta_4$  corresponds to  $(1 - \phi) \times (1 - \lambda)$  rather than  $1 - \phi$ , and needs to be estimated with other structural parameters.

### G.5 Sensitivity of Switching With Respect to Fees: $\alpha \times \rho$

In Equation (G5), I estimate the sensitivity of switching providers with respect to lag residualized fees, after absorbing employer fixed effects. The coefficient  $\beta_5$  corresponds to a mixture of employer fee sensitivity  $\alpha$  and attention  $\rho$  when they decide whether to conduct RFPs. I separate these two parameters using Equation (G6).

$$1\{y_{it} \neq y_{i,t-1}\} = \beta_5 \times \tilde{f}_{ij,t-1} + \Gamma_i + \epsilon_{ijt} \tag{G5}$$

This auxiliary model highlights the importance of controlling for employer fixed effects  $\Gamma_i$ . Some employers may have higher  $\tilde{f}_{ij,t-1}$  potentially because of higher  $\Gamma_i$ , if they offer more services in their 401(k) plans. When considering whether to run RFPs, employers focus on fees relative to their production costs, or  $\tilde{f}_{ij,t-1} - \Gamma_i$ . When estimating Equation (G6), if we only use  $\tilde{f}_{ij,t-1}$ , we would underestimate  $\beta_5$  due to an attenuation bias.<sup>39</sup> Appendix Table A5 reports  $\beta_5$  with and without employer fixed effects for the full sample. After controlling for employers fixed effects,  $\beta_5$  becomes three times larger.

VARIABLES	(1)	(2)
Lag residualized fee	$0.03^{***}$ (0.00)	$0.09^{***}$ (0.01)
Observations Employer FE	$95,\!538$	$95,538 \\ { m X}$

Appendix Table A5: Sensitivity of Switching With Respect to Fees

Regression coefficient of switching on residualized fees.

#### G.6 Fee Reductions When Employers Switch Providers: $\alpha$

I use the average fee reductions when employers switch providers to separate the average fee sensitivity  $\alpha$  from attention  $\rho$ . Fee reductions when employers switch providers include two components. The first component is the accumulated fee growth relative to the trend since the previous RFP. After recovering the probability of RFP using Equation (G3) and measuring fee growth  $g_j$  in the data, this component is known. The second component captures higher markups charged by incumbent providers at RFPs. Comparing incremental markups with the probability that employers choose their incumbents at RFPs (which is known after Equation (G3)) identifies the average fee sensitivity  $\alpha$ . Suppose  $\alpha_i$  is high on average. Even though employers face switching costs, their incumbents cannot substantially raise markups, and fee reductions from switching providers are thus lower.

$$\Delta f_{ijt} = \beta_6 + \epsilon_{ijt} \quad \text{if } y_{it} \neq y_{i,t-1} \tag{G6}$$

To provide additional evidence, I examine how fee reductions from switching providers vary across employer characteristics. I regress fee changes on whether employers switch

<sup>&</sup>lt;sup>39</sup>One caveat is that this method cannot address time-varying unobserved production costs, so I potentially still underestimate employer sensitivity of switching providers with respect to fees.

providers, interacted with a set of employer covariates.

$$\Delta \tilde{f}_{ijt} = \beta^0 + \beta^{sw} Switch_{it} + X'_{it}\beta + Switch_{it} \times X'_{it}\beta^{sw,x} + \epsilon_{ijt}$$

All covariates are standardized. Appendix Table A6 reports coefficient  $\beta^{sw,x}$ . Employers reduce fees on average when they switch. Less negative  $\Delta \tilde{f}_{ijt}$  (positive  $\beta^{sw,x}$ ) suggests employers are more fee sensitive. Employer covariates include average total plan assets, which I also use to group employers into size quartiles, and the average number of plan participants. Larger employers should be more sophisticated and have higher fee sensitivities. I indeed estimate positive  $\beta^{sw,x}$  in columns (1) and (2). In addition, I include

	(1)	(2)	(3)
VARIABLES			
Provider switch	-9.67***	-9.78***	-5.33***
$\times$ plan asset	(0.24) $4.70^{***}$	(0.25)	(0.96)
$\times$ num participants	(0.27)	2.22***	
$\times$ employer contribution		(0.26) $1.27^{***}$	
$\times$ participation rate		(0.26) $1.14^{***}$	
$\times$ ESG score		(0.29)	5.11***
			(0.91)
Observations	$95,\!538$	88,291	3,622

Appendix Table A6: Fee Reduction When Switching Providers

Regression coefficient of changes in fees on provider switches interacted with a set of employer characteristics

the fraction of employer contribution and the rate of employee participation in 401(k) plans. Employers can choose to match all or a fraction of employees' contributions to their 401(k) accounts. A higher employer contribution suggests employers value their employees' retirement benefits more. Not all employees participate in 401(k) plans, and a higher participation rate suggests that the employers put more effort into educating their employees about saving for retirement. In both cases, these employers should be more fee sensitive, and I estimate positive  $\beta^{sw,x}$  on these covariates in column (2).

