#### NBER WORKING PAPER SERIES

# INVESTING IN LENDING TECHNOLOGY: IT SPENDING IN BANKING

Zhiguo He Sheila Jiang Douglas Xu Xiao Yin

Working Paper 30403 http://www.nber.org/papers/w30403

NATIONAL BUREAU OF ECONOMIC RESEARCH 1050 Massachusetts Avenue Cambridge, MA 02138 August 2022, Revised October 2023

For helpful comments, we thank Doug Diamond, Francesco D'Acunto (discussant), Isil Erel (discussant), Mark Flannery, Yueran Ma, Jun Pan (discussant), Nitzan Tzur-Ilan (discussant), João Granja, Nicola Pierri, Manju Puri, and Raghu Rajan. We are thankful for valuable comments from seminar and conference participants at the Richmond Fed, University of Connecticut, University of Florida, University of Chicago, City University of New York, 5th IMF Annual Macro-Financial Conference, EFA, SFS Calvacade, and the 22nd Annual Conference in Digital Economics, and AFA 2023 in New Orleans. Zhiguo He acknowledges financial support from the John E. Jeuck Endowment at the University of Chicago Booth School of Business. All errors are our own. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2022 by Zhiguo He, Sheila Jiang, Douglas Xu, and Xiao Yin. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Investing in Lending Technology: IT Spending in Banking Zhiguo He, Sheila Jiang, Douglas Xu, and Xiao Yin NBER Working Paper No. 30403 August 2022, Revised October 2023 JEL No. G21,G51,O12,O32

#### **ABSTRACT**

Banks' lending technology hinges on their handling of soft and hard information in dealing with different types of credit demand. Through assembling a novel dataset on banks' investment in information technologies (IT), this paper provides concrete empirical evidence on how banks adapt their lending technologies. We find investment in communication IT is associated with improving banks' ability to produce and transmit soft information, while investment in software IT helps enhance banks' hard information processing capacity. We exploit policies that affect geographic regions differentially to show causally that banks respond to an increased demand for small business credit (mortgage refinance) by increasing their spending on communication (software) IT spending. We also find that the entry of fintech induces commercial banks to increase their investment in IT—more so in the software IT category.

Zhiguo He University of Chicago Booth School of Business 5807 S. Woodlawn Avenue Chicago, IL 60637 and NBER zhiguo.he@chicagobooth.edu

Sheila Jiang
Department of Finance and Real Estate
Warrington College of Business
University of Florida
P.O. Box 117168
Gainesville, FL 32611
Sheila.Jiang@warrington.ufl.edu

Douglas Xu University of Florida Warrington College of Business 1454 Union Rd Stuzin Hall 315 Gainesville, FL 32611 douglas.xu@warrington.ufl.edu

Xiao Yin
University College London
30 Gordon St, London
Drayton House
London WC1H 0AX
United Kingdom
xiao.yin@ucl.ac.uk

## 1 Introduction

Commercial banks have long relied on cutting-edge technology to deliver innovative products such as ATMs and online banking, streamline loan making processes, and improve back-office efficiency. According to a 2012 Mckinsey Report, across the globe commercial banks spend about 4.7% to 9.4% of their operating income on information technology (IT); for comparison, insurance companies and airlines only spend 3.3 percent and 2.6 percent of their income on IT, respectively. Recently, the impact of information technology on the banking sector and on financial stability has been a headline topic in policy discussions (Banna and Alam, 2021; Pierri and Timmer, 2020).

Although the financial services industry—especially the banking industry—is increasingly becoming a tech-like industry, the academic literature lags behind in understanding the economics of IT spending in banking. Which banks, large or small, have invested more in IT? Do banks adapt their information technologies in response to different credit demand shocks? How do traditional banks react to the entry of fintech in recent years? We take the first step toward understanding the key empirical patterns on these issues, and further explore the mechanism that underlies the connection between these expenditures and the core functioning of banking.

To place our research in the literature, think about the information transmission between a loan officer and a borrower, or between layers of loan officers within a bank organization. As highlighted by Stein (2002), a less hierarchical structure within a bank facilitates the effective transmission of "soft" information. At the same time, fast-developing technologies in recent decades provide more options for the banking sector to cope with such problems. So, can information technologies reduce frictions in communicating soft information and potentially improve banks' credit approval decisions? Likewise, with the explosive development of big data analytics which combine "hard" information such as credit scores with other alternative

<sup>&</sup>lt;sup>1</sup>For instance, First Citizens National Bank implemented its employee intranet to strengthen internal communications in February 2019. For details, see this article.

data, have traditional banks started adopting these technologies?

Our study relies on a comprehensive dataset, the Harte Hanks Market Intelligence Computer Intelligence Technology database, which has been used in the literature on the economic implications of technology in non-financial sectors (e.g., Bloom et al., 2014; Forman et al., 2012). This dataset, which aligns well with the regulatory Y-9C dataset in measuring total IT spending, provides detailed branch-level information on specific spending categories. We focus on two major categories in banks' IT spending: 2 software and communication. Software IT products mainly aim to improve information processing accuracy and speed through automation, specialized programming and AI technologies, etc. Communication IT products facilitate smoother exchanges of information within bank branching networks, across banks, and with borrowers.

In Section 3 we start by documenting that IT expenditure in the U.S. banking sector has been growing rapidly over the last decade. Growth in IT spending varies by bank size: large banks' IT spending increased steadily, while there was almost no growth in IT spending for the smallest banks. Another noticeable distinction between large and small banks is that the latter, who presumably engage in more small business lending, consistently allocate a higher share of their IT budget towards communication technology than the former. As we will elaborate, this pattern points to the role communication IT plays in conducting small business lending.

We then examine the relationship between banks' IT spending and their lending activities. Among the three main categories of loan types in Call Reports, the share in commercial and industrial (C&I) loans is positively associated with the lenders' communication spending, but uncorrelated with their software spending. In contrast, the share of personal loans is positively associated with the lenders' software spending, but not with communication

<sup>&</sup>lt;sup>2</sup>Section 2.2 explains in detail the four major categories of IT expenditure in the Harte Hanks data set—hardware, software, communication, services—in the context of the banking industry. Representative examples of software include desktop applications (e.g., Microsoft Office), information management software, and risk and payment management software. Examples of communication technology include radio and TV transmitters, private branch exchanges, video conferencing, etc.

spending. Going one step further, within C&I loans, we show that small business lending stands out as a sub-category that drives the overall positive association with communication IT spending, whereas within personal loans, mortgage refinance is the main contributor to the positive correlation between personal loans and software spending. As different types of loans often require different technologies in dealing with relevant information, these positive associations (or the lack thereof) offer important guidance in understanding banks' IT spending profiles from the perspective of lending technology.

Aside from broad credit categories of loan portfolios, we also explore how banks' IT investment is shaped by other factors affecting their business operations. Regarding the complexity of internal hierarchical structure, banks with more internal layers tend to have a higher communication spending. Further, hierarchical complexity has an impact on the responsiveness of banks' IT spending to their loan profiles—a more complex hierarchical structure makes banks' communication spending respond more to their small business lending, but displays no systematic effect on the relation between software spending and mortgage refinancing.<sup>3</sup> Finally, in the context of the syndicated lending market, frequent lead lenders spend significantly more on communication than participant lenders, as lead banks take more direct responsibility in interacting with borrowers.

In Section 4 we delve deeper into the underlying economics behind the connection between banks' IT investment and their lending activities. Conceptually, we distinguish two fundamentally different types of lending technologies. The first heavily relies on the gathering and augmentation of soft information from borrowers; in the context of Berger and Udell (2002), "relationship lending" is a concrete example of the first type. The second type of lending technology relies primarily on the processing and quantification of hard information; leading examples include "transactions lending," i.e., loans that are based on a specific credit scoring system and quantified financial statement metrics (Berger and Udell, 2006).

<sup>&</sup>lt;sup>3</sup>This asymmetric pattern is consistent with the notion of "hierarchical friction" in Stein (2002): A lower level of hierarchical complexity helps facilitate the within-organization transmission of soft information, which is more relevant for small business lending than mortgage refinancing.

We formulate our first hypothesis along the dimension of soft information. Increased demand for loans that involve intensive soft information production/transmission (e.g., small business loans) should lead banks to invest more in communication technologies, say video conferencing, as they not only enable banks to more effectively collect soft information from entrepreneurs but also allow for a smoother transmission of such otherwise hard-to-verify soft information within a bank organization. Taking advantage of an arguably exogenous demand shifter, we show an increase in banks' small business credit demand—due to a higher ex-ante exposure of local counties to the policy shock exploited in our analysis—leads to a positive and significant growth in banks' communication spending, without much impact on the bank's software spending.<sup>4</sup> Furthermore, we find that such responses are more pronounced for banks operating in regions with more young firms, for whom effective transmission of soft information is particularly important.

In our second hypothesis, a positive demand shock for loans that rely heavily on hard information processing (e.g., mortgage refinancing) should push banks to engage in more IT investment in software (which facilitates lenders in processing such existing data). For causal identification, we utilize the cross-county variation in interest savings of outstanding mortgages to construct a shifter to the mortgage refinance demand across different regions.<sup>5</sup> We show that an increase in mortgage refinancing by a bank (due to its local exposure to high refinance savings) results in a higher software spending, without any significant impact on its communication spending.

The last part of our analysis concerns how the entry of fintech lenders into local credit markets affects banks' IT spending. In the past decade, we have witnessed a growing pene-

<sup>&</sup>lt;sup>4</sup>Our empirical identification is based on the "Small Business Health Care Tax Credit," which was introduced in 2010 and then experienced a significant policy change in 2014. We construct the instrumental variable using counties' exposure to policy change in 2014, with details explained in Section 4.2.2.

<sup>&</sup>lt;sup>5</sup>We take advantage of the low-interest episode from 2011 to 2015, during which nationwide average mortgage interest rates decreased from 6.5% to 3.5%. When interest rates drop, the mortgage prepayment option is in the money (Eichenbaum et al., 2022; He and Song, 2022), implying a greater mortgage refinance demand by local households. In addition to mortgage repayment saving, we also consider an alternative instrumental variable (IV) constructed as the weighted mortgage rate gap for outstanding mortgages in a given county.

tration of fintech into the traditional banking sector. Utilizing the staggered entry of Lending Club into seven states after 2010 as an experimental setting, we investigate how the traditional banking sector reacted to the penetrating fintechs. Right after the regulatory approval of Lending Club's operation in a state, banks operating in that state saw a significant increase in their IT investment. Importantly, shocked banks' software spending experienced an economically and statistically significant growth (around 7%), whereas the change in their communication spending was insignificant.

There is also significant heterogeneity across bank size groups in their technology spending reactions in response to fintech entry. In particular, the increased IT investment is predominantly observed among large banks, whereas small banks barely respond. Our findings suggest an overall "competition reaction" from the traditional banking sector, in that banks—particularly larger ones—are catching up with fintech challengers. Consistent with this competition interpretation, such "catching up" behavior by commercial banks is especially noticeable in improving their automating and information processing technology through increased software spending, which is precisely the domain of lending technology in which fintech lenders have a comparative advantage.

#### Related Literature

Bank lending technology and the nature of information. Berger and Udell (2006) provide a comprehensive framework of the two fundamental types of bank lending technology, i.e., relationship lending and transactions lending, in the SME lending market; see also related work by Bolton et al. (2016). A key difference between these two types of lending is related to the role played by information as highlighted by Stein (2002), who provides an explanation for why soft information production favors an organizational structure with

<sup>&</sup>lt;sup>6</sup>Relatedly, consistent with large banks' dominance in the personal loan lending market, we also find that the response to Lending Club's entry is particularly pronounced for banks more specialized in personal loan lending.

fewer hierarchical layers.<sup>7</sup>

We contribute to this literature by linking banks' IT spending to their lending technology, especially with regard to the distinction between soft information production/transmission and hard information processing. We further establish causal linkages from the informational components in credit demand to banks' IT spending. It is, to our knowledge, the first attempt in the literature to show how credit demand shocks drive banks' investment in their information-driven lending technologies.<sup>8</sup>

Information technology in the banking industry. Our paper belongs to the literature on the interaction between the development of information technology and the evolution of the banking industry. For instance, Berger (2003) shows that progress in both information and financial technologies led to significant improvement in banking services, and Petersen and Rajan (2002) document that communication technology greatly increased the lending distance of small business loans. Using the number of computers per employee as a measurement for IT adoption, two recent papers show that IT adoption helps banks weather financial crisis (Pierri and Timmer, 2022) and spur entrepreneurial activities (Ahnert et al., 2021). Our paper, with the aid of detailed IT spending data, studies the specific economic mechanisms that connect banks' lending technology with their IT spending.<sup>9</sup>

Our paper is closely related to Modi et al. (2022) who construct IT spending data using the Call Reports; we compare our sample with them in Section 2.1. While their analysis also investigates the linkages between banks' IT spending and their lending behaviors (e.g., mortgage lending and reactions to monetary policies), our analysis differs in our focus on

<sup>&</sup>lt;sup>7</sup>Along these lines, Liberti and Mian (2009) find empirically that greater hierarchical distance leads to less reliance on subjective information and more on objective information. Paravisini and Schoar (2016) document that credit scores, which serve as "hard information," improve the productivity of credit committees, reduce managerial involvement in the loan approval process, and increase the profitability of lending.

<sup>&</sup>lt;sup>8</sup>Previous literature has shown that credit supply positively affects non-financial firms' technology adoption or innovation (Amore et al. (2013), Chava et al. (2013), Bircan and De Haas (2019)).

<sup>&</sup>lt;sup>9</sup>There is also a vast theoretical literature (Hauswald and Marquez, 2003, 2006; Vives and Ye, 2021) on the interactions among information technology and credit market competition; see Freixas and Rochet (2008) for a review. For instance, in the framework of credit market competition where the specialized lender acquires additional "soft" information, He et al. (2023a) study the role of information span where loan quality is determined by multi-dimensional characteristics.

linking different categories of IT investment to banks' lending activities associated with different types of information nature (i.e., soft information vs. hard information). Furthermore, our analysis aims to establish a causal linkage between banks' adaptation of their lending technology and credit demand shocks, which is not the focus of Modi et al. (2022).

Fintech entry and banks' IT spending. The emergence of fintech reflects the recent developments in information technologies. Our study aligns closely with studies on how the emergence of the fintech industry is affecting (or has affected) the traditional banking sector. While a common theme of this research has mostly focused on bank-fintech competition during which traditional banks are largely viewed as a passive player, little attention has been paid to the banks' active responses; we take the latter angle by studying whether and how traditional banks are catching up with penetrating fintech lenders. Along a similar line, Modi et al. (2022) also document that banks with more fintech exposure in the mortgage market tend to spend more on IT, and that their lending behaviors are also likely to resemble fintech companies.

Micro-level evidence on technology adoption. Our paper also broadly contributes to the literature studying firms' technology adoption behavior using micro-level data. Using the same IT spending data as this paper, Forman et al. (2012) study the impact of firms' technology adoption on regional wage inequality, Bloom et al. (2014) investigate the effect of information technology on firms' internal control, and Ridder (2019) explores how software adoption explains the decline in business dynamism and the rise of market power.<sup>12</sup>

<sup>&</sup>lt;sup>10</sup>Related works include but are not limited to Jagtiani and Lemieux (2017), Buchak et al. (2018a), Fuster et al. (2019), Frost et al. (2019), Hughes et al. (2019), Stulz (2019), and Di Maggio and Yao (2020).