Lastly, I look at employers' ESG scores from S&P Global. ESG scores are only available in the later part of the sample, and I use the average score for each employer for all periods. Employers with higher ESG scores also have positive  $\beta^{sw,x}$ , suggesting they are more fee-sensitive. It is reasonable that employers who value environmental, social and governance objectives also prioritize their employees' benefits in their 401(k) plans.

As another way to see how  $\beta_6$  captures fee sensitivities, I estimate a version of my

model where I set  $\alpha_i$  to be a very high value such that providers essentially engage in perfect competition. In Equation (10), as  $\alpha_i$  approaches infinity, providers set fees at RFPs equal to costs offset by future gains per period, and the average markups would be zero. I estimate structural parameters ( $\kappa_{sc}, \kappa_{rfpc}, \rho, \phi, \xi_{other}$ ) other than fee sensitivities to match the rest of auxiliary model coefficients. Appendix Figure A6 shows that fee reductions from switching providers are much smaller under perfect competition.



Appendix Figure A6: Fee changes when switching providers

#### G.7 Variance of Fee Reductions From Switching: $\sigma$

Since the average fee reductions when employers switch providers identifies the average fee sensitivity  $\alpha$ , the variance of fee reductions from switching should reflect the variance of fee sensitivity  $\sigma^2$ .

$$Var(\Delta f_{ijt}) = \beta_7 \text{ if } y_{it} \neq y_{i,t-1}$$
 (G7)

However, the variance of fee reductions also includes the variance of changes in unobserved production costs as well as measurement errors in fees. To see this, I expand residualized fees  $\tilde{f}_{ijt}$  into the state variable capturing markups  $f_{ijt}$ , an employer fixed effect for constant production costs  $\Gamma_i$ , and a time-varying component of production costs which can also include measurement errors in fees  $\epsilon_{ijt}$ . Although fee changes difference out  $\Gamma_i$ , the variance of fee changes still captures  $Var(\Delta \epsilon_{ijt})$ , which I assume to be independent from markups.

$$\tilde{f}_{ijt} = f_{ijt} + \Gamma_i + \epsilon_{ijt}$$
$$\Delta \tilde{f}_{ijt} = \Delta f_{ijt} + \Delta \epsilon_{ijt}$$
$$Var(\Delta \tilde{f}_{ijt}) = Var(\Delta f_{ijt}) + Var(\Delta \epsilon_{ijt}) \text{ if } y_{it} \neq y_{i,t-1}$$

To make progress, I introduce another auxiliary model to account for  $Var(\Delta \epsilon_{ijt})$ . Note that markup changes when employers switch providers  $\Delta f_{ijt}$  are negative in my model, after controlling for provider quality and cost fixed effects. Any positive  $\Delta \tilde{f}_{ijt}$  has to come from  $\epsilon_{ijt}$ . If  $Var(\Delta \epsilon_{ijt})$  is small, it is unlikely to see large positive changes in residualized fees. Hence, I add the eighth auxiliary model to recover  $Var(\Delta \epsilon_{ijt})$ .

$$E[\Delta \tilde{f}_{ijt} | \Delta \tilde{f}_{ijt} > 0, Switch_{it}] = \beta_8 \tag{G8}$$

My structural model does not have  $\Delta \epsilon$ , so I simulate markups and analytically incorporate noise terms assuming  $\Delta \epsilon \sim \mathcal{N}(0, \sigma_{\epsilon})$  and that  $\Delta \epsilon_{ijt}$  is independent from markups.

I expand  $\beta_8$  as the following. I can use my structural model to simulate the distribution of markup changes conditional on switching  $dF(\Delta f_{ijt}|Switch_{it})$ , and evaluate the standard normal CDF  $\Phi$  and PDF  $\phi$  at  $\frac{-\Delta f_{ijt}}{\sigma_{\sigma}}$  to compute the structural counterpart of  $\beta_8$ 