<sup>&</sup>lt;sup>11</sup>This fast-growing literature includes Lorente et al. (2018); Hornuf et al. (2018); Calebe de Roure and Thakor (2019); Tang (2019); Erel and Liebersohn (2020); Aiello et al. (2020); Gopal and Schnabl (2022); He et al. (2023b), and Huang (2022).

<sup>&</sup>lt;sup>12</sup>While we use detailed IT "budget" data from Harte Hanks, several papers use its IT "installment" data that report firm-level IT product installment; see Section 2.1 for details of the differences between these two datasets. For instance, Charoenwong et al. (2022) study installment of IT products catering to compliance requirements, and Pierri and Timmer (2022) investigate whether banks installed with more PCs per employee can better survive a financial crisis. We use the budget data which reports detailed dollar amounts for various IT categories, which are crucial for our study.

## 2 Data and Background

We explain our main data sources in this section, together with detailed descriptions of various categories of IT spending.

## 2.1 Data Source for Bank IT Spending and Sample Construction

The data on banks' IT spending comes from the Harte Hanks Market Intelligence Computer Intelligence Technology database, which covers over three million establishment-level observations from 2010 to 2019 obtained while conducting IT-related consulting for firms. Harte Hanks collects and sells this information to technology companies, who then use it for marketing purposes or to better serve their clients. Firms have incentives to report their IT spending data truthfully to Harte Hanks as they also want to receive tailored advice for better IT services in the future.

Our paper focuses on commercial banks.<sup>13</sup> The sample consists of 1,450 commercial banks in the U.S., which covers more than 80% of the U.S. banking sector in terms of asset size (Figure A1). The sample is more representative for large banks, as shown in Table 1 which reports our coverage by bank asset size group. For the three groups of relatively large banks (with assets above \$1 billion), the coverage in frequency and assets are both over 80%. However, for small banks with size below \$100 million, our sample covers only 9.47% (11.36%) of the total number (assets) of commercial banks in the U.S. system.

Table 2 displays the summary statistics of banks' IT spending. In our sample, total IT spending as a share of net income ranges from 1.7% (25<sup>th</sup> percentile) to 8.5% (75<sup>th</sup> percentile), suggesting a large cross-sectional variation across banks. Median IT spending as a share of net income is 5.2%, consistent with a 2012 McKinsey survey (Figure A5) reporting that banks' IT spending as a share of net operating income ranges from 4.7% to 9.4%.

<sup>&</sup>lt;sup>13</sup>The Harte Hanks dataset has been utilized by a broad literature of economic studies. For instance, Forman et al. (2012) investigate firms' IT adoption and regional wage inequality; Bloom et al. (2014) study the impact of information communication technology on firms' internal control; and Tuzel and Zhang (2021) study labor-technology substitution at establishment level, based on this dataset.

Matching procedure and matching quality. We provide detailed descriptions of the matching algorithm for our sample construction in Appendix B.1.1. This matching mainly involves mapping sites in Harte Hanks to bank branches, where we take bank names from the Call Report. To evaluate the matching quality, we conduct several cross-checks of our sample with other data sources, especially Call Reports, regarding certain key bank-level variables. More specifically, Appendix B.1.2 shows a close alignment of our sample with the Call Report in terms of the total number of banks' branches at both the bank level and bank-county level, as well as banks' total revenue and total number of employees.

Cross-validation of bank IT expense measure. To confirm the reliability of Harte Hanks data in measuring banks' IT spending, we conduct thorough validity checks against various alternative sources, including the "FR Y-9C Consolidated Financial Statements for Holding Companies" which contain regulatory data on IT spending at the bank holding company (BHC)-year level. As explained in Appendix B.2, we follow Kovner et al. (2014) to construct BHCs' IT expenses by summing up two standardized "other noninterest expenses" in the Y-9C dataset, the "Data processing expenses" and the "Telecommunication," together with the unstandardized write-in items reported in "other noninterest expenses" containing IT-related keywords. For the top 50 BHCs sorted by assets, Figure B9 offers a BHC-by-BHC comparison of IT spending between Y-9C and Harte Hanks (adjusted for BHC subsidiaries), at years that Y-9C IT expenses are not missing or zero. <sup>14</sup> Further, for the overall matched sample, Table B4 reports that a regression of the logarithm of the IT spending in Harte Hanks on that in Y-9C has a slope coefficient close to one (0.935) and a constant close to zero (0.037). Overall, we find a decent match between Y-9C regulatory filings and Harte Hanks data (adjusted for BHC subsidiaries), suggesting a high quality of the Harte Hanks data in measuring banks' IT expenses.

<sup>&</sup>lt;sup>14</sup>An important feature of the regulatory Y-9C and Call Report data is that the reporting of banks' IT expenses is censored at certain threshold, i.e., IT expenses below that threshold are often reported as zero or missing values; this contributes to a higher IT spending measure in our sample than either Kovner et al. (2014) or Modi et al. (2022). More details regarding the reporting rules of these regulatory data are provided in Data Appendix B.2.1.

We also compare our IT spending measure with that in Modi et al. (2022), who construct the IT spending using only IT-related write-in expenses in Call Reports. We find a decent correlation ( $\rho = 0.77$ ) between the two IT expense measures. Furthermore, Figure B12 separates banks into different size groups as in Modi et al. (2022) and shows similar time trends in banks' IT expenses during 2010-2019 across bank size groups in these two data sets.<sup>15</sup> Data Appendix B.3 provides several further comparisons with our sources on the empirical measures for banks' IT investment, from which overall consistent results are obtained.<sup>16</sup>

Data collection practice by Harte Hanks. Our analysis uses the "IT budget data" offered by Harte Hanks. According to the official description provided by the data collection team of Harte Hanks, the construction of this data set is mainly based on data collected from surveys conducted at the site-year level; in addition, the IT budget data reflects the purchases of ready-to-use IT products or services, and does not include expenditures on IT-related R&D activities which might take a longer time to accomplish. Furthermore, since the usage of IT products or services (e.g., software programs) is often licence-based, the IT expense is therefore likely to be spread across branches of a given bank based on branch-level usage, rather than concentrated at the headquarters of the bank. As one piece of supporting evidence, we find no significant differences between the IT expenses at bank headquarters and local bank branches. Data Appendix B.4 provides more details for supplemental materials and analysis regarding the IT data collection practices of Harte Hanks.

<sup>&</sup>lt;sup>15</sup>It is worth noting that our IT expense measure, which aligns well (in dollar amounts) with the aggregated IT spending in Y-9C following the method in Kovner et al. (2014), is systemically higher than that in Modi et al. (2022). This is because Kovner et al. (2014) sum the two standardized "other noninterest expenses" and unstandardized write-in items, while Modi et al. (2022) only account for the unstandardized write-in items of "other noninterest expenses." Furthermore, the censored reporting in Call Report (and in Y-9C) also contributes to a smaller IT expense measure in Modi et al. (2022) compared to the measure offered by Harte Hanks.

<sup>&</sup>lt;sup>16</sup>These comparisons include i) the industry-level Bureau of Economic Analysis (BEA) data, with the focus being the detailed composition of bank' IT spending; and ii) the data storage cost in Feyen et al. (2021).

<sup>&</sup>lt;sup>17</sup>The Harte Hanks data set has two major parts—the "IT installment" data that contains information on whether a firm installs certain IT product and the earliest installment date, and the "IT budget data" which contains information on the dollar amounts of detailed categories of IT investment. While the data vendor gets some of the installment information using algorithms extracting installment information from firms' public reports, job postings, etc., the IT budget data is based primarily on surveys. Our analysis in this paper only uses the IT budget data by Harte Hanks.

## 2.2 IT Investment Categorization

Our dataset offers a detailed decomposition of banks' IT investments in four major categories specified by Harte Hanks: *hardware*, *software*, *communication*, and *services*. We now explain these categories, with formal definitions given in (a) to (d) of Figure A6.

Software is defined as software programs purchased from third parties, including those offered as an SaaS from a multi-tenant shared-license server accessible by a browser. More specifically, the category of software covers desktop applications, information management software, processing software, and risk and payment management software. For desktop applications, one representative example is Microsoft Office. Processing software specializes in automatically processing information from loan applicants' document packets through specialized programming and AI technologies with improved accuracy and speed, which would otherwise be done manually by loan officers. Risk management software provides ongoing risk assessment after loans have been issued, through augmenting borrowers' repayment status as well as real-time industrial and economic conditions. 20

Communication is defined as the network equipment that banks operate to support their communication needs. It includes routers, switches, private branch exchanges, radio and TV transmitters, Wi-Fi transmitters, desktop telephone sets, wide-area networks, local-area network equipment, video conferencing systems, and mobile phone devices. For effective project evaluation, these machines allow bankers to conveniently talk to and see borrowers who seek credit. In addition, communication equipment such as private branch exchanges facilitates the exchange of information, opinions, and decisions within the bank branching

<sup>&</sup>lt;sup>18</sup>These software products are easy to grasp by bank employees who are then able to conduct basic calculations and visualizations of data associated with lending businesses. For example, on Mendeley.com, the job postings for loan officers or project managers by many banks require applicants to be proficient with Microsoft Office.

<sup>&</sup>lt;sup>19</sup>Examples of processing software include Trapeze Mortgage Analytics, Treeno Software, and Kofax. These software products feature document assembly enhancement, digitization, and information classification.

<sup>&</sup>lt;sup>20</sup>These software products, e.g. Actico, ZenGRC, Equifax, and Oracle ERP, allow banks to better monitor loans in progress. Other software products include security trading systems and operating systems that are typically bundled with the specific software products.

networks.

Hardware as a form of IT investment includes classic computer hardware such as PCs, monitors, printers, keyboards, USB devices, storage devices, servers, and mainframes. In terms of lending services, hardware investment can complement and facilitate both the gathering of borrower information and the processing of that information. This is because hardware devices, such as PCs and servers, help provide storage and transmission of data, in addition to serving as the carriers of software and toolboxes.

Services are defined as project-based consulting services (including, say, IT strategy and security assessments) or systems integration services that vendors provide to banks, which are often provided by IT outsourcing companies on a contractual basis. Similar to hardware, services work as complements to other categories of information technology investment to facilitate banks' lending. Examples include Aquiety, a Chicago-based IT service company that provides cybersecurity services to banks and other firms; and Iconic IT, a New-York based IT service company that provides software and hardware procurement, together with installment and upgrade services.

Table 2 reports summary statistics on the detailed structure of banks' IT spending profiles. By size, software and services are the top two among all categories of IT spending, each constituting around 33% of total IT budget; while hardware (communication) constitutes about 17% (9%) of total IT budget. We conduct analysis on banks' IT spending at bank-year level in Section 3, while in Section 4 the analysis is at bank-county-year level, in which we aggregate the branch-level spending information of each bank at the county level.

#### 2.3 Other Datasets

To supplement our study on banks' lending technologies and their relation to IT spending, we combine loan-level information from multiple sources.

Bank balance sheet. We obtain bank-level balance sheet information from Call Reports; for detailed matching procedures, see Appendix C.1.1. In our OLS analysis investigating the correlation between IT spending and loan portfolios, we calculate the dependent variable as IT spending in Harte Hanks as a share of "Revenue" (RIAD4000 of Call Report) at bank-level. In the identification part that involves bank-county analysis (Section 4 and 5), since there is no bank-county revenue in regulatory data, we use "revenue" scaled by number of employees in Harte Hanks to control for the profitability at the bank-county level.

Loans and local characteristics. We obtain syndicated loan information on the frequency of a bank acting as lead bank in syndicated loan packages from LPC Dealscan. Small business loan origination data are from the Community Reinvestment Act (CRA), which is at the bank-county-year level covering the sample period of 2010–2019. Mortgage refinance information is available through the Home Mortgage Disclosure Act (HMDA) from 2010–2019, and we use the county-level average mortgage interest rate before 2010 obtained from Freddie Mac as the demand shifter for mortgage refinancing.

Bank hierarchical structure. We obtain banks' hierarchical structure information from Mergent Intellect platform, which covers 97 million public and private businesses including their locations and industry classifications.<sup>21</sup> We restrict our sample to entities with the two-digit SIC code of "60," which designates "Depository Institutions."

The database provides the complete family trees of the companies, with detailed information on its family members. Importantly, this database classifies each family member of a company into one of the three categories of location types: "Headquarters," "Single Location," and "Branch." We define a bank as having n layers of hierarchical structure if the bank has n types of locations in the family tree, where  $n \in \{1, 2, 3\}$ . To give some concrete examples, Wells Fargo has all three location types and hence is classified as three hierarchical layers; North Valley Bank with headquarters located in Corning (OH) and seven branches is

<sup>&</sup>lt;sup>21</sup>Huvaj and Johnson (2019) use this database to study the impact of firms' organizational structure on their innovation activities.

classified as two layers; and First Place Bank located in Warren (OH) with one single location is classified as only one layer. For each bank, we match the banks in Mergent Intellect with banks in our sample based on bank names and the city where the banks' headquarters are located (see Appendix C.1.1 for more details). While the number of distinct "location types" in the Mergent Intellect dataset can provide information on the hierarchical complexity of a bank, it is admittedly a somewhat coarse empirical measure and could underestimate the hierarchical complexity, especially for large banks.

# 3 Empirical Patterns of Banks' IT Spending

We start our analysis by reporting some basic statistics of banks' investment in IT over the last decade as well as across bank size. We further show that banks' IT investment is related to their lending activities and organization structures.

## 3.1 Banks' IT Investment: Trend and Cross-Section

In Panel a) of Figure 1, we plot the evolution of IT spending as a share of total revenue of banks in our sample, in the manufacturing sector (2-digit SIC "20-39"), and in all other non-depository sectors (2-digit SIC "not 60"). As is evident from the figure, the IT investment in the U.S. banking sector (represented by our sample) has witnessed faster growth during the past two decades compared to other industries. Also, banks' IT spending saw faster growth starting in 2016, which could be potentially driven by the release of a "white paper" by the Office of the Comptroller of the Currency on March 16, 2016. This white paper set forth the regulators' perspective on supporting responsible innovation across all sizes of banks, 22 which might have pushed banks to be more aggressive in embracing technology investment as part of their strategic planning (see this article).

<sup>&</sup>lt;sup>22</sup>In this white paper, the Office of the Comptroller of the Currency defines "responsible innovation" as "the use of new or improved financial products, services, and processes to meet the evolving needs of consumers, businesses, and communities in a manner that is consistent with sound risk management and is aligned with the bank's overall business strategy."

Banks of different sizes often behave differently in systematic ways. In Panel b) of Figure 1, we follow FDIC bank size classifications to break banks into five size groups, and present the growth trend for banks' IT investments for each group as a fraction of non-interest expenses.<sup>23</sup> We observe that large banks invest more in IT than their small peers do (for more detailed statistics, see Table A2);<sup>24</sup> and IT spending in large banks (with asset size \$10–250 billion and above \$250 billion) has been steadily growing, though there is also an apparent growth in the smaller groups (with asset size below \$10 billion).<sup>25</sup> The study in Modi et al. (2022) also confirms the empirical pattern that larger banks tend to invest more in IT and experience higher rates of growth in IT spending. While we do not have a conclusive answer for why such heterogeneity exists, our analysis of how banks (of different sizes) react to the entry of fintech in Section 5 touches on this issue directly.

Another noticeable pattern in Table A2 is that small banks tend to allocate a higher fraction of their IT budget towards communication technology than large banks do: the average communication spending over total spending decreases from 12.6% for the smallest group (below \$100 million) to 6.2% for mega banks (above \$250 billion). For software spending, however, there are no significant differences across bank size groups. We will come back to this contrast in Section 4, where we connect banks' IT spending categories to their lending activities that involve different ways of handling information.

<sup>&</sup>lt;sup>23</sup>The magnitude of IT budget as a share of non-interest expenses in this figure is also in line with Hitt et al. (1999), who report banks' IT spending could be as high as 15% of non-interest expenses in their survey. The trend of IT spending as a share of total revenue, as is shown by the solid line in Figure 1, shares a consistent pattern with IT spending as a share of non-interest expenses.

<sup>&</sup>lt;sup>24</sup>Table A2 shows a robust cross-sectional pattern that IT spending (say, scaled by non-interest income) increases with bank size. This could be due to the fact that small banks often cannot afford IT purchases with significant lump-sum costs.

<sup>&</sup>lt;sup>25</sup>Medium- and smaller-size banks (asset size bins bellow \$10 billion) saw growth from 2010 to 2013, slowed down in 2015, and then have picked up again since 2016. One possible explanation for the temporary slowdown in IT spending in 2015 is that banks chose to "wait and see" in 2015 before the release of the white paper in 2016 (see first paragraph in Section 3.1).

## 3.2 Empirical Patterns of Bank IT Investment

We now present the first set of empirical results that relate banks' IT investment to their operations, from three angles: 1) specialization in loan making; 2) the role that a bank plays in syndicated loans; and 3) the complexity of bank's internal hierarchical structure.

#### 3.2.1 Loan Specialization

Banks provide three major types of loans: commercial and industrial (C&I) loans, personal loans, and agricultural loans. Lending to different types of borrowers often involves distinct ways of dealing with borrower-specific information. Therefore, if banks specialize in different types of loans, one should expect them to differ in their IT investment profiles. Specifically, we run the following bank-level regression:

$$\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,10\text{-}19} = \alpha_i + \beta \frac{\text{Type L loan}}{\text{Total loan}}_{i,10\text{-}19} + \gamma \mathbf{X}_i + \epsilon_i. \tag{1}$$

Here, the outcome variable of interest is  $\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,10\text{-}19}$ , which is the average investment in a specific type of IT spending as a share of bank i's revenue over 2010-2019. The "Revenue" in the denominator of the dependent variables is the total income in Call Report (RIAD4000). The main explanatory variable  $\frac{\text{Type L loan}}{\text{Total loan }}_{i,10\text{-}19}$ , which is the average share of a specific type of loan relative to bank i's total loan size, captures bank i's loan specialization. Control variables, which are measured over 2010-2019 at the bank level, include net income, total deposits, total equity, total salaries (all scaled by total assets), and revenue per employee. In all of our regression analysis (except for Section 5, where the main independent variable is the dummy variable indicating post-entrance), we standardize the dependent variables and the regressors.

Table 3 reports the estimation results of (1) for C&I loans, together with detailed regression outcomes for control variables. We apply the same methodology for personal loans and

<sup>&</sup>lt;sup>26</sup>For robustness, we also conduct analyses with an alternative measure of banks' IT spending intensity scaled by banks' deposits, with qualitatively similar results (Table A3).

agricultural loans. For exposition purposes, Table 4 reports key regression coefficients (i.e., those of specific IT spending shares).

#### A. Commercial and Industrial (C&I) loans

Specialization in C&I loans is most positively associated with banks' spending in communication technology (Table 3 column 2). A one standard deviation (13 percentage points) increase in loan portfolio share allocated to C&I loans predicts a \$0.24 million higher expenditure on communication per year. For detailed calculation, see Table A6. Our economic magnitude calculation for Table 3 follows the same method.

A higher degree of specialization in C&I loans also predicts more spending on hardware (column 3), though the magnitude is slightly smaller than that of communication spending. The coefficient of software spending, however, is insignificant (column 1).

Within C&I loans. Rows 2-3 of Table 4 further decompose C&I loans into "Small Business Loans," which are measured by a bank's small business lending reported in the CRA, and "other C&I loans." While the share of small business loans in a bank's portfolio is positively associated with communication spending, it is negatively related to the bank's software spending. In contrast, "other C&I loans" (e.g., loans to large firms) are positively associated with software spending, but not with communication spending. Panel A in Table 5 further shows that this positive association between small business loans and communication spending is more statistically significant for small banks.

#### B. Personal loans

Row 5 of Table 4 reports the associations between shares in personal loans and banks' IT spending. Contrary to the pattern we observe for C&I loans, a higher share of loan portfolio allocated to personal loans appears to predict more spending on software only. Quantitatively, a one standard deviation increase in personal loans share (an increase of about 7 percentage points) predicts a 0.0617 standard deviation increase in software spending as a share of total revenue; in dollar terms, this amounts to an increase of \$1.53 million in

software spending per year. On the other hand, a higher personal loans share does not have qualitatively significant predictive power on communication, hardware, or services budgets.

Within personal loans. Paralleling our analysis of small business loans within C&I loans, we decompose personal loans into two subcategories: mortgage refinancing and everything else. It is mortgage refinancing—but not other kinds of personal loans—that positively correlates with banks' software spending. This finding motivates our study in Section 4 to pay particular attention to mortgage refinancing as a specific type of lending activity in which the processing of hard information plays a critical role.

Additionally, the richness of mortgage data allows us to gain further insights by distinguishing "refinancing an existing loan" from "originating a new loan," with results reported in Row 5 and 7 of Table 4. We postpone more detailed discussion to Section 4.3.

#### C. Agricultural loans

As shown in Row 8 of Table 4, the agriculture loan presence seems to be positively associated with all categories of IT spending, although there is no statistically significant correlation between agriculture loan proportion and a particular type of IT investment.

## 3.2.2 Complexity of Hierarchical Structure

Another important factor that may affect a bank's efficacy in handling information is the internal organization structure of a bank (Stein, 2002). In the first row of Panel B in Table 4, we use the measure of hierarchical layers defined in Section 2.3 as our main proxy for hierarchical complexity. When the number of banks' hierarchical layers increases, banks spend more across all IT categories, especially on communication. Increasing from 1 hierarchical layer to 3 layers implies \$0.31 million more in communication spending each year. This result is under the specification with bank-size group fixed effects included, implying that hierarchical complexity predicts higher communication spending beyond bank size.<sup>27</sup>

 $<sup>^{27}</sup>$ Recall that in Section 3.1 we show that smaller banks tend to allocate a larger portion of their IT budget to communication spending. Our finding therefore suggests that despite its high correlation with bank size, the complexity of a bank's internal hierarchical structure has an additional impact on its IT spending on

For robustness, we also proxy banks' hierarchical complexity using the logarithm of the total number of offices, with qualitatively similar results reported in the second row of Panel B.

As we will explain in Section 4.2.1, one can relate these findings to Stein (2002), from the perspective of within-organization transmission of information that is difficult to verify and relay. Despite a crude empirical measure of hierarchical complexity, our paper establishes a direct link between hierarchical complexity and banks' IT investment for information production and transmission.

#### 3.2.3 Role in Syndicated Lending

Aside from specialization in different types of loans or having different levels of hierarchical complexity, banks may also differ in the role they play in dealing with information when conducting lending. For instance, in the context of syndicated lending, lead lenders and participant lenders perform drastically different tasks. Table 4 Panel C presents the same regression as in Eq. (1), except we replace the right-hand side variable with "%Lead bank/Total syndicate," defined as the percentage frequency that a bank shows up as lead bank in the syndicated loan market. We find that communication, hardware, and services show a strong positive correlation with changes in lead bank frequency in the syndicated loan market, with communication spending having the largest magnitude. A one standard deviation increase in the lead bank frequency is associated with \$0.27 million more in the bank's annual communication budget. These findings, as we will elaborate in Section 4.2, can be attributed to the distinct responsibilities for handling information assumed by lead and participant banks.

top of the bank size effect. Put differently, one cannot simply use the size of a bank as an empirical proxy for its hierarchical complexity.

## 4 Economics of Banks' IT Investment

Having demonstrated the basic patterns of IT investment in the U.S. banking sector and its interaction with various banking business operations, we now move on to our central question: What are the economics behind banks' IT spending decisions, and how do they relate to—and contribute to—the development of banks' lending technology? We start with a conceptual discussion of lending technologies based on the nature of information handling. By mapping different types of IT investment onto various dimensions of lending technology, this framework helps us understand various empirical patterns shown in Section 3. We then study two credit demand shocks that involve different kinds of information handling, and establish their causal impact on banks' lending technology adoption behaviors.

## 4.1 Lending Technology, Information Handling and IT Spending

We view a bank's lending technology as its ability to deal with borrower-specific information throughout the lending process. Broadly speaking, banks engage in two types/stages of activities in loan making: information production/transmission and information processing. More specifically, information production/transmission—broadly related to soft information in Stein (2002)—refers to the stage in which information on borrowers is gathered and then relayed to those who make decisions later. On the other hand, information processing—broadly related to hard information—is more about the stage in which lenders assemble and examine existing information on borrowers for better decision making.

Communication IT and soft information production/transmission. When facing borrowers with whom lenders have never dealt, or whose information is relatively opaque, for effective information gathering bankers often need to communicate with their borrowers—either through face-to-face meetings, or seeing borrowers' projects for themselves. Once this first-hand information has been gathered, which often can be subjective and thus difficult to convey to others, effective transmission of such information within the organization can also

affect banks' lending efficiency.

One concrete example of how communication technology can help in the two aforementioned dimensions is video conferencing, which has become an important means for loan officers to interact with customers and colleagues during the past decade. In the past, banks opened new checking accounts and originated loans only through brick-and-mortar branches and in-person visits; now, they also use video conferencing, as it makes the direct—yet virtual—contact between loan officers and borrowers more efficient.<sup>28</sup> Moreover, video conferencing within an institution has also been welcomed by the banking sector for its advantages in facilitating effective internal communication and collaboration among employees.<sup>29</sup>

Software IT and hard information processing. Once information has been produced (by the lender itself) or is readily accessible (via a third party), the next concern for the lender is how to use this information. In the context of credit allocation, banks need to properly evaluate the borrowers' creditworthiness to determine loan amounts and rates. For borrowers whom bankers already know from previous interactions or with transparent information, lending decisions simply boil down to efficient processing of the existing information.

Accurate evaluations of borrowers' credit risk often require complicated modeling and simulations, which are impossible without modern software tools. Nowadays banks have actively adopted new software-based technologies to store, organize, and analyze large chunks of loan applicants' data.<sup>30</sup> One popular form of software technology product is credit scoring software for banks making refinancing decisions,<sup>31</sup> which primarily involve the processing and assessment of existing information that lenders already possesses through past interactions. In fact, the recent penetration of fintech companies, which specialize in utilizing software

<sup>&</sup>lt;sup>28</sup>See "Liveoak" for a real world example of a communication tool designed for banking services.

<sup>&</sup>lt;sup>29</sup>See this article from Bankingdive for a detailed description of how video conferencing helps within-bank communication.

<sup>&</sup>lt;sup>30</sup>For example, "nCino" is operating system software that allows financial institutions to replace manual collection of loan/account applications with automated and AI-based solutions. "Finaxtra" and "Turnkey" are both comprehensive loan origination systems that offer solutions for the whole lending process.

<sup>&</sup>lt;sup>31</sup>Some concrete examples of credit scoring software include SAS Credit Scoring, GinieMachine, and RND-Point. To use such software, banks usually just need to import borrowers' demographic and historical data, based on which the software calculates credit scores and conducts statistical tests using AI and machine learning methodologies, saving banks from tedious manual work and expediting the processing.

and algorithm-driven lending approaches, has been particularly pronounced in the mortgage refinancing market (Buchak et al., 2018b; Fuster et al., 2019).

In the next two sections we will explore in detail the lending technology adoption along two dimensions—those targeting the production and transmission of soft information (Section 4.2), and those targeting hard information processing (Section 4.3). In short, communication devices facilitate the gathering and dissemination of soft information, whereas software is for efficiently utilizing "hard" information. From this point on, we focus on two particular categories—communication and software—when examining banks' IT investment behavior.<sup>32</sup>

## 4.2 Bank IT Spending and Soft Information

#### 4.2.1 Soft Information Production/Transmission in Bank Lending

Small business lending. Lending to small business borrowers is one concrete example in which the efficient production and transmission of soft information is essential. Sahar and Anis (2016) document that in the context of lending to small- and medium-size enterprises, direct contact with borrowers and frequent visits to their work sites allow loan officers to collect and produce soft information. Agarwal et al. (2011) highlight that soft information, such as what the borrower plans to do with the loan proceeds, is always the product of multiple rounds of lender-borrower interactions.

That small business lending involves intensive soft information production and transmission is consistent with Section 3.2.1, where we show that banks which specialize in small business lending (as measured by small business loans over total loans) incur more spending on communication IT. As smaller banks generally extend more loans to small businesses (Berger and Udell, 2006; Chen et al., 2017), this helps explain the robust pattern that smaller

<sup>&</sup>lt;sup>32</sup>We will shortly show in Section 4.2 and 4.3 that these two categories of banks' IT spending have a more direct link to banks' dealing with different types of borrower-specific information, a fact already hinted at by the empirical patterns of bank IT spending documented in Section 3.2.

banks have higher fractions of communication IT spending shown in Table A2. Indeed, Table 5 shows that small banks' communication spending is significantly more associated with their small business loans compared with large banks.<sup>33</sup>

Hierarchical complexity. Recall that in Section 3.2.2 we find banks with a more complex hierarchical structure tend to have higher communication IT spending. This is in line with Stein (2002), who argues that a low hierarchical complexity facilitates the within-organization transmission of soft information, making it easier to issue loans requiring soft information (e.g., small business loans). Digging one step further, Table 5 Panel B shows that, given the same percentage increase in small business loans, banks with a more complex hierarchical structure respond with a greater increase in their communication spending. This is consistent with "hierarchical frictions" in soft information transmission: When banks face a need (or choose) to make more small business loans, which implies a demand for improving their soft information handling capability, those with a more complex internal hierarchical structure have to spend more on communication IT so as to overcome such frictions.<sup>34</sup>

Finally, as a placebo test, one should expect no systematic impact of banks' hierarchical complexity on the correlation between their software spending and mortgage refinancing activities, which is indeed confirmed in Table 5 Panel B. Overall, our empirical findings on banks' hierarchical complexity corroborate previous works studying banking organization structure and information production (Degryse et al., 2008; Levine et al., 2020; Skrastins and Vig, 2018), and more research needs to be done on this topic.

<sup>&</sup>lt;sup>33</sup>The relatively lower communication spending by large banks is also consistent with recent empirical findings that large banks, who have deeper pockets than small banks, more frequently invest in or acquire fintech startups (Hornuf et al., 2021; Cornelli et al., 2022). As fintech businesses specialize in transforming the soft information embedded in the alternative data of consumers into credit scores (a form of hard information), large banks' reliance on communication technology in small business lending is lower.

<sup>&</sup>lt;sup>34</sup>Our finding echoes previous work on credit decision making. For instance, Paravisini and Schoar (2016) document that business loan decisions are often made by committee; when decisions cannot be made after committee discussions, the committee will refer to managers in an upper layer, say regional managers. The greater the hierarchical complexity, the higher the "transaction cost" involved for loan decisions.

Lead lender in syndicated loans. The syndicated loan market also provides a special environment to explore the relationship between communication technology and soft information production/transmission. In syndicated lending, the nature of interactions between lenders and borrowers depends crucially on whether the lender is a lead bank or a participant bank (Sufi, 2007). Compared to participant banks, the lead bank is mandated by borrowers to organize other lending participants, conduct compliance reports, and negotiate loan terms. After the loan is issued, it also has the responsibility to conduct monitoring, distribute repayments, and provide overall reporting among all lenders within the deal. In this regard, performing the job of lead bank involves significantly heavier effort in information generation and sharing as well as coordinating negotiations, during which effective communication plays a central role. These differences between lead and participant banks are empirically verified in Section 3.2.1: There is a strong correlation between the frequency of a bank serving as a lead arranger in syndicated loans and its communication IT spending (Table 4 row 4).