$$\begin{split} E[\Delta \tilde{f}_{ijt} | \Delta \tilde{f}_{ijt} &> 0, Switch_{it}] \\ &= \frac{E[(\Delta f_{ijt} + \Delta \epsilon_{ijt})1\{\Delta f_{ijt} + \Delta \epsilon_{ijt} > 0\}|Switch_{it}]}{E[\Delta f_{ijt} + \Delta \epsilon_{ijt} > 0|Switch_{it}]} \\ &= \frac{\int_{\Delta f} \Pr(\Delta \epsilon_{ijt} > -\Delta f_{ijt})E[\Delta f_{ijt} + \Delta \epsilon_{ijt}|\Delta \epsilon_{ijt} > -\Delta f_{ijt}]dF(\Delta f_{ijt}|Switch_{it})}{\int_{\Delta f} \Pr(\Delta \epsilon_{ijt} > -\Delta f_{ijt})dF(\Delta f_{ijt}|Switch_{it})} \\ &= \frac{\int_{\Delta f} \left(1 - \Phi(\frac{-\Delta f_{ijt}}{\sigma_{\sigma}})\right)\Delta f_{ijt} + \sigma_{\epsilon}\phi(\frac{-\Delta f_{ijt}}{\sigma_{\sigma}})dF(\Delta f_{ijt}|Switch_{it})}{\int_{\Delta f} 1 - \Phi(\frac{-\Delta f_{ijt}}{\sigma_{\sigma}})dF(\Delta f_{ijt}|Switch_{it})} \end{split}$$

#### G.8 Alternative Assumption: Variance of Residualized Fees: $\sigma$

For robustness, I consider an alternative auxiliary model where I directly fit the variance of residualized fees, ruling out variation in unobserved production costs  $Var(\Delta \tilde{f}_{ijt}) = Var(\Delta f_{ijt})$ . Because I condition on employers size quartiles, this assumption is treating employers with similar size as having the same economies of scale. Any unobserved services or plan features do not lead to differences in costs.

$$Var(\tilde{f}_{ijt}) = \beta_{7b} \tag{G7b}$$

# H Hedonic Regression

To residualize fees, I use a two-way fixed effect regression in Equation (12), which I repeat below

$$fee_{ijt} = \underbrace{X_{it}\beta}_{\text{Observed}} + \underbrace{\Gamma_j + \tau_{jt} + \tau_{s(i)t}}_{\text{Provider FE}} + \underbrace{\Gamma_i + \epsilon_{ijt}}_{\tilde{f}_{ijt}: \text{Residualized fee}}$$
(A10)

I use a balanced panel of employers so time trends do not reflect sample selections. I rely on employer switches to separate employer and provider fixed effects (Abowd et al., 1999). Due to limited number of observations in 2009 when many employers have not started reporting Form 5500, I include year 2010 to 2019.

First,  $X_{it}$  captures observed plan characteristics. I assume employers determine these characteristics as menu composition and types of services from the providers, and abstract away from how providers steer menu composition. Most providers have access to funds across all major asset classes and investment styles, and can be relatively indifferent to employers' requests. In Appendix Table A9, I show that provider fixed effects explain generally less than 20% of the variation in menu compositions, whereas provider and employer fixed effects together explain around 70%.

Specially, I include compositions of funds from broad asset classes, such as bonds and large, mid, and small cap equities. To further control for menu composition within broad asset classes, I also include the relative fee of the Morningstar category within each asset class to account for whether participants have more specific investment requests that come with higher expense ratios. To construct this relative fee measure, I regress the Morningstar category average fees on asset class fixed effects and obtain the residuals. These residuals measure whether funds in this Morningstar category are on average more expensive within the asset class. To aggregate at plan  $\times$  year level, I compute the average of the relative fee of each fund.

I do not control for the fraction of index funds, because offering index funds with lower expense ratios can be driven by both markups and heterogeneous services. It is reasonable for employers to offer more active funds catering to participants' preferences, and higher expense ratios of active funds suggest higher production costs. Alternatively, employers may allow providers to charge high markups by designing plan menus with expensive active funds. Also, Appendix Table A9 shows that provider fixed effects explain a larger fraction of variation in fraction of index funds on plan menus, suggesting that employers main jointly choose providers and the fraction of index funds on their plan menus.

 $X_{it}$  also includes whether the main provider offers additional services such as actuary, advisory, insurance, and asset management. The data includes limited information about services. I only observe recordkeepers' service codes such as advisory, asset management, and insurance.

I also control for the across time variation of employer's size in  $X_{it}$  to reflect how changes in size affect fees. The employer fixed effects  $\Gamma_i$  absorb the variation in fees due to persistent differences in economies of scale or fee sensitivities across employers.