#### 4.2.2 Banks IT Spending and Demand Shock on Small Business Loans

We now present the first piece of causal evidence on banks' adaptation of their lending technology by studying their IT investment responses when hit by a positive demand shock in small business loans. As small business lending is associated with intensive soft information production/transmission, we predict that banks will increase their spending on communication technology (soft information), but not on software (hard information).

Our identification strategy relies on a policy shock that affected small businesses' credit demand, which hit the U.S. banking sector heterogeneously across different regions. The "Small Business Health Care Tax Credit" was initially enacted in 2010 as part of the "Affordable Care Act." The program, whose details are available here, offers a tax credit to small business employers who pay health insurance premium on behalf of employees. From 2010-

<sup>&</sup>lt;sup>35</sup>Due to the vast reporting and coordination efforts, lead banks often charge an initiation fee, which can be as high as 10% (Ivashina, 2005).

2013 (the first phase), the tax credit was up to 35% for qualified small businesses (QSBs); and in 2014, the tax credit increased from the 35% to 50% for QSBs (the second phase). To qualify, the employer needed to i) have 25 or fewer employees; ii) pay average wages less than \$50,000 a year per full-time equivalent; iii) pay at least 50% of its full-time employees' premium costs; and finally, iv) have provided a health plan to employees that is qualified under SHOP requirements.

In addition to raising the tax credit from 35% to 50%, in 2014 the government also launched the Small Business Health Options Program (SHOP) Marketplace to offer small business owners a transparent and convenient platform/exchange to compare and shop for insurance packages. Qualified employers were required to purchase insurance packages via the Marketplace, which directly lists health plan choices certified for the tax credit in which the employers could enroll their employees. The Marketplace was initially planned to be launched at the end of 2013; however, there was a delay in the launch of the marketplace till November of 2014 so that the tax credit could be applied to coverage starting from 2015. <sup>36</sup>

We utilize the tax credit hike in 2014 (i.e., the second phase) to identify the impact of soft information demand on banks' technology adoption. Because 2010 is right after the implementation of "the Recovery Act" in 2009, during which numerous other stimulus policies were launched to aid in post-recession economic recovery, this proximity in timing may contaminate the identification. Perhaps more importantly, several surveys revealed that the first phase of the tax credit was not well implemented; some small businesses think the tax credit in the first phase is not high enough, or were not even aware of the policy after its implementation.<sup>37</sup> On the other hand, after the of tax credit hike and the launch of the SHOP Marketplace in 2014, there is a significant decrease in the number of uninsured small

<sup>&</sup>lt;sup>36</sup>See the IRS's FAQ regarding the tax claim rules, and see this report for the delay of the Marketplace.

<sup>&</sup>lt;sup>37</sup>This summary explains the low participation of small businesses in the first few years after 2010: "SHOP programs were operational nationwide, but many features were not initially available, and enrollment had been lower than anticipated. Many small businesses did not enroll because they were apprehensive about joining an unestablished program." Relatedly, the 2012 GAO Report summarized that "the small employers do not likely view the credit as a big enough incentive to begin offering health insurance and to make a credit claim." Regarding small businesses' awareness of the policy, this survey posted in 2011 summarized small businesses' unfamiliarity with the policy.

## business employees.<sup>38</sup>

There are many channels through which this program could boost credit demand from small businesses. First of all, the policy is economically significant: the increased tax credit on average can induce a 14% of savings in terms of total net profit.<sup>39</sup> Thanks to the increased program subsidy, some small business owners who previously could not afford employee health coverage were now likely to provide it; and some may even have chosen to expand their businesses by hiring more employees given the lower effective labor cost.<sup>40</sup> More importantly, the nature of the timing for the tax "rebate" incentivizes small businesses to apply for extra business loans, as all business owners would need to borrow in advance to cover employee health packages and then repay the loan once they have claimed the credit the following year. In turn, banks would be handling additional soft information—such as the employee hiring, health plans, etc.—to screen for genuinely credit-worthy borrowers.

The key to our identification is that the fraction of total establishments that are qualified for the tax policy right before the program launch date varies substantially across different counties. Since the qualified small business share is a key determinant for credit demand from local small businesses, such variation thus helps us identify the impact of the small business credit demand shock on local banks' behavior. As the policy only explicitly targets small

businesses, its impact on other types of local credit demand would be indirect or limited.

38This report finds that uninsured small business employees decreased by around 30% during 2014-2016 compared to 2013.

<sup>&</sup>lt;sup>39</sup>We provide further evidence for the positive impact of the tax credit hike on QSBs by studying their growth in number of both establishments and employees after the tax policy. The details of these analyses on the mechanism and economic magnitude of the policy impact are provided in Appendix Table A7, Table A8, and Table A9.

<sup>&</sup>lt;sup>40</sup>Similar expansion of factor input is also documented in Agrawal et al. (2020), who show that following an R&D tax credit to small businesses, which resembles the health insurance tax credit in our paper, firms responded by increasing their R&D spending significantly. Further, Gao et al. (2023) show that following insurance premium increase, firms reduce employment. More broadly, for reactions from small businesses after the implementation of corporate tax cuts or the launch of subsidies, see Cerqua and Pellegrini (2014), Rotemberg (2019), and Ivanov et al. (2021).

Empirical design: 2SLS regression We run the following 2SLS regression:

$$\Delta \ln(\text{CRA})_{i,c,post} = \tilde{\alpha}_i + \mu_1 \left( \frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}} \right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

$$\Delta \ln \text{IT}_{i,c,post} = \alpha_i + \beta \Delta \widehat{\ln(\text{CRA})}_{i,c,post} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}.$$
(2)

In the first-stage regression, the outcome variable  $\Delta \ln(\text{CRA}_{i,c,post})$  is the change in the logarithm of bank i's small business loans in county c in the three-year time window, before and after the policy change in 2014. The instrumental variable  $\frac{\# \text{Qualified small business est}}{\text{Total } \# \text{ of establishments}} \frac{\text{est}}{c,pre}$  is the proportion of total business establishments that have fewer than or equal to 20 employees, averaged between 2011 and 2013 before the shock. In the second stage, we regress  $\Delta \ln \text{IT}_{i,c,post}$ , which is the change in logarithm of a specific type of IT spending of bank i in county c during 2014-2017 compared to the period of 2011-2013, on the fitted value from the first stage.

The instrument in (2), i.e., the QSB share before the policy shock, is a slow-moving object that reflects the status of the local economy. Our identification assumption is that, conditional on the control variables, the QSB share affects the cross-county growth rate in banks' IT spending around the policy shock only through affecting the small business loans extended in the local economy. Table A8 shows that the growth in small business around the policy year was mostly concentrated in those businesses that were qualified for

 $<sup>^{41}</sup>$ Recall that only employers with 25 or fewer employees are qualified for this program. However, the "County Business Pattern" database provides categorization of small businesses sizes (number of employees) based on the following cut-offs:  $\leq 5$ , 5-9, 10-19, 20-49, 50-99, 100-249, 250-499, 500-1000, and  $\geq$  1000. Due to this data limitation, we chose the closest cut-off, which is "fewer than or equal to 20."

 $<sup>^{42}</sup>$ We make two points. First, the QSB share, together with the growth rate of bank IT spending and local small business loans, are all independent of scale; this helps alleviate the concern that the heterogeneity in policy exposure might be correlated with the size of the local economy. Second, we use  $\ln(\text{Spending})$  and  $\ln(\text{Loan})$  in all of our regression analyses, by removing all observations with zero budget or zero lending. Since we aggregate branch-level observations to bank-county level, the occurrence of zero budget/lending is low—the total amount of observations with zero budget/lending is only 0.6% in our sample. One could use  $\ln(1+x)$  instead of  $\ln(x)$  to include observations with zeros; but because the aggregated IT spending (quoted in USD) and CRA loans (quoted in 1K USD) range from a couple of thousands to millions (software spending and CRA lending have medians of 25,000 and 1,400, respectively), which are much larger than 1,  $\ln(1+x)$  and  $\ln(x)$  are close to each other. Indeed, these two specifications yield quantitatively similar results (in the second stage, we get 0.67 as opposed to 0.76.)

the tax policy, which corroborates our first stage results that counties with a higher "QSB share" experienced faster growth in small business credit around 2014. Finally, the parallel trend assumption requires that heterogeneity in the qualified small business share explains divergent paths in local banks' IT spending only after the policy, which we empirically verify shortly.

We have included a rich set of pre-shock control variables in regression (2). Bank fixed effects absorb any unobserved heterogeneity that may also induce banks in areas with more qualified small businesses to be on a higher IT spending growth path. Revenue per employee at the bank-county level proxies for investment opportunity of a bank in the local economy. We also add a set of county-level economic characteristics, which include county size (proxied by the logarithm of total number of establishments) and local economic situation (proxied by population growth rate, changes in unemployment, labor force participation ratio, the share of non-tradable sector business establishments, and real GDP per capita).

Besides adding controls and fixed effects, we also perform several placebo tests along various dimensions. In Appendix Table A11, we hypothetically postulate the tax policy event to take place in year 2018 and then examine the effect "QSB share" of the dynamics of small business credit around these pseudo event years. In Appendix Table A12, we test whether the variation in the "QSB share" also drives differences in the growth rate in other types of credit (e.g., mortgage origination or refinancing) around the tax policy time. In either test, we find little effect from the variation in "QSB share," in contrast to its significant impact on the small business credit growth around the policy year of 2014.

Estimation results We report the estimation results of (2) in the first three columns of Table 6. Standard errors are clustered at the county level. Column (1) shows the regression estimates in the first-stage regression with a strong first stage result: the F-statistic of 13.71 is above the conventional threshold for weak instruments (Stock and Yogo, 2005).

We find a positive and statistically significant response in banks' communication investment across counties in the second stage. In particular, banks who were facing a one standard deviation higher growth in their small business loan making—due to a higher policy exposure captured by QSB share—experienced a 0.67 std higher growth in their communication spending; in dollar terms, this translates into an increase of \$40,298 in communication IT spending. On the other hand, one standard deviation higher growth in small business loans lead to 0.057 standard deviation slower growth in software spending and is statistically insignificant, suggesting that banks did not respond in increasing their software spending (which is more pertinent to dealing with ready-to-use hard information). Note, by including bank fixed effects, our results come from "within-bank but cross-county" variations. Overall, this asymmetric impact on banks' IT adoption behavior is consistent with our hypothesis that small business lending relies more on soft information handling rather than on processing hard information.

Bank responses and "young firm share" With the premise that young small businesses often lack credit records and thus need extra interaction for loan officers to gather relevant soft information, we test whether the tax policy leads to a larger impact on banks' IT spending response in counties with more young small business borrowers. Specifically, we expand the 2SLS regression in Eq. (2) by introducing an interaction term  $\Delta \ln(\text{CRA})_{i,c,post} \times \text{High young}$ , where the dummy "High young" takes value 1 for counties whose proportion of small businesses younger than 1 year was above median among all counties in 2013. Table 7 shows that communication spending for banks in the "High young share" counties is the main driver of the overall positive causal impact, while the response of banks in "Low young share" counties is statistically insignificant.

<sup>&</sup>lt;sup>43</sup>For detailed calculations, see Appendix Table A6. To put this number into perspective, according to our estimation banks can earn around \$167K extra revenue (given the \$8.366 million extra increase in CRA loans from the first stage), which makes the increased communication IT spending seemingly small. This is expected because soft information relies not only on IT products, but also loan officers who gather and transmit information using the communication IT. Accompanying the communication IT spending, banks will also need to hire extra loan officers (or compensate more hours) when they increase their labor input to deal with the increased small business lending. Since our data does not contain compensation to employees, our calculation only provides a lower bound of the estimation of banks' expenses in response to small business lending growth.

**Dynamic treatment effects** We now study the dynamics of bank IT spending responses to the policy shock; this also helps us evaluate the validity of the IV by examining the pretrend patterns of banks' IT spending. We run the following regression with observations of bank i at county c in year t:

$$\ln IT_{i,c,t} = \alpha_{i,t} + \alpha_{i,c} + \sum_{s \in [-3,3], s \neq -1} \beta_s \times \mathbb{1}_{\{\text{t-2014} = s\}} \times \text{High QSB exposure}_{pre} + \Pi_t \times \mathbf{X}_{i,c,t} + \epsilon_{i,c,t}$$

where  $\alpha_{i,t}$  and  $\alpha_{i,c}$  are bank-year and bank-county fixed effects, and "High QSB exposure" is an indicator variable which equals 1 (0) if the average  $\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}} \frac{\text{est}}{c,pre}$  between 2011 and 2013 sits in the top (bottom) tercile. Note here we allow the coefficients on control variables to be time-varying.

Figure 2 plots the set of estimated coefficients  $\{\hat{\beta}_s\}$ , which measures the intent-to-treat (ITT) effects of the policy change on  $\ln IT_{i,c,t}$  through heterogeneous exposure as captured by QSB share, with the base year as 2013. Prior to the policy shock, the time trends of both types of IT spending display no significant differences for banks located in high-exposure counties versus those in low-exposure counties. Since the tax credit hike in 2014, the communication spending of banks located in high-exposure counties see a continual growth for two consecutive years (left panel), while the software spending of banks (right panel) in high-exposure counties and low-exposure counties demonstrates no difference before 2014 and remain similar after. Finally,  $\{\hat{\beta}_s\}$  for communication IT spending (left panel) starts to decrease around 2016; this might be due to the "capital" nature of IT investment.<sup>44</sup>

Comparison: OLS estimates We report the OLS estimates in Columns (4)–(5) of Table 6. Qualitatively, OLS estimates are similar to those obtained from the 2SLS method; but in terms of magnitude, they are significantly smaller. One explanation for such a downward bias in OLS estimators could be a potential "omitted variable" problem, in which counties

<sup>&</sup>lt;sup>44</sup>That is, having built up their "IT capital" stock after two years of high "flow" spending right after the policy shock, bank branches no longer need to install more IT equipment even if the demand for small business credit remains high in these high-exposure regions. This would then translate into a reduction in the flow of IT spending.

experiencing faster growth in small business loans are those with even faster growth in some unobservable economic variables—say, mortgage refinancing demand—that drive local banks to spend less on communication, leading the OLS estimator to be downward biased.

## 4.3 Bank IT Spending and Hard Information

## 4.3.1 Hard Information Processing in Bank Lending: Mortgage Refinancing

Unlike the lending activities analyzed in Section 4.2 where soft information handling is key, in other situations banks' ability to extend profitable credit is determined by how efficiently they can deal with hard information. As mentioned earlier, mortgage refinancing is the stereotypical type of loan that relies heavily on efficient processing of readily accessible hard information. The discussion in Section 4.1 suggests that banks' software spending should be positively correlated with mortgage refinancing, an empirical fact that we have shown in Table 4 row 5 in Section 3.2.1.

We move one step further and conduct a similar analysis within the mortgage lending business, by splitting it into mortgage origination and mortgage refinancing. For each bank we construct "Refinance/Origination" over the period of 2010-2019, and Table 4 row 7 shows that banks with a greater "Refinance/Origination" spend more on software, while there is no significant effect on communication spending.

The close linkages between banks' software spending and their engagement in mortgage refinancing is also consistent with a recent strand of literature studying fintech lenders' penetration into credit markets. As documented in Fuster et al. (2019), the expansion of fintech lenders—who often serve as the suppliers of new banking software products and typically rely on readily available hard information—is particularly pronounced in the refinancing segment of the mortgage, auto loans, and student loan markets. Later in Section 5, we confirm that software indeed stands out as the major category of IT spending in which commercial banks respond to the entry of fintech lenders.

#### 4.3.2 Bank IT Spending and Demand Shock on Mortgage Refinancing

Paralleling Section 4.2.2, we ask: How would banks respond when hit by credit demand shocks that mostly involve processing hard information, say mortgage refinancing? We expect banks to increase their spending on software (hard information), but not on communication (soft information).

For exogenous sources of cross-sectional variation in mortgage refinance demand, following Di Maggio et al. (2017) and Eichenbaum et al. (2022) we construct an IV for county-level refinancing propensity by utilizing the post-crisis low interest rate period. The nationwide mortgage rate decrease prompted existing homeowners to refinance their mortgages, and an important determinant of homeowners' refinancing propensity was the pre-crisis mortgage characteristics in place before the low-interest episode kicked in.<sup>45</sup>

We consider two ways to construct the instrumental variable. The first follows Eichenbaum et al. (2022) by constructing the IV as the "dollar amount difference." Specifically, for each loan j in county c with unmatured balance in year t between 2011-2016, we calculate the interest savings under the new mortgage rate compared to the old mortgage rate:

```
\begin{split} \Delta \text{Payment}_{j,c,t} &= \text{(Total Interest Payment}|\text{mortgage rate}_j) \\ &- \text{(Total Interest Payment}|\text{new mortgage rate} \text{ }^{\text{FICO, maturity, zip}}_t) \end{split}
```

where the total interest payments are calculated from the amortization schedule with the remaining loan balance as principal. The new mortgage rates are constructed by the bucket of "zip × maturity × FICO" based on new origination in year t, and then matched to each loan j. We then calculate the average  $\Delta$ Payment<sub>c</sub> by taking the average of all savings of unmatured loans j and over years 2011-2016.<sup>46</sup> In words, we calculate the average total

<sup>&</sup>lt;sup>45</sup>Berger et al. (2021) show that effectiveness of monetary policy is crucially dependent upon the previous levels of mortgage rates.

<sup>&</sup>lt;sup>46</sup>We construct the payment savings based on the 2011-2016 sample, because the Federal Funds rate and mortgage rate remained at the low level till 2016 (Figure A2).

remaining mortgage savings under old versus new interest rates at the county level.<sup>47</sup> The variation in local homeowners' refinancing savings thus serves as an exogenous shifter on the mortgage refinance demand faced by local banks,

Though frequently utilized by the previous literature,  $\Delta Payments_c$  is likely correlated with the remaining loan balances of a county, which are in turn correlated with the average loan size or house price level of a county. Therefore  $\Delta Payments_c$  may correlate with local banks' IT spending due to other channels beyond mortgage refinance demand. As an alternative IV, we construct the average mortgage rate gap between the rates at origination and the current rate for the unmatured existing mortgages in a county c:

$$\begin{split} \Delta \text{Mortgage rate}_{c,t} &= \sum_{j} (\text{mortgage rate}_{j,c} - \text{new mortgage rate}_{t}^{(\text{FICO, maturity, zip})}) \\ &\times \frac{\text{Total loan amount}_{j,c}}{\text{Total loan amount during 1999-2010}_{c}}. \end{split}$$

We then take the average of  $\Delta$ Mortgage rate<sub>c,t</sub>—which captures the average mortgage rate savings instead of dollars at year t—for mortgage borrowers across years 2011-2016. In words, for a given county we compute the weighted average of mortgage interest rate gaps, with weights as the loan volume at the initiation.

**Empirical design and estimation results** We aim to identify whether banks' software investment increases given a greater mortgage refinance demand compared with mortgage origination, with the following standard 2SLS specification:

$$\ln(\text{Refinance/Origination})_{i,c} = \tilde{\alpha}_i + \mu_1 \Delta \text{Payments}_c / \Delta \text{Mortgage rate}_c + \mu_2 \mathbf{X}_{i,c} + \tilde{\epsilon}_{i,c},$$

$$\ln(\text{Software})_{i,c} \text{ or } \ln(\text{Communication})_{i,c} = \alpha_i + \beta \ln(\text{Refinance/Origination})_{i,c} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}.$$
(3)

Similar to before, our control variables include banks' revenue per employee and deposit market share of the bank in a county. County level control variables include the unemploy-

<sup>&</sup>lt;sup>47</sup>We remove loans that were defaulted on or prepaid to ensure that the measure captures only refinance propensity from local households with outstanding loans.

ment rate, labor force participation rate, population growth rate, logarithm of number of establishments, and logarithm of small business loans. We include bank fixed effects and cluster standard error at county level.

Table 8 reports our estimation results. In the first stage of column (1), the instrumental variable " $\Delta$ Mortgage rate<sub>c</sub>" predicts mortgage refinancing activities quite well, with a high F-statistics (10.81). For the second-stage, Columns (2) and (3) show that a one standard deviation increase in mortgage refinancing relative to mortgage origination—driven by its local exposure to high refinance interest savings—leads to a 0.315 standard deviation increase in software spending. In dollar terms, this translates to an increased software spending of \$133.26K. This increase is of a similar magnitude, though smaller than, the corresponding \$455.5K revenue increase from mortgage refinancing,  $^{48}$  and it is worth emphasizing that our current data does not provide a comprehensive estimate of costs across other dimensions.

Columns (4)-(6) show the results using  $\Delta$ Payment<sub>c</sub> as the instrumental variable. Note, while the difference between coefficient estimates with these two IVs is statistically insignificant (consistent with the premise that both IVs give us unbiased estimates), we believe  $\Delta$ Mortgage rate<sub>c</sub> satisfies the exclusion restriction condition better. In Appendix Table A15, we tabulate the correlations between  $\Delta$ Payment<sub>c</sub> and  $\Delta$ Mortgage rate<sub>c</sub> with major county-level economic variables. As shown,  $\Delta$ Mortgage rate<sub>c</sub> exhibits statistically and economically insignificant correlations with most of the county characteristic variables, while  $\Delta$ Payment<sub>c</sub> has positive correlations with some of them (though of relatively small magnitude).

Finally, by including bank fixed effects, our result is identified from within-bank-cross-county variations. In addition, communication spending does not demonstrate statistically significant changes in response to the refinancing demand shocks, which supports our premise that mortgage refinancing is a stereotypical lending activity that hinges on efficient processing of readily accessible hard information instead of producing new information.

 $<sup>^{48}</sup>$ This increase of revenue is implied by a one standard deviation increase in  $\ln(\text{Refinance/Origination})$ ; detailed calculation is provided in Appendix Table A6.

Comparison: IV estimates and OLS estimates We conduct the OLS version of the 2SLS regression in Eq. (3) and report the results in the last two columns of Table 8, with quantitatively smaller OLS estimators. Similar to our analysis of small business credit demand in Section 4.2.2, an "omitted variable" issue can explain such downward biases in OLS estimators. Here, counties seeing more mortgage refinances issued by local banks might also have other loan demands that recovered more significantly during the post-crisis period (say, small business loans), which might then tilt local banks IT budget towards other types of IT spending (say, communication as shown in Section 4.2), lowering their spending on software. Our instrumental variable used in the 2SLS method addresses this issue.

## 5 Bank IT Spending and Fintech Entry

In recent years, the emergence and expansion of fintech lenders have drawn heightened public attention to the competition between fintech lenders and traditional banks. Via the angle of examining commercial banks' IT spending, we aim to study a widely debated question: Has the traditional banking sector started reacting to the fast-growing fintech industry? If yes, how?

## 5.1 How Should Banks React to Fintech Entry?

Existing studies suggest that fintech lenders' services involve better use of technology and little human interaction. This tech-intensive feature improves customer experience and likely reduces lending-associated costs (Buchak et al., 2018a; Fuster et al., 2019).

While fintech lenders have been quickly gaining market share in various markets over the past decade, it remains unclear how the incumbent commercial banks should react. For instance, when banks and non-bank lenders offer complementary services, it is possible for banks to strategically shift investment towards areas with fewer activities from fintech lenders. Furthermore, from an information channel, the emergence of fintech lenders who

 $<sup>^{49}</sup>$ Table A4 shows the results of the same OLS specification with bank, year and county fixed effects and bank×year and county fixed effects.

have comparative advantages in information handling in certain markets would render traditional bank lenders more adversely selected in these markets. Both would imply a "falling back" of traditional banks from the markets with fintech entry and a lowered investment in the IT category that fintech lenders have comparative advantages in.

On the other hand, incumbent banks might instead choose to protect their market share and compete against these new fintech entrants, suggesting a potential "catching-up" behavior of the traditional banking sector. Which economic force dominates is an empirical question that we now aim to answer.

## 5.2 Entry of Lending Club and Local Bank IT Investment

To causally identify banks' response in their IT spending towards the increasing presence of fintech lenders, we employ a difference-in-difference strategy that relies on the staggered entrance of Lending Club into different states.

Staggered entry of Lending Club As one of the leading players in the fintech industry, Lending Club launched its platform in 2007. Since 2008, Lending Club has been pursuing regulatory approval to conduct peer-to-peer lending in all 50 states. By October 2008, 41 states moved relatively fast to approve its entry; and between 2010 and 2016, another nine states approved Lending Club's entrance at different times.<sup>50</sup> Table 9 summarizes the timing of Lending Club's staggered entrance into different states.

Following Wang and Overby (2017) and Kim and Stähler (2020), we first drop the 41 states who approved Lending Club's entry in 2008.<sup>51</sup> For Kansas and North Carolina, the

<sup>&</sup>lt;sup>50</sup>As explained by Wang and Overby (2017), Lending Club launched its platform in 2007. In April 2008, Lending Club entered a "quiet" period, in which it suspended peer-to-peer lending until it registered with federal and state regulators as a licensed lender (or loan broker). During this quiet period, Lending Club funded some loans with its own money, and pursued regulatory approval to resume peer-to-peer lending in all 50 states. Six months later, it had received approval in 40 states, plus the District of Columbia by October 2008. For nine states, it received approval at different times between 2010 and 2016. For one state (Iowa), it had not received approval as of February 2021.

<sup>&</sup>lt;sup>51</sup>Given that a majority of states approved Lending Club around the same time period (2008-Q4), a potential concern of endogeneity arises: as these approvals occurred shortly after applications by Lending Club who might have seen a rising opportunity from entering, these approvals might coincide with some unobserved changes in economic conditions happening during the same time.

actual approval was in 2010Q4. Since 2010 is the starting year of our Harte Hanks dataset, 2010 as a pre-treatment period is contaminated for these two states. We hence also exclude these two states, leaving us with a total of seven states for our staggered entrance analysis.

Importantly for our identification purpose, the variation in the approval time since 2010—presumably due to variations in administrative efficiency and potential political issues across states—allows us to get around several major endogeneity concerns regarding the entry of Lending Club. For instance, if Lending Club were to chose to enter the local markets with rising credit demand, then any observed change in local commercial banks' IT investment behavior could not have been convincingly attributed to the entry of their fintech challenger.

For the states in our sample, after its entrance, the personal loan issuance market share of Lending Club across states has a median of 4.85%, with 1.79% (10.29%) being its  $25^{th}$  ( $75^{th}$ ) percentile (Appendix Table A17). These statistics suggest that i) at the state level, Lending Club's presence in the personal loan market features significant variations; and ii) in states where Lending Club actively operates, it makes a nontrivial contribution to the local personal loan market. Both of these two facts are important for our empirical identification, in which the key variation (driven by differences in approval time) operates at the state level.

From the perspective of incumbent banks, Table A18 shows that personal loans represent a significant portion of banks' interest income among all categories of loans, especially for larger banks (banks with more than \$10 billion in assets): around 20% of interest income comes from personal loans. This indicate that banks have a compelling reason to react when fintech lenders emerge in one of their most profit-generating loan segments.

Empirical design and results Our empirical method follows the staggered difference-indifference design as in Wang and Overby (2017). The regression specification is

$$\ln(\text{IT Spending})_{i,c,t} = \alpha_{i,t} + \alpha_c + \beta \times LC_{i,c,t} + \mu_t X_{i,t} + \epsilon_{i,c,t}, \tag{4}$$

where IT Spending  $\in$  {Software, Communication}. We include the bank-year and county fixed effects, denoted by  $\alpha_{i,t}$  and  $\alpha_c$  respectively; and controls  $X_{i,t}$  are in the caption of Table 10. LC<sub>i,c,t</sub> is a dummy variable that is equal to one if Lending Club entered the state where county c is located in year t for bank i;  $\beta$  hence measures the average treatment effect of Lending Club entry on bank technology spending. Estimations are weighted by Lending Club loan volume after entry, and the standard error is clustered at county level.

Columns (1) to (2) in Table 10 Panel A report the results for software and communication spending, respectively. Consistent with the "catching-up" story, column (1) shows that, after Lending Club entered country c, banks on average increase their software IT spending in country c by around 7.6 percentage points, and this estimate is statistically significant. In contrast, communication spending right after Lending Club's entry displays no statistically significant changes compared to pre-entrance.

Figure 3 graphically explores the dynamics of banks' IT spending within the 3-year time window around the fintech entrance year, from the following estimation:

$$\ln \operatorname{IT}_{i,c,t} = \alpha_{i,c} + \mu_t + \sum_{s \in [-3,3], s \neq -1} \beta_s \times \mathbb{1}_{s=t-\text{entrance year}} + \Pi_t \mathbf{X}_{i,c,t} + \epsilon_{i,c,t},$$

where fixed effects and controls are the same as in Eq. (4) and the coefficients on the controls are allowed to be time varying. The estimated  $\{\hat{\beta}_s\}$  and the 95% confidence intervals are plotted. Importantly for our identification, there is no statistically significant pre-trend in either type of IT spending before the fintech entrance, which allows us to plausibly attribute changes in banks' IT spending to the penetration of fintech into the local economy. Consistent with Table 10, a bank's software spending displayed a significantly sharper increase than communication IT spending after the fintech entry.

Recent literature points out the bias in a staggered two-way fixed effects (TWFE) setting, even if the assumption of parallel trends holds. For robustness, we use the interacted TWFE

design as in Callaway and Sant'Anna (2021).<sup>52</sup> As shown in columns (4) to (5) of Table 10, the estimates are similar to, albeit a little larger than, those in columns (1) and (2).

Heterogeneity in responses across bank sizes In Panel B of Table 10, we explore whether banks of different sizes respond differently to fintech entry. Similar to our specification in Table 5, large (small) banks are defined as lenders with asset size above (below) the median size in our sample. We find that large banks increased software spending by 6.2 percentage points more compared to small banks after Lending Club's entry, and the difference is statistically significant. On the other hand, large banks cut their communication spending by 5.8 percentage points compared to small banks following the fintech entry, which is statistically significant. Note, via a different instrument variable, Modi et al. (2022) also document that large banks increase their IT spending when facing competition from fintech lenders, but our data allow us to speak to the underlying mechanism of such response by separating different categories of IT spending, showing that it is mainly driven by "hard" information considerations.