Second, I residualize time trends to control for the decline in production costs. Time trends are provider specific  $\tau_{jt}$  to allow for different rates of decline in production costs, and also specific to each employer size quartile  $\tau_{s(i)t}$  to allow for different rates of decline across different sizes. For example, Vanguard might have already reached low production cost at the beginning of the sample due to large asset under management in its mutual funds. So incremental reduction can be smaller compared to other providers. With a stationary model, I assume that the employer does not have a higher preference for Vanguard in the early part of my sample (higher  $\nu_{Vanguard}$ ). To relax this assumption requires a non-stationary model with calendar time as another state variable. This extended version is much more computationally demanding, where additional economic insights can be limited. The stable market share across main providers shown in Figure 1b suggests that it is reasonable to assume employers have relatively constant preferences across providers.

Third, I residualize provider fixed effects  $\Gamma_j$  to control for persistent differences in quality or production costs across providers. However, in my structural model, providers with different  $\nu_j$  and  $c_j$  can also charge different markups. This is not a major concern since I do not use my model to target any average fee moments. I tried including provider fixed effects in residualized fees in previous versions. I obtained similar estimates for structural parameters, but have to introduce two sets of fixed effects for provider quality  $\nu_j$  and production costs  $c_j$ .

In Appendix Table A7 and Appendix Table A8, I report coefficients of regression estimates. I report two specifications, and use column (2) to residualize fees in my main analysis. All variables are standardized except for provider fixed effects and time trends. Fees tend to be higher if the employers offer more mid/small cap equity funds, international equity funds, and alternative funds. Fees are lower with more cash funds and large-cap equities. Within each asset class, fees are higher if employers offer Morningstar categories where funds typically have higher expense ratios.

Comparing across providers, Vanguard has the lowest fees, followed by Empower and Fidelity. Principal Financial Group has the highest fees. In column (1), I only include time trend without interaction, which suggests that fees decline by 2.7 bps each year on average. In column (2), I interact time trends with provider fixed effects and employer size quartile fixed effects. Fees decline by 1 to 4 bps per year across different providers. Relative to the largest providers in the 4th quartile, smaller providers see greater fee reductions.

The hedonic regressions have R-squared over 80%, after including employer fixed

effects. But these employer fixed effects can capture both markups and production costs. Hence, I only include employer fixed effects to help estimate provider fixed effects and I still keep them as part of the residualized fees.

	(1)	(2)
VARIABLES		
Frac international	40.65***	35.05***
	(2.68)	(2.66)
Frac alternative	52.84***	49.89***
	(3.85)	(3.81)
Frac cash	-43.22***	-40.95***
	(3.02)	(3.01)
Frac bond	3.14	4.24*
	(2.40)	(2.38)
Frac small cap equity	31.15***	34.01***
	(3.42)	(3.39)
Frac mid cap equity	$22.95^{***}$	22.80***
	(2.78)	(2.75)
Frac large cap equity	-4.67**	$-5.91^{***}$
	(1.93)	(1.91)
Frac target date fund	$-5.81^{***}$	-6.63***
	(1.60)	(1.58)
Category fee within asset class	$2.05^{***}$	$2.05^{***}$
	(0.06)	(0.06)
Actuary service	0.78	1.42
	(3.43)	(3.38)
Advisory service	$1.16^{***}$	0.17
	(0.38)	(0.37)
Insurance service	-3.57***	-3.02***
	(0.51)	(0.51)
Asset management service	-0.21	-0.11
	(0.34)	(0.33)
Ln number of participants	$-1.81^{***}$	$-1.67^{***}$
	(0.34)	(0.34)
Observations	30,535	30,535
R-squared	0.82	0.83
Employer FE	X	Х

Appendix Table A7: Two Way Fixed Effects Regression

Two way fixed effects regression of total fees including investment and recordkeeping fees on asset category composition, number of participants, provider specific time trends, provider and employer fixed effects. Coefficients are other variables are reported in the second part of in Appendix Table A8 I use a balanced panel of plans from 2010-2019.

	(1)	(2)
VARIABLES		
	0.00	0.01
ADP	-0.89	-0.91
_	(1.03)	(1.02)
Empower	-13.05***	-12.62***
	(0.80)	(0.79)
Fidelity	$-6.71^{***}$	$-6.49^{***}$
	(0.60)	(0.60)
Principal Fin	$8.62^{***}$	$7.47^{***}$
	(1.12)	(1.12)
Vanguard	-21.23***	-21.09***
-	(1.06)	(1.04)
Year	-2.71***	· · · ·
	(0.03)	
Year $\times$ Other	( )	-2.48***
		(0.05)
$Year \times ADP$		-2.90***
		(0.15)
Year × Empower		-2 36***
		(0.13)
Vear × Fidelity		-2 69***
		(0.06)
Voor × Principal Fin		2 07***
Tear × T Interpar Fin		-3.97
Voor V Vonguard		0.19)
Teal × valiguard		-0.95
Veen V employee 1st quantile		(0.08)
Tear × employer 1st quartile		-0.90
Varia v and large for d and startile		(0.13)
Year $\times$ employer 2nd quartile		-0.04
		(0.08)
Year $\times$ employer 3rd quartile		-0.69
		(0.06)
Observations	20 525	20 525
Deservations Deservations	30,333	30,333
n-squared Error larger EE	0.82	0.83
Employer FE	А	Λ

Appendix Table A8: Two Way Fixed Effects Regression (Continued)

Two way fixed effects regression of total fees including investment and recordkeeping fees on asset category composition, number of participants, provider specific time trends, provider and employer fixed effects. Coefficients are other variables are reported in the first part of in Appendix Table A7 I use a balanced panel of plans from 2010-2019.

Appendix Table A9: Variation of Employer Characteristics Explained by Provider and Employer Fixed Effects

	Provider FE	Provider & employer FE
Frac large cap equity	0.076	0.814
Frac mid cap equity	0.081	0.796
Frac small cap equity	0.089	0.771
Frac bond	0.191	0.791
Frac cash	0.079	0.666
Frac international	0.254	0.830
Frac alternative	0.049	0.757
Frac target date fund	0.252	0.860
Morningstar category average fee within asset class	0.150	0.769
Fund fee within category	0.322	0.817
Frac index funds	0.321	0.835

Table reports  $R^2$  of regressing employer  $\times$  year level menu characteristics on provider fixed effects and both provider and employer fixed effects.

# I Appendix Tables and Figures



Appendix Figure A7: Sample RFP Timeline

Source: SageView, "401(k) RFP Guide: A step-by-step resource to selecting the right advisor for your retirement program"

Appendix Table A10: Probability of RFPs and Offline Parameter Estimates

Employer quartiles	1st	2nd	3rd	4th
Prob of RFP	0.30	0.23	0.18	0.15
Conditional choice probability (%)				
ADP	0.08	0.08	0.02	0.02
Empower	0.17	0.17	0.14	0.14
Fidelity	0.17	0.17	0.23	0.23
Principal Fin	0.04	0.04	0.06	0.06
Vanguard	0.05	0.05	0.06	0.06
Other	0.49	0.49	0.48	0.48
Production cost decline / fee state variable growth rates (bps)				
ADP	2.91	2.76	2.62	2.39
Empower	1.94	1.79	1.65	1.42
Fidelity	2.44	2.29	2.15	1.92
Principal Fin	3.93	3.78	3.64	3.41
Vanguard	1.16	1.01	0.87	0.64
Other	2.24	2.09	1.95	1.72

Table displays the average probability of running RFPs in the first row. The middle panel displays conditional choice probabilities at RFPs. The bottom panel displays production cost decline for each provider  $g_j$  measured using the data.
Employer quartiles		1	st	21	nd	3	rd	4	$^{\mathrm{th}}$
		Model	Data	Model	Data	Model	Data	Model	Data
Pr of switching from large providers	eq. (G1)	0.054	$0.030 \\ (0.002)$	0.032	0.024 (0.001)	0.022	0.018 (0.001)	0.019	0.014 (0.001)
Pr of switching from other providers	eq. (G2)	0.059	$0.065 \\ (0.003)$	0.067	0.067 (0.002)	0.066	$\begin{array}{c} 0.067 \\ (0.002) \end{array}$	0.049	$\begin{array}{c} 0.049 \\ (0.002) \end{array}$
Fee autocorrelation ratio	eq. ( <b>G3</b> )	0.597	$\begin{array}{c} 0.599 \\ (0.042) \end{array}$	0.721	$\begin{array}{c} 0.721 \\ (0.025) \end{array}$	0.788	$0.788 \\ (0.024)$	0.818	$\begin{array}{c} 0.819 \\ (0.024) \end{array}$
Pr of no menu update	eq. (G4)	0.334	$\begin{array}{c} 0.333 \\ (0.004) \end{array}$	0.272	$\begin{array}{c} 0.272 \\ (0.004) \end{array}$	0.247	$\begin{array}{c} 0.247 \\ (0.003) \end{array}$	0.270	$\begin{array}{c} 0.270 \\ (0.003) \end{array}$
Sensitivity of switching wrt fee	eq. (G5)	0.090	$0.119 \\ (0.016)$	0.096	$0.098 \\ (0.015)$	0.091	$\begin{array}{c} 0.091 \\ (0.015) \end{array}$	0.057	$\begin{array}{c} 0.057 \\ (0.013) \end{array}$
Fee changes from switching	eq. (G6)	-0.138	-0.142 (0.009)	-0.093	-0.094 (0.008)	-0.076	-0.076 (0.007)	-0.044	-0.044 (0.006)
Variance of fee changes from switching	eq. (G7)	0.087	$\begin{array}{c} 0.079 \\ (0.003) \end{array}$	0.065	$0.066 \\ (0.002)$	0.045	$\begin{array}{c} 0.048 \\ (0.002) \end{array}$	0.032	$\begin{array}{c} 0.034 \\ (0.002) \end{array}$
Positive fee changes from switching	eq. (G8)	0.174	$0.179 \\ (0.008)$	0.164	$0.163 \\ (0.006)$	0.137	$\begin{array}{c} 0.136 \\ (0.006) \end{array}$	0.120	$0.118 \\ (0.006)$