The asymmetric impact on the IT spending reactions by different sized banks is intriguing, and suggests that the specialty (regarding information handling) of the newly entered fintech is more relevant for the market segments served by large banks. This is consistent with Balyuk et al. (2020), who find that fintech lending often substitutes lending made by large banks rather than small banks. Given that small banks engage more in relationship-based small business lending, the entry of Lending Club—who is equipped with superior hard information processing capacity—will not strongly affect these banks' profit making. Furthermore, that large banks cut their communication spending is also consistent with the recent literature studying how fintech entry affects credit market outcomes. For instance, as documented by Balyuk et al. (2020), credit extended by fintech entrants often substi-

 $<sup>^{52}</sup>$ In this method, we run separate regressions in (4) for each group of states that are treated at the same year, with the not-yet-treated as the comparison group, and then aggregate  $\beta$  to form the aggregated average treatment effect of the treated (ATT). For aggregation, we weight the cohort-specific treatment effect by the total volume of loans made through lending club within the three years after Lending Club's entry. Standard errors are based on Bootstrapping with 50 draws.

tutes for loans by out-of-market banks (which are often large ones), as opposed to those by small/in-market banks. As a consequence of large banks' retracted engagement in out-of-market lending, which often rely on the support of communication IT, one should naturally expect them to reduce their communication IT spending.

That banks' IT spending responses are size-dependent is also consistent with the notion that the entry of fintech lenders helps convert soft information to hard information.<sup>53</sup> Linking this "hardening soft-information" effect to our analysis where the focus is placed on bank lenders' decision making, one should expect large banks—rather than small ones who specialize more in relationship-based soft information handling—to reallocate their investment away from communication to software due to a decreased (increased) need of dealing with soft (hard) information.

Finally, recall that in Section 3.2.1 we document a strong correlation between banks' software IT spending and their specialization in the personal loan lending, which is what Lending Club mainly focuses on. Consistent with this, Table 10 Panel C shows that banks with higher personal loan shares respond more significantly to the entry of Lending Club.

Summary and discussion We find that the fintech entry induces banks—especially large ones—to "catch up" and invest to adapt their lending technology. To the best of our knowledge, this is the first piece of direct evidence that the entry of fintech lenders spurs incumbent banks to invest more in their technology to catch up. Furthermore, consistent with existing literature (say, Berg et al., 2021) that highlights the comparative advantage of fintech lenders in processing hard information and making prompt decisions, we show that most "catching up" from traditional banks takes the form of ramping up their *software* IT spending.

We have discussed in Section 5.1 the potential channels through which the entry of fintech lenders affects local commercial banks' IT investment decisions. Our empirical findings support a competition story that, following fintech entry, large banks respond by increasing their IT spending in the relevant categories, presumably to protect their market share.

<sup>&</sup>lt;sup>53</sup>For instance, Beaumont et al. (2019) show that borrowers with better fintech-access are more likely to purchase and pledge hard-information-heavy assets as collateral to obtain new bank credit.

Behind this increased investment in IT could be a "winner's curse" channel that banks need to upgrade their lending technology for fear of being adversely selected by the newly entered fintech competitors, once they have decided to continue operating in the same market segment. However, to fully assess this channel one would need to investigate the composition change of banks' customers induced by the entry of fintech lenders, as well as the dynamics of market share composition. We leave these endeavors to future research.

## 6 Conclusion

Development of information technologies over the past several decades has dramatically revolutionized the way lending is conducted by the banking sector. In this paper, we provide the first comprehensive study of banks' IT spending, which we view as banks' investment to improve their lending technology, especially their ability to deal with soft information and hard information.

The detailed IT spending profiles available in our unique dataset enable us to uncover several novel findings. First, at the aggregate level, we document an overall fast-growing trend in banks' IT spending in the last decade. Second, as a key step in linking banks' IT spending to the development of their lending technology, we show that different types of information technology are closely related to the nature of information embedded in different types of lending activities. More specifically, the production and transmission of "soft" information, which plays a crucial role in conducting small business lending or performing the role of a "lead" bank in syndicated lending, is strongly associated with banks' communication spending. By contrast, "hard" information processing, which is most relevant for conducting mortgage refinancing, is strongly associated with banks' software spending.

We conduct a set of event-based analyses whose answers inform us of how banks adapt their lending technologies in response to economic shocks on their operating environment, including credit demand shocks and the entry of fintech. These causal analyses, to the best of our knowledge, provide the first piece of evidence on the endogenous lending technology adoption in the banking literature.

Our findings open up several important follow-up questions. How does endogenous technology adoption in the banking sector transform the banking/credit market structure? How do technology upgrades in the banking sector affect banks' deposit-taking activities, loan outcomes, properties of the credit cycle, and monetary policy transmission? We leave these questions to future research.

## References

- Agarwal, S., Ambrose, B., Chomsisengphet, S., and Liu, C. (2011). The role of soft information in a dynamic contract setting: Evidence from the home equity credit market. *Journal of Money, Credit and Banking*, 43:633–655.
- Agrawal, A., Rosell, C., and Simcoe, T. (2020). Tax credits and small firm r&d spending. *American Economic Journal: Economic Policy*, 12(2):1–21.
- Ahnert, T., Doerr, S., Pierri, N., and Timmer, Y. (2021). Does IT Help? Information Technology in Banking and Entrepreneurship. *IMF Working Paper*, 21/214.
- Aiello, D., Garmaise, M. J., and Natividad, G. (2020). Competing for deal flow in local mortgage markets. *Working Paper*.
- Amore, M. D., Schneider, C., and Žaldokas, A. (2013). Credit supply and corporate innovation. Journal of Financial Economics, 109(3):835–855.
- Avraham, D., Selvaggi, P., and Vickery, J. (2012). A structural view of u.s. bank holding companies. *Economic Policy Review*, (07):65–81.
- Balyuk, T., Berger, A., and Hackney, J. (2020). What is fueling fintech lending? the role of banking market structure. *Working Paper*.
- Banna, H. and Alam, M. R. (2021). Is digital financial inclusion good for bank stability and sustainable economic development? evidence from emerging asia. *Working Paper*.
- Barkai, S. (2020). Declining labor and capital shares. The Journal of Finance, 75(5):2421–2463.
- Beaumont, P., Tang, H., and Vansteenberghe, E. (2019). The role of fintech in small business lending. *Working Paper*.
- Berg, T., Fuster, A., and Puri, M. (2021). Fintech lending. Working Paper 29421, National Bureau of Economic Research.
- Berger, A. N. (2003). The economic effects of technological progress: Evidence from the banking industry. *Journal of Money, Credit and Banking*, 35(2):141–176.

- Berger, A. N. and Udell, G. F. (2002). Small business credit availability and relationship lending: The importance of bank organisational structure. *The Economic Journal*, 112(477):F32–F53.
- Berger, A. N. and Udell, G. F. (2006). A more complete conceptual framework for sme finance. Journal of Banking and Finance, 30(11):2945 – 2966.
- Berger, D., Milbradt, K., Tourre, F., and Vavra, J. (2021). Mortgage prepayment and path-dependent effects of monetary policy. *American Economic Review*, 111(9):2829–78.
- Bircan, C. and De Haas, R. (2019). The Limits of Lending? Banks and Technology Adoption across Russia. *The Review of Financial Studies*, 33(2):536–609.
- Bloom, N., Garicano, L., Sadun, R., and Reenen, J. V. (2014). The distinct effects of information technology and communication technology on firm organization. *Management Science*, 60(12).
- Bolton, P., Freixas, X., Gambacorta, L., and Mistrulli, P. E. (2016). Relationship and Transaction Lending in a Crisis. *The Review of Financial Studies*, 29(10):2643–2676.
- Buchak, G., Matvos, G., Piskorski, T., and Seru, A. (2018a). Fintech, regulatory arbitrage, and the rise of shadow banks. *Journal of Financial Economics*, 130(3):453–483.
- Buchak, G., Matvos, G., Piskorski, T., and Seru, A. (2018b). Fintech, regulatory arbitrage, and the rise of shadow banks. *Journal of Financial Economics*, 130(3):453–483.
- Calebe de Roure, L. P. and Thakor, A. V. (2019). P2P Lenders versus Banks: Cream Skimming or Bottom Fishing? SAFE Working Paper No. 206.
- Callaway, B. and Sant'Anna, P. H. (2021). Difference-in-differences with multiple time periods. Journal of Econometrics, 225(2):200–230. Themed Issue: Treatment Effect 1.
- Cerqua, A. and Pellegrini, G. (2014). Do subsidies to private capital boost firms' growth? A multiple regression discontinuity design approach. *Journal of Public Economics*, 109:114–126.
- Charoenwong, B., Kowaleski, Z. T., Kwan, A., and Sutherland, A. (2022). RegTech. Working Paper.
- Chava, S., Oettl, A., Subramanian, A., and Subramanian, K. V. (2013). Banking deregulation and innovation. *Journal of Financial Economics*, 109(3):759–774.
- Chen, B. S., Hanson, S. G., and Stein., J. C. (2017). The decline of big-bank lending to small business: Dynamic impacts on local credit and labor markets. *NBER Working Paper Series*, *No.* 23843.
- Cornelli, G., Doerr, S., Franco, L., and Frost, J. (2022). Funding for fintechs: patterns and drivers. Working Paper.
- Degryse, H., Laeven, L., and Ongena, S. (2008). The Impact of Organizational Structure and Lending Technology on Banking Competition. *Review of Finance*, 13(2):225–259.
- Di Maggio, M., Kermani, A., Keys, B. J., Piskorski, T., Ramcharan, R., Seru, A., and Yao, V. (2017). Interest rate pass-through: Mortgage rates, household consumption, and voluntary deleveraging. *American Economic Review*, 107(11):3550–88.

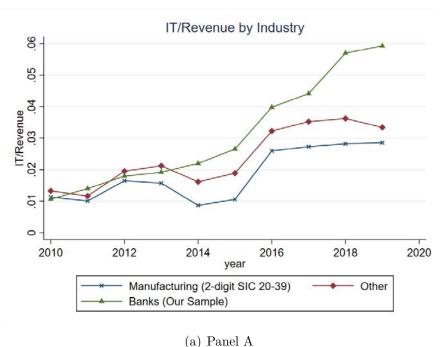
- Di Maggio, M. and Yao, V. (2020). Fintech borrowers: lax Screening or cream-skimming? The Review of Financial Studies.
- Eichenbaum, M., Rebelo, S., and Wong, A. (2022). State-dependent effects of monetary policy: The refinancing channel. *American Economic Review*, 112(3):721–61.
- Erel, I. and Liebersohn, J. (2020). Does finTech substitute for banks? Evidence from the paycheck protection program. Working Paper.
- Feyen, E., Frost, J., Gambacorta, L., Natarajan, H., and Saal, M. (2021). Fintech and the digital transformation of financial services: implications for market structure and public policy. *BIS Working Paper 117*.
- Forman, C., Goldfarb, A., and Greenstein, S. (2012). The internet and local wages: A puzzle. *American Economic Review*, 102(1):556–75.
- Freixas, X. and Rochet, J.-C. (2008). *Microeconomics of Banking, 2nd Edition*. MIT Press Books. The MIT Press.
- Frost, J., Gambacorta, L., Huang, Y., Shin, H. S., and Zbinden, P. (2019). Investment in ict, productivity, and labor demand: The case of argentina. *BIS Working Papers*.
- Fuster, A., Plosser, M., Schnabl, P., and Vickery, J. (2019). The role of technology in mortgage lending. *The Review of Financial Studies*, 32(5).
- Gao, J., Ge, S., Schmidt, L., and Tello-Trillo, C. (2023). How do health insurance costs affect firm labor composition and technology investment? Working Paper.
- Gopal, M. and Schnabl, P. (2022). The Rise of Finance Companies and FinTech Lenders in Small Business Lending. *The Review of Financial Studies*, 35(11):4859–4901.
- Hauswald, R. and Marquez, R. (2003). Information technology and financial services competition. *The Review of Financial Studies*, 16(3):921–948.
- Hauswald, R. and Marquez, R. (2006). Competition and Strategic Information Acquisition in Credit Markets. *The Review of Financial Studies*, 19(3):967–1000.
- He, Z., Huang, J., and Parlatore, C. (2023a). Multi-dimensional information with specialized lenders. Available at SSRN 4557168.
- He, Z., Huang, J., and Zhou, J. (2023b). Open banking: Credit market competition when borrowers own the data. *Journal of Financial Economics*, 147(2):449–74.
- He, Z. and Song, Z. (2022). Agency mbs as safe assets. Technical report, National Bureau of Economic Research.
- Hitt, L., Frei, F., and Harker, P. (1999). How financial firms decide on technology. *Brookings Wharton Papers on Financial Services*.
- Hornuf, L., Klus, M. F., Lohwasser, T. S., and Schwienbacher, A. (2018). How do banks interact with fintechs? forms of alliances and their impact on bank value. *CESifo Working Paper*.

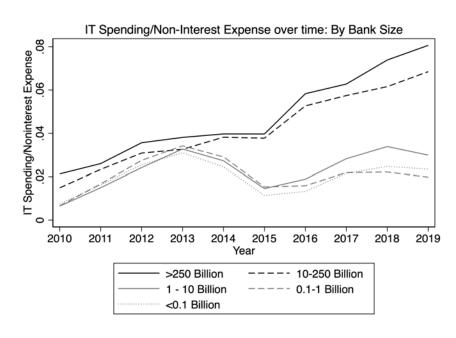
- Hornuf, L., Klus, M. F., Lohwasser, T. S., and Schwienbacher, A. (2021). How do banks interact with fintech startups? *Small Business Economics*, 57(3):1505–1526.
- Huang, J. (2022). Fintech expansion. Available at SSRN 3957688.
- Hughes, J., Jagtiani, J., and Moon, C.-G. (2019). Consumer lending efficiency:commercial banks versus a fintech lender. FRB of Philadelphia Working Paper No. 19-22.
- Huvaj, M. N. and Johnson, W. C. (2019). Organizational complexity and innovation portfolio decisions: Evidence from a quasi-natural experiment. *Journal of Business Research*, 98:153–165.
- Ivanov, I., Pettit, M. L., and Whited, T. (2021). Taxes depress corporate borrowing: Evidence from private firms. *Working Paper*.
- Ivashina, V. (2005). Structure and pricing of syndicated loans. Working Paper.
- Jagtiani, J. and Lemieux, C. (2017). Fintech lending: Financial inclusion, risk pricing, and alternative information. FRB of Philadelphia Working Paper No. 17-17.
- Kim, J.-H. and Stähler, F. (2020). The impact of peer-to-peer lending on small business loans. Working Paper.
- Kovner, A., Vickery, J., and Zhou, L. (2014). Do big banks have lower operating costs? *Economic Policy Review*, 20(2).
- Kwon, S., Ma, Y., and Zimmermann, K. (2022). 100 years of rising corporate concentration. Working Paper.
- Lerner, J., Seru, A., Short, N., and Sun, Y. (2021). Financial innovation in the 21st century: Evidence from u.s. patents. Working Paper 28980, National Bureau of Economic Research.
- Levine, R., Lin, C., Peng, Q., and Xie, W. (2020). Communication within Banking Organizations and Small Business Lending. *The Review of Financial Studies*, 33(12):5750–5783.
- Liberti, J. M. and Mian, A. R. (2009). Estimating the effect of hierarchies on information use. *The Review of Financial Studies*, 22(10):4057–4090.
- Lorente, C., Jose, J., and Schmukler, S. L. (2018). The fintech revolution: A threat to global banking? Research and Policy Briefs 125038, The World Bank.
- McGrattan, E. R. (2017). Intangible capital and measured productivity. Working Paper 23233, National Bureau of Economic Research.
- Mian, A. and Sufi, A. (2014). What explains the 2007–2009 drop in employment? *Econometrica*, 82(6):2197–2223.
- Modi, K., Pierri, N., Timmer, Y., and Pería, M. S. M. (2022). The anatomy of banks' it investments: Drivers and implications. *IMF working paper*.
- Paravisini, D. and Schoar, A. (2016). The incentive effect of scores: Randomized evidence from credit committees,. *NBER Working Paper*, 19303.