Appendix Table A11: Fits of Auxiliary Model Coefficients

Table displays auxiliary model coefficients estimated using data generated by the structural model and coefficients estimated using observed data. Standard errors of coefficients estimated using observed data are reported in parenthesis. These standard errors are bootstrapped using 250 samples with replacement.

Employer quartiles		1st		2	2nd		3rd		4th	
		Coef	SE	Coef	SE	Coef	SE	Coef	SE	
RFP cost bps \$ thousand	$\kappa_{rfp}/\bar{\alpha}$	$\begin{array}{c} 16.0 \\ 10 \end{array}$	(0.1)	$\begin{array}{c} 14.4\\ 22 \end{array}$	(0.1)	$8.6 \\ 31$	(0.1)	0.8 18	(0.0)	
Switching cost bps \$ thousand	$\kappa_{sw}/\bar{\alpha}$	$\begin{array}{c} 68.6\\ 44 \end{array}$	(0.3)	$51.6\\80$	(0.3)	$25.5 \\ 91$	(0.2)	$1.9 \\ 42$	(0.0)	
Fee sensitivity: mean Coefficient Elasticity	$\bar{\alpha}$	$\begin{array}{c} 0.07 \\ 1.0 \end{array}$	(0.00)	$\begin{array}{c} 0.11\\ 1.4 \end{array}$	(0.00)	$0.22 \\ 2.9$	(0.00)	$4.92 \\ 62.2$	(0.08)	
Fee sensitivity: st.dev. Coefficient Elasticity	$\sigma_{lpha}$	$0.09 \\ 1.0$	(0.00)	$\begin{array}{c} 0.14 \\ 1.4 \end{array}$	(0.00)	$0.29 \\ 2.8$	(0.00)	$14.35 \\ 53.5$	(0.33)	
RFP attention Coefficient	ρ	0.94	(0.00)	0.95	(0.00)	0.69	(0.00)	0.26	(0.00)	
Pr menu update	$\phi$	0.50	(0.00)	0.62	(0.00)	0.68	(0.00)	0.63	(0.00)	
Net quality of other providers Relative to other providers	$\xi/\bar{\alpha}$	-20.8	(0.1)	-23.5	(0.2)	-16.1	(0.2)	-1.2	(0.0)	
ADP		-1.3		1.2		0.5		0.2		
Principal Fin.		-2.8		0.3		1.7		0.3		
Empower		1.4		3.1		3.7		0.4		
Fidelity		1.4		2.9		4.7		0.4		
Vanguard		-2.7		0.5		2.7		0.5		

Appendix Table A12: Structural Parameter Estimates, No Unobserved Production Costs

Table reports parameter estimates when I directly fit variance of residualized fees. Standard errors are computed using the asymptotic variance formula, where the variance-covariance matrix of auxiliary model coefficients is computed using the bootstrap method over 250 iterations. In each iteration, I sample employers with replacement and estimate auxiliary model coefficients.

Employer quartiles		1st		21	2nd		3rd		4th	
		Model	Data	Model	Data	Model	Data	Model	Data	
Pr of switching from large providers	eq. (G1)	0.057	$\begin{array}{c} 0.030 \\ (0.002) \end{array}$	0.035	$\begin{array}{c} 0.024 \\ (0.001) \end{array}$	0.026	$0.018 \\ (0.001)$	0.016	0.014 (0.001)	
Pr of switching from other providers	eq. (G2)	0.058	$\begin{array}{c} 0.065 \\ (0.003) \end{array}$	0.056	$\begin{array}{c} 0.067 \\ (0.002) \end{array}$	0.067	$\begin{array}{c} 0.067 \\ (0.002) \end{array}$	0.046	$\begin{array}{c} 0.049 \\ (0.002) \end{array}$	
Fee autocorrelation ratio	eq. (G3)	0.604	$\begin{array}{c} 0.599 \\ (0.042) \end{array}$	0.716	$\begin{array}{c} 0.721 \\ (0.025) \end{array}$	0.788	$\begin{array}{c} 0.788 \\ (0.024) \end{array}$	0.819	$\begin{array}{c} 0.819 \\ (0.024) \end{array}$	
Pr of no menu update	eq. (G4)	0.336	$\begin{array}{c} 0.333 \\ (0.004) \end{array}$	0.272	0.272 (0.004)	0.245	$\begin{array}{c} 0.247 \\ (0.003) \end{array}$	0.270	$\begin{array}{c} 0.270 \\ (0.003) \end{array}$	
Sensitivity of switching wrt fee	eq. (G5)	0.095	$\begin{array}{c} 0.119 \\ (0.016) \end{array}$	0.104	$0.098 \\ (0.015)$	0.086	$\begin{array}{c} 0.091 \\ (0.015) \end{array}$	0.059	$\begin{array}{c} 0.057 \\ (0.013) \end{array}$	
Fee changes from switching	eq. (G6)	-0.143	-0.142 (0.009)	-0.091	-0.094 (0.008)	-0.076	-0.076 (0.007)	-0.044	-0.044 (0.006)	
Variance of fee	eq. (G7b)	0.087	0.088 (0.002)	0.053	$\begin{array}{c} 0.053 \\ (0.001) \end{array}$	0.037	$\begin{array}{c} 0.037 \\ (0.001) \end{array}$	0.025	$\begin{array}{c} 0.025 \\ (0.001) \end{array}$	

Appendix Table A13: Fits of Auxiliary Model Coefficients, No Unobserved Production Costs

Table displays auxiliary model coefficients estimated using data generated by the structural model and coefficients estimated using observed data, when I directly fit variance of residualized fees. Standard errors of coefficients estimated using observed data are reported in parenthesis. These standard errors are bootstrapped using 250 samples with replacement.

Employer quartiles		1	st	2	2nd		3rd		4th	
		Coef	SE	Coef	SE	Coef	SE	Coef	SE	
RFP cost bps \$ thousand	$\kappa_{rfp}/\bar{\alpha}$	2.1 $1$	(0.0)	1.9	(0.0)	1.9 7	(0.2)	$0.2 \\ 5$	(0.0)	
Switching cost bps \$ thousand	$\kappa_{sw}/\bar{\alpha}$	$\begin{array}{c} 13.4\\9\end{array}$	(0.0)	$\begin{array}{c} 10.1 \\ 16 \end{array}$	(0.1)	$8.4\\30$	(0.3)	$\begin{array}{c} 1.4\\ 31 \end{array}$	(0.1)	
Fee sensitivity: mean Coefficient Elasticity	$\bar{\alpha}$	$\begin{array}{c} 0.57\\ 8.4 \end{array}$	(0.00)	$0.94 \\ 12.2$	(0.01)	$1.21 \\ 13.7$	(0.06)	$\begin{array}{c} 8.43 \\ 60.0 \end{array}$	(0.82)	
Fee sensitivity: st.dev. Coefficient Elasticity	$\sigma_{lpha}$	$\begin{array}{c} 0.74 \\ 6.5 \end{array}$	(0.00)	$1.23 \\ 9.8$	(0.01)	$1.59 \\ 10.7$	(0.06)	$11.05 \\ 62.9$	(2.04)	
RFP attention Coefficient	ρ	0.71	(0.01)	0.50	(0.02)	0.49	(0.10)	0.01	(0.00)	
Pr menu update	$\phi$	0.49	(0.00)	0.60	(0.00)	0.66	(0.00)	0.68	(0.00)	
Net quality of other providers Relative to other providers	$\xi/ar{lpha}$	-2.8	(0.0)	-6.3	(0.0)	-6.1	(0.2)	-1.0	(0.1)	
ADP		-0.3		4.5		3.6		0.7		
Principal Fin.		-1.5		3.8		4.5		0.8		
Empower		1.1		5.4		5.3		0.9		
Fidelity		1.1		5.4		5.8		1.0		
Vanguard		-1.2		4.0		4.5		0.8		

Appendix Table A14: Structural Parameter Estimates, Myopic

Table reports parameter estimates when I assume providers and employers are myopic. Standard errors are computed using the asymptotic variance formula, where the variance-covariance matrix of auxiliary model coefficients is computed using the bootstrap method over 250 iterations. In each iteration, I sample employers with replacement and estimate auxiliary model coefficients.