- Petersen, M. A. and Rajan, R. G. (2002). Does distance still matter? the information revolution in small business lending. *The Journal of Finance*, 57(6):2533–2570.
- Pierri, N. and Timmer, Y. (2020). Tech in fin before fintech: Blessing or curse for financial stability? *IMF Working Paper, No. 20/14.*
- Pierri, N. and Timmer, Y. (2022). The importance of technology in banking during a crisis. *Journal of Monetary Economics*.
- Ridder, M. D. (2019). Market Power and Innovation in the Intangible Economy. Working Paper.
- Rotemberg, M. (2019). Equilibrium effects of firm subsidies. *American Economic Review*, 109(10):3475–3513.
- Sahar, L. and Anis, J. (2016). Loan officers and soft information production. Cogent Business & Management, 3(1):1199521.
- Skrastins, J. and Vig, V. (2018). How Organizational Hierarchy Affects Information Production. *The Review of Financial Studies*, 32(2):564–604.
- Stein, J. C. (2002). Information production and capital allocation: Decentralized versus hierarchical firms. *The Journal of Finance*, 57(5):1891–1921.
- Stock, J. and Yogo, M. (2005). Testing for Weak Instruments in Linear IV Regression, pages 80–108. Cambridge University Press, New York.
- Stulz, R. M. (2019). FinTech, BigTech, and the future of banks. *Journal of Applied Corporate Finance*, 31(4):86–97.
- Sufi, A. (2007). Information asymmetry and financing arrangements: Evidence from syndicated loans. *The Journal of Finance*, 62(2):629–668.
- Tang, H. (2019). Peer-to-Peer lenders versus banks: Substitutes or complements? The Review of Financial Studies, 32(5):1900–1938.
- Tuzel, S. and Zhang, M. B. (2021). Economic stimulus at the expense of routine-task jobs. *The Journal of Finance*, 76(6):3347–3399.
- Vives, X. and Ye, Z. (2021). Information technology and bank competition. Working Paper.
- Wang, H. and Overby, E. (2017). How does online lending influence bankruptcy filings? evidence from a natural experiment. *Management Science*, 2022:15937.

## Figure 1. IT Spending: Time Trend

Panel A shows the evolution of IT spending as a share of total revenue of banks in our sample, the manufacturing sector, and all other industries constructed using IT budget and revenue from Harte Hanks. "Manufacturing" sector is defined as establishments with 2-digit SIC code 20-39. "Other" sector is defined as all industries other than "Depository institutions" (2-digit SIC code=60). "Banks (our sample)" refers to banks in our sample. These ratios are calculated by aggregating total IT spending and then scaling it against the total revenues sourced from Harte Hanks. In Panel B, the vertical axis is banks' total IT spending scaled by non-interest expenses. The asset size groups are categorized based on a bank's average asset size during 2010 and 2019. Non-interest expenses are calculated using banks' balance sheet item "RIAD4093" in the Call Report.



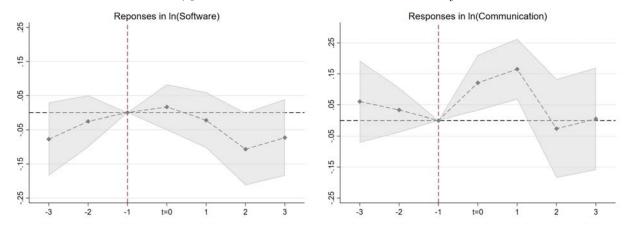


## Figure 2. Bank IT Spending Around Small Business Tax Credit Policy

This figure reports the event studies of IT spending around the small business tax credit event. The specification is

$$\ln IT_{i,c,t} = \alpha_{i,t} + \alpha_{i,c} + \sum_{s \in [-3,3], s \neq -1} \beta_s \times \mathbb{1}_{\{\text{t-2014} = s\}} \times \text{High QSB exposure}_{pre} + \Pi_t \times \mathbf{X}_{i,c,t} + \epsilon_{i,c,t}$$

where for bank i at county c in year t,  $\alpha_{i,t}$  are the bank-year fixed effects,  $\alpha_{i,c}$  are the bank-county fixed effects.  $1_{\{t-2014=s\}}$  is a dummy variable that is equal to one if the distance between year t and the event year (2014) is s. "High exposure pre" is equal to one if the average  $\left(\frac{\# \text{ Qualified small business est}}{\text{Total }\# \text{ of establishments}}\right)_{c,pre}$  is within the top tercile between 2011-2013; "High exposure" is equal to zero if the average  $\left(\frac{\# \text{ Qualified small business est}}{\text{Total }\# \text{ of establishments}}\right)_{c,pre}$  in the bottom tercile between 2011-2013. Bank control variables include banks' revenue per employee of the bank in a county. County control variables include unemployment rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of small businesses in non-tradable sector, and GDP per capita. Shaded regions are the 95% confidence interval of the estimated  $\beta_s$ . Standard errors are clustered at the county level.

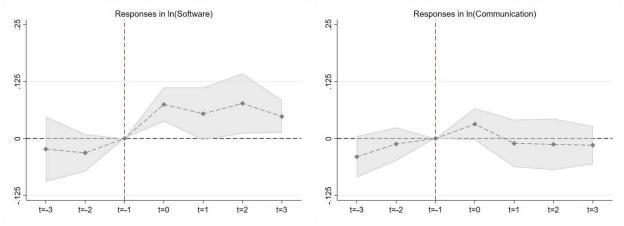


## Figure 3. Bank IT Spending Around Fintech Entrance

This figure reports the event studies of IT spending around the entrance of Lending Club. The specification is

$$\ln IT_{i,c,t} = \alpha_{i,c} + \mu_t + \sum_{s \in [-3,3], s \neq -1} \beta_s \times \mathbb{1}_{t-\text{entrance year} = s} + \Pi_t \mathbf{X}_{i,c,t} + \epsilon_{i,c,t}$$

where for bank i at county c in year t,  $\alpha_{i,c}$  are the bank-county fixed effects,  $\mu_t$  are the year fixed effects.  $\mathbbm{1}_{t-\text{entrance year}=s}$  is a dummy variable that is equal to one if the number of years between the observation year t and the Fintech entrance year into the state where county c is located is s. For the left panel, the left-hand-side variable is logarithmic spending on software IT. For the right panel, the left-hand-side variable is logarithmic spending on communication IT. Bank control variables include banks' revenue per employee of the bank in a county. County control variables include unemployment rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of small businesses in non-tradable sector, and GDP per capita. Shaded regions are the 95% confidence interval of the estimated  $\beta_s$ . Standard errors are clustered at the county level.



## Table 1. Sample Coverage

This table demonstrates the sample coverage of banks across five categories of banks' size groups. The Call Report bank population is constructed by applying the commercial bank restriction ("Charter Type" being 200) following FFIEC definition. The first two columns show the number of banks and the average asset sizes of banks in our sample, across five size groups. Column 3 and column 4 show the total number of banks and average asset sizes of all banks in the Call Report. Column 5 shows the percentage of sample coverage in terms of frequency compared with the population in Call Report, and column 6 shows the percentage of sample coverage in terms of total asset size compared with the population in Call Report.

Coverage of data	Sample		Call Report		Freq %	Asset %
Average Assets 2010-2019 (Billion)	Num banks	Ave Assets	Num banks	Ave Assets		
>\$250 Billion	6	1196.15	6	1196.15	100%	100%
\$10 Billion–\$250 Billion	98	43.82	106	43.69	92.45%	92.72%
\$1 Billion—\$10 Billion	418	2.95	590	2.78	70.85%	85.62%
\$100 Million-\$1 Billion	734	0.42	4161	0.32	17.64%	23.43%
<\$100 Million	194	0.06	2048	0.05	9.47%	11.36%

## Table 2. IT Spending Summary Statistics

This table shows the summary statistics of banks' IT Spending. Total IT Spending is the sum of all types of IT spending in millions of dollars. "IT Spending/Revenue" is total IT Spending scaled by banks' total gross income ("Revenue" is RIAD4000 of Call Report); "IT Spending/Non-interest expense" is total IT spending scaled by non-interest expenses ("Non-interest expenses" is RIAD4093 in Call Report); "IT spending/Net income" is total IT spending scaled by total income minus the gross total expenses ("Net income" is total income minus the sum of RIAD4073 and RIAD4093 in Call Report). The different categories of IT spending are the four categories of IT spending scaled by total IT spending.

	Mean	S.d.	p(25)	Median	p(75)
Total IT Spending (Million)	11.125	160.239	0.024	0.159	0.796
No. of IT Employees	178.756	1828.766	5.000	20.682	56.912
IT Spending/Income	0.020	0.039	0.006	0.012	0.021
IT Spending/Net income	0.068	0.113	0.017	0.037	0.084
IT Spending/Expenses	0.022	0.027	0.008	0.014	0.026
IT Spending/Noninterest expense	0.051	0.036	0.009	0.018	0.035
Communication/Total	0.089	0.108	0.028	0.051	0.110
Software/Total	0.334	0.172	0.219	0.315	0.468
Hardware/Total	0.172	0.111	0.066	0.161	0.235
Services/Total	0.327	0.129	0.243	0.329	0.415
Other/Total	0.066	0.104	0.009	0.022	0.111

## Table 3. C&I Loans and Banks' IT Spending

This table presents the results of the regression of banks' C&I loan on the four major categories of banks' IT spending and a vector of control variables at bank-year level. The sample period is 2010 to 2019.

$$\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,10\text{-}19} = \alpha + \beta \frac{\text{C\&I Loan}}{\text{Total loan}}_{i,10\text{-}19} + \gamma \mathbf{X} + \epsilon_i$$

C&I Loan/Total Loan is commercial and industrial loan of bank i scaled by total loan between 2010-2019, Software/Rev is software spending scaled by total revenue, Communication/Rev is communication spending scaled by total revenue, Hardware/Rev is Hardware spending scaled by total Revenue, Services/Rev. Control variables include net income scaled by total assets, deposits scaled by total assets, revenue per employee, salaries scaled by total assets and equity scaled by total assets. Both the left-hand side and the right-hand side variables are taken using the average values across 2010-2019 for each bank i. Fixed effects include bank size group and banks' headquarters state fixed effects. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	Software/Revenue	Communication/Revenue	Harware/Revenue	Services/Revenue
	(1)	$\frac{}{(2)}$	(3)	(4)
C& I loans/Total loan	0.0380	0.0813***	0.107***	0.0386
	(0.0276)	(0.0270)	(0.0273)	(0.0279)
Net income/Total Assets	-0.126***	-0.155***	-0.183***	-0.0904***
	(0.0316)	(0.0308)	(0.0312)	(0.0318)
Deposits/Assets	-0.0928	-0.284*	-0.242	-0.103
	(0.175)	(0.171)	(0.173)	(0.177)
Revenue/Employee	-0.313***	-0.437***	-0.397***	-0.333***
	(0.0542)	(0.0529)	(0.0536)	(0.0546)
Salaries/Assets	0.0683	-0.147***	-0.0550	0.0104
	(0.0452)	(0.0441)	(0.0447)	(0.0455)
Equity/Assets	0.140**	$0.0925^{*}$	0.0647	0.124**
	(0.0574)	(0.0560)	(0.0567)	(0.0578)
Size group FE	Y	Y	Y	Y
State FE	Y	Y	Y	Y
AdR-squared	0.0925	0.127	0.110	0.0649
N	1442	1442	1442	1442

Standard errors in parentheses

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

## Table 4. Bank Characteristics and Banks' IT Spending

This table presents the results of correlation between banks' IT spending and banks' characteristics. The regression specification is as follows.

$$\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,10\text{-}19} = \alpha + \beta \frac{\text{Type L loan}}{\text{Total loan}}_{i,10\text{-}19} \text{ or (Bank Char)} + \gamma \mathbf{X} + \epsilon_i$$

Panel A shows how the banks' loan specialization correlates with banks' IT spending. Type L loan/Total Loan is the average of a specific type of loan scaled by total loan. Among them, Personal loan/Total Loan is the sum of personal loans and real estate loans to 1-4 family units scaled by total loan; Agriculture/Total loan is the agricultural loan scaled by total loan; CRA/Total loan is the sum of small business loans reported in CRA scaled by total loan; "Other C&I/Total loan" is the total C&I loan minus small business loans reported in CRA, scaled by total loan; "Mortgage refinance" is the total amount of mortgage refinance reported in HMDA scaled by the bank's total loan; "Other personal loans" is the deduction of "Mortgage refinance" from "Personal and mortgage loans." "Refinance/Origination" is the dollar amount of mortgage refinance scaled by dollar amount mortgage origination of a bank. Software/Rev is software spending scaled by total revenue, Communication/Rev is communication spending scaled by total revenue, Hardware/Rev is Hardware spending scaled by total revenue, Services/Rev is services spending scaled by total revenue. Panel B shows how a bank's hierarchical structure correlates with its IT spending. "Hierarchical layer" is the number of types of its locations as defined in Section 2.3. "In(num offices)" is the logarithmic of total number of offices. Control variables include net income scaled by total assets, deposits scaled by total assets, revenue per employee, salaries scaled by total assets and equity scaled by total assets. Fixed effects include bank size group, and banks' headquarter state fixed effects. Panel C shows how a bank's role in the syndicated loan market correlates with its IT spending. "Lead bank is the frequency of a bank's showing up as a lead bank in the syndicated loan market as a share of total number of syndicated loans lent out. All of the loan profile variables are calculated as the average of the loan profile of a bank between 2010 and 2019. \*\*\*\*, \*\*\*, and \* denote si

	Panel A:	Loan Specialization		
	Software/Revenue	Communication/Revenue	Hardware/Revenue	Services/Revenue
	(1)	(2)	(3)	(4)
C& I loan/Total loan	0.0380	0.0813***	0.107***	0.0386
	(0.0276)	(0.0270)	(0.0273)	(0.0279)
CRA/Total loan	-0.190***	$0.124^{***}$	0.0662**	0.0482
	(0.0307)	(0.0303)	(0.0309)	(0.0314)
Other C&I loan/Total loan	$0.0645^{**}$	0.0653**	$0.0995^{***}$	0.0330
	(0.0275)	(0.0269)	(0.0273)	(0.0278)
Personal loan	0.0617**	$0.0507^{*}$	0.0162	-0.000905
	(0.0294)	(0.0287)	(0.0292)	(0.0297)
Refinance/Total loan	0.0763***	0.0369	0.0466	-0.000550
	(0.0311)	(0.0305)	(0.0309)	(0.0314)
Other personal loan/Total loan	$0.0530^{*}$	$0.0522^{*}$	0.0103	0.00167
C F/	(0.0292)	(0.0285)	(0.0290)	(0.0295)
Refinance/Origination	0.0682**	0.0257	0.0482	0.0442
	(0.0329)	(0.0315)	(0.0322)	(0.0335)
Agriculture loan/Total loan	0.0238	0.0504	0.00356	0.0150
	(0.0337)	(0.0331)	(0.0336)	(0.0343)
P	anel B: Hierarchica	al Complexity and IT Sp	ending	
Hierarchical layer	0.0244	0.0721**	0.0331	0.0371
·	(0.0354)	(0.0342)	(0.0351)	(0.0354)
ln(num of offices)	$0.0519^{'}$	0.0842**	0.0278	0.0487
,	(0.0413)	(0.0394)	(0.0406)	(0.0413)
	Panel C: Banks'	Role in Syndicated Lend	ling	
% Lead bank/Total syndicate	0.0751	0.0932**	0.0475	0.0162
	(0.0468)	(0.0453)	(0.0467)	(0.0473)

## Table 5. Bank Characteristics and Banks' IT Spending: Size- and Hierarchical-Dependence

This table presents the results of the dependence of correlation between banks' IT spending with their lending activities on the size and hierarchical complexity of banks. The regression specification is as follows.