Employer quartiles		1st		21	2nd		3rd		4th	
		Model	Data	Model	Data	Model	Data	Model	Data	
Pr of switching from large providers	eq. (G1)	0.044	$\begin{array}{c} 0.030 \\ (0.002) \end{array}$	0.025	0.024 (0.001)	0.018	0.018 (0.001)	0.014	$\begin{array}{c} 0.014 \\ (0.001) \end{array}$	
Pr of switching from other providers	eq. (G2)	0.048	$0.065 \\ (0.003)$	0.068	0.067 (0.002)	0.067	$\begin{array}{c} 0.067 \\ (0.002) \end{array}$	0.049	$\begin{array}{c} 0.049 \\ (0.002) \end{array}$	
Fee autocorrelation ratio	eq. ( <b>G3</b> )	0.599	$\begin{array}{c} 0.599 \\ (0.042) \end{array}$	0.721	$\begin{array}{c} 0.721 \\ (0.025) \end{array}$	0.787	$\begin{array}{c} 0.788 \\ (0.024) \end{array}$	0.819	$\begin{array}{c} 0.819 \\ (0.024) \end{array}$	
Pr of no menu update	eq. (G4)	0.333	$\begin{array}{c} 0.333 \\ (0.004) \end{array}$	0.272	$\begin{array}{c} 0.272 \\ (0.004) \end{array}$	0.246	$\begin{array}{c} 0.247 \\ (0.003) \end{array}$	0.270	$\begin{array}{c} 0.270 \\ (0.003) \end{array}$	
Sensitivity of switching wrt fee	eq. (G5)	0.121	$\begin{array}{c} 0.119 \\ (0.016) \end{array}$	0.097	$0.098 \\ (0.015)$	0.092	$\begin{array}{c} 0.091 \\ (0.015) \end{array}$	0.057	$\begin{array}{c} 0.057 \\ (0.013) \end{array}$	
Fee changes from switching	eq. (G6)	-0.141	-0.142 (0.009)	-0.094	-0.094 (0.008)	-0.076	-0.076 (0.007)	-0.044	-0.044 (0.006)	
Variance of fee changes from switching	eq. (G7)	0.090	$\begin{array}{c} 0.079 \\ (0.003) \end{array}$	0.066	$0.066 \\ (0.002)$	0.047	$\begin{array}{c} 0.048 \\ (0.002) \end{array}$	0.029	$\begin{array}{c} 0.034 \\ (0.002) \end{array}$	
Positive fee changes from switching	eq. (G8)	0.172	$\begin{array}{c} 0.179 \\ (0.008) \end{array}$	0.162	$0.163 \\ (0.006)$	0.136	$\begin{array}{c} 0.136 \\ (0.006) \end{array}$	0.121	$\begin{array}{c} 0.118 \\ (0.006) \end{array}$	

Appendix Table A15: Fits of Auxiliary Model Coefficients, Myopic

Table displays auxiliary model coefficients estimated using data generated by the structural model and coefficients estimated using observed data, where I assume providers and employers are myopic. Standard errors of coefficients estimated using observed data are reported in parenthesis. These standard errors are bootstrapped using 250 samples with replacement.





Figure illustrates model generated markups for employers in each size quartile. Panel (a) separates plan fees into markups and production costs. Panel (b) decompose markups into components set by providers at RFPs and growth in markups between RFPs. Panel (c) plots distribution of markups across states. Panel (d) shows the decomposition of fee dispersion.

Appendix Table A16: Plan Fees and Transaction Costs under Counterfactual Policies, No Unobserved Production Costs

	Pl	an fees	Transaction cost		
Counterfactual	(bps)	(\$ billion)	(bps)	(\$ billion $)$	
Status quo	64.5	5.80	10.2	0.29	
Employers ignore transaction costs	68.8	5.49	103.1	4.11	
RFP/switch that minimize fees	59.6	5.28	19.9	1.16	
High RFP attention $(\rho = 0.95)$	63.0	5.41	10.0	0.30	
Consolidate plans, production cost only	58.0	5.65	10.2	0.29	
Consolidate plans	48.7	5.40	0.5	0.05	
High fee sensitivity (90th pct)	51.0	4.93	19.0	0.53	

Table reports outcomes under the status quo and counterfactual policies. I report both average plan fees or transaction costs measured in basis points across employers and also the aggregated dollar amount of annual plan fees or transaction costs. I directly fit variance of residualized fees assuming no unobserved production costs within employer size quartile.

Appendix Figure A9: Plan Fees and Transaction Costs under Counterfactual Policies, No Unobserved Production Costs



Figure plots counterfactual plan fees across different transaction costs. The y-axis corresponds to average plan fees, and the x-axis corresponds to employers' transaction costs. I directly fit variance of residualized fees assuming no unobserved production costs within employer size quartile.