$$\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,10\text{-}19} = \alpha + \beta \times (\text{Bank Char.}) \times \left(\frac{\text{CRA}}{\text{Total loan}}_{i,10\text{-}19} \text{or } \frac{\text{Refinance}}{\text{Total loan}}_{i,10\text{-}19}\right) + \gamma \mathbf{X} + \epsilon_i$$

In Panel A, small (large) banks are defined as the banks with asset size below (above) median asset size in our sample. In Panel B, "High layer" is defined equal to 1 if a bank has 2 or 3 hierarchical layers. the number of types of its locations as defined in Section 2.3. "Size group FE" refers to the fixed effects of the five bank asset groups defined in Section 3.1 or Table 1. Control variables include net income scaled by total assets, deposits scaled by total assets, revenue per employee, salaries scaled by total assets and equity scaled by total assets. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

Panel A: Bank Size and IT Spending Software/Revenue Communication/Revenue							
	Software/ Revenue	Communication/ Revenue					
	(1)	(2)					
Refinance/Total loan	0.0817*						
	(0.0428)						
Small×Refinance/Total loan	0.0315						
	(0.0538)						
CRA/Total loan	, ,	0.0282					
		(0.0413)					
Small×CRA/Total loan		0.162**					
·		(0.0699)					
Small	-0.0310	-0.286***					
	(0.0824)	(0.0894)					
Size group FE	Y	Y					
State group FE	Y	Y					
R-squared	0.103	0.133					
N	1432	1432					

	Software/Revenue	Communication/Revenue
	(1)	(2)
Refinance/Total loan	0.0409	
	(0.0361)	
High layer×Refinance/Total loan	0.0836	
	(0.0516)	
CRA/Total loan		0.301**
		(0.121)
High layer×CRA/Total loan		0.0870*
		(0.0517)
High layer	0.0532	0.0242
	(0.0521)	(0.0506)
Size group FE	Y	Y
State group FE	Y	Y
R-squared	0.0894	0.127
N	1426	1426

## Table 6. Soft Information and Banks' IT Spending

This table presents the results of 2SLS and OLS discussed in Section 4.2.2. The first three columns show the results for the following specification:

$$\Delta \ln(\text{CRA})_{i,c,post} = \tilde{\alpha}_i + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$
$$\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \widehat{\Delta \ln(\text{CRA})}_{i,c,post} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$$

Column (4) and column (5) show the following OLS specification:

$$\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \Delta \ln(\text{CRA})_{i,c,post} + \mu_c + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln (CRA)_{i,c,post}$  is the change in average natural log of small business loans reported in CRA of bank i at county c during the years 2014-2017 compared with 2011-2013. Bank control variables include pre-shock revenue per employee of the bank in a county. County level control variables include the pre-shock unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of nontradable sector small business establishments, and GDP per capita. Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	First stage	ln(Software)	ln(Communication)	$\ln(\text{Software})(\text{OLS})$	ln(Communication)(OLS)
	(1)	(2)	(3)	(4)	(5)
	1.032***				
e,pre	(0.251)				
$\widehat{\Delta \ln(CRA)}$		-0.057	0.670**		
		(0.305)	(0.328)		
$\Delta \ln(CRA)$				0.004	$0.019^*$
				(0.010)	(0.011)
Bank FE	Y	Y	Y	Y	Y
Clustered	Y	Y	Y	Y	Y
Controls	Y	Y	Y	Y	Y
F-stat	13.708				
AdR-squared	0.427	-0.179	-0.522	0.120	0.102
N	19,848	19,848	19,848	19,848	19,848

## Table 7. Dependence of Banks' Shock Response on "Young Firm Share"

This table presents the impact of soft information demand and banks' IT spending for counties with differential share of younger small businesses. The 2SLS regression specifications are as follows:

$$\Delta \ln(\text{CRA})_{i,c,post} = \hat{\alpha}_i + \eta_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \eta_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

$$\Delta \ln(\text{CRA})_{i,c,post} \times \text{High young} = \tilde{\alpha}_i + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

$$\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \Delta \ln(\widehat{\text{CRA}})_{i,c,post} + \beta_1 \times \Delta \ln(\widehat{\text{CRA}})_{i,c,post} \times \text{High young} + \beta_2 \text{High young} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln (CRA)_{i,c,post}$  is the change in average natural log of small business loans reported in CRA of bank i at county c during the years 2014-2017 compared with 2011-2013. Bank control variables include pre-shock revenue per employee. "High young" counties are defined as counties whose proportion of small businesses younger than 1 year old was above median among all counties in 2013. County level control variables include the pre-shock unemployment growth rate, labor force participation rate, population growth rate, log of total number of establishments, share of nontradable sector small business establishments, and GDP per capita. Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	First stage $\Delta \ln(\text{CRA})$	$\begin{array}{c} {\rm First\ stage} \\ \Delta \ln({\rm CRA}) \times {\rm High\ young} \end{array}$	Second stage ln(Software)	$\begin{array}{c} {\rm Second\ stage} \\ {\rm ln}({\rm Communication}) \end{array}$	OLS ln(Software)	OLS ln(Communication)
	(1)	(2)	(3)	(4)	(5)	(6)
$\sqrt{NQSB_{pre}}$	0.024*	0.028**				
	(0.012)	(0.014)				
$\%QSB_{pre} \times \text{High young}$	-0.021**	$0.025^{**}$				
	(0.010)	(0.010)				
$\Delta \ln(\widehat{\mathrm{CRA}})$			-0.429	-0.321		
			(0.617)	(0.685)		
$\Delta \widehat{\ln(CRA)} \times High young$			0.687	$1.534^{*}$		
, , , ,			(0.844)	(0.928)		
High young			0.017	0.025	-0.042*	-0.056**
			(0.070)	(0.075)	(0.023)	(0.022)
$\Delta \ln(\mathrm{CRA})$					0.003	0.017
					(0.014)	(0.016)
$\Delta \ln(\text{CRA}) \times \text{High young}$					0.004	0.019
					(0.021)	(0.021)
Bank FE	Y	Y	Y	Y	Y	Y
Clustered	Y	Y	Y	Y	Y	Y
F-stat	12.350	12.350				
AdR-squared	0.321	0.291	-0.330	-1.289	-0.179	-0.178
N	19,234	19,234	19,234	19,233	19,234	19,233

57

Standard errors in parentheses

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

## Table 8. Hard Information and Banks' IT Spending

This table presents the results of the regressions discussed in Section 4.3.2.

The first six columns show the results for the 2SLS specification below:

ln(Refinance/Origination)<sub>i,c</sub> = 
$$\tilde{\alpha}_i + \mu_1 \times \Delta \text{Mortgage rate}_c(\text{or } \Delta \text{Payments}_c) + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$
  
ln(Type S Spending)<sub>i,c</sub> =  $\alpha_i + \beta \times \widehat{\text{ln}(\text{Refinance/Origination})_{i,c}} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$ 

Column (7) and (8) show the results of the OLS specification below:

$$\ln(\text{Type S Spending})_{i,c} = \alpha_i + \beta \times \ln(\text{Refinance})_{i,c} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $ln(Type S Spending)_{i,c}$  is the average logarithmic of banks' IT spending during 2011 and 2016.  $ln(Refinance/Origination)_{i,c}$  is the average logarithmic of amount of mortgage refinance loan relative to mortgage origination issued by bank i in county c during 2011 and 2016.  $\Delta$ Payments<sub>c</sub> is the hypothetical amount of interest payments that could be saved due to the interest rate decrease, if local households chose to refinance their mortgages during the year of 2011 and 2016. ΔMortgage rate<sub>c</sub> is the average differences of mortgage rate of unmatured loans in 2011-2016 and the prevailing mortgage rates of newly issued mortgages in a county c. Bank control variables include banks' revenue per employee of the bank in a county. County level control variables include unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, logarithmic of total small business loan, share of nontradable sector establishments, and GDP per capita. Fixed effects include bank fixed effects.

Standard errors are clustered at county level. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

		2SLS			2SLS		0]	LS
	First stage	$\ln(\text{Software})$	$\ln(\text{Comm})$	First stage	$\ln(\text{Software})$	$\ln(\text{Comm})$	$\ln(\text{Software})$	$\ln(\mathrm{Commu})$
	(1)	(2)	(3)	(4)	$\frac{}{(5)}$	(6)	(7)	(8)
$\Delta$ Mortgage rate <sub>c</sub>	1.824*** (0.622)							
$\Delta \mathrm{Payment}_c$	,			$0.819^{***}$ $(0.237)$				
ln Refinance/Origination		$0.315^*$ (0.167)	0.239 $(0.150)$		$0.373^*$ $(0.225)$	0.296 $(0.211)$		
$\ln(\text{Refinance/Origination})$		, ,			, ,	,	0.024*** (0.006)	$0.025^{***}$ $(0.006)$
Bank FE	Y	Y	Y	Y	Y	Y	Y	Y
Clustered	Y	Y	Y	Y	Y	Y	Y	Y
F-stat	10.81			13.82				
AdR-squared	0.356	-0.349	0.072	0.423	0.447	0.576	0.449	0.423
N	$14,\!626$	14,626	14,626	14,626	14,626	14,626	14,626	14,626

Table 9. Staggered Entry of Lending Club to 9 States after 2010

State	Approval year
All states, except the states listed below	2008
Kansas	2010 Q4
North Carolina	2010 Q4
Indiana	2012 Q4
Tennessee	2013 Q1
Mississippi	2014 Q2
Nebraska	2015  Q2
North Dakota	2015  Q2
Maine	2015 Q3
Idaho	2016 Q1
Iowa	Not approved as of 2022-Q1

## Table 10. Fintech Entry and Banks' Lending Technology Adoption

This table presents the effect of Lending Club's entrance on local banks' IT spending. The regression equation is as follows

$$\ln(\text{ITSpending})_{i,c,t} = \alpha_{i,c} + \alpha_t + \beta \times \text{LC}_{i,c,t} + \gamma_t \mathbf{X}_{i,t} + \epsilon_{i,c,t},$$

where  $\alpha_{i,c}$  and  $\alpha_t$  are the bank-county and year FE, respectively. Column (1) and (2) of Panel A show the baseline results. Bank control variables include banks' revenue per employee of the bank in a county. County level control variables include unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of nontradable sector establishments, and GDP per capita. Standard errors are based on 50 Bootstrapped samples. Panel B presents the differential responses to Fintech entrance of banks with different sizes. "Large banks" are defined as banks with asset size above median of all the asset sizes in the sample. Panel C presents the differential responses to Fintech entrance of banks with different personal loan share. "High personal loan" banks are defined as banks for which the personal loan as a share of total loan is above median among all banks in the sample. The estimations in Panel column 3 and column 4 of the three panels are based on the interacted TWFE method as in Callaway and Sant'Anna (2021). Standard errors are in the parentheses and are clustered at the county level. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

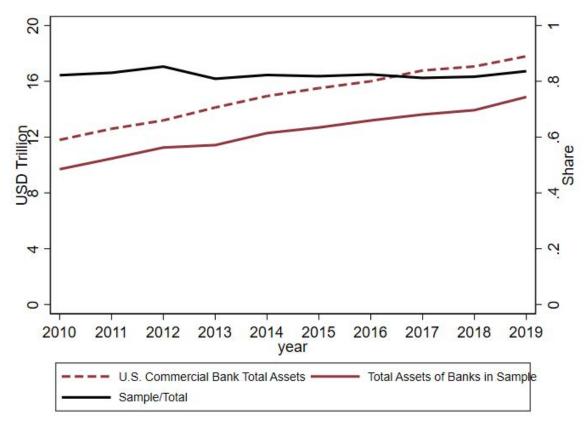
Panel A		Baseline	Callaway an	d Sant'Anna (2021)
	ln(Software)	ln(Communication)	ln(Software)	ln(Communication)
	(1)	(2)	(3)	(4)
After	0.076***	0.001	0.080**	0.007
	(0.023)	(0.018)	(0.042)	(0.040)
Fixed Effects	Bank×Coun	ty, Year, Size group		
AdR-squared	0.808	0.790		
N	13,406	13,406		
Panel B		Baseline	Callaway an	d Sant'Anna (2021)
	ln(Software)	ln(Communication)	ln(Software)	ln(Communication)
	(1)	(2)	(3)	(4)
After	0.051	0.037	0.043	0.090
	(0.032)	(0.030)	(0.048)	(0.052)
$After \times Large$	0.062*	-0.058*	0.097***	-0.167***
	(0.032)	(0.035)	(0.040)	(0.061)
Fixed Effects	Bank×Coun	ty, Year, Size group		
Clustered	Y	Y		
AdR-squared	0.777	0.96		
N	13,406	13,406		
Panel C		Baseline	Callaway and Sant'Anna (20	
	ln(Software)	ln(Communication)	ln(Software)	ln(Communication)
	(1)	(2)	(3)	(4)
After	0.054**	0.007	0.058	0.050
	(0.026)	(0.023)	(0.050)	(0.052)
After×High personal loan	0.055**	-0.002	0.078**	-0.041
	(0.025)	(0.026)	(0.043)	(0.054)
Fixed Effects	Bank×Coun	ty, Year, Size group		
Clustered	Y	Y		
AdR-squared	0.837	0.774		
N	13,406	13,406		

Online Appendix
— "Investing in Lending Technology: IT Spending in Banking"

## A Supplemental Materials for Additional Analysis

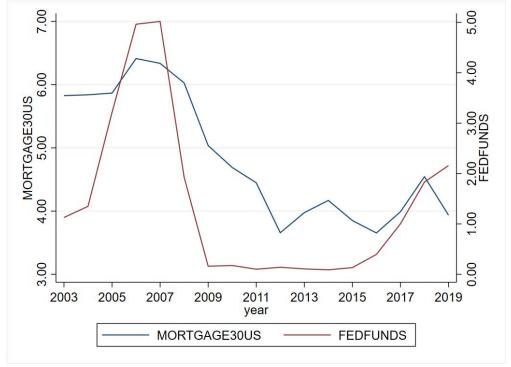
## Figure A1. Total Asset of Banks in Sample

This figure shows the sum of total asset size of all banks in our sample from 2010 to 2019 U.S. The red dashed line is the sum of all commercial banks' asset size in U.S., data source is Board of Governors of the Federal Reserve System (US), Total Assets, All Commercial Banks [TLAACBW027SBOG]. The red solid line is the sum of total asset sizes of banks in our sample. The black solid line is the sample bank size out as a share of total nation-wide banks' total asset size.



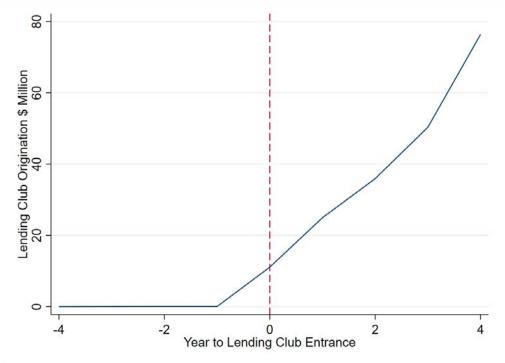
## Figure A2. "Low Mortgage Rate Episode"

The figure show the time-series of aggregate mortgage interest rate and the Federal Funds Rates. "MORT-GAGE30US" is the 30-Year Fixed Rate Mortgage Average in the United States from Freddie Mac. "FED-FUNDS" is the effective Federal Funds Rate by Board of Governors of the Federal Reserve System.



## Figure A3. Lending Club Loan Origination Around Approval of Entrance

This figure shows the loan origination volume of Lending Club around the year that Lending Club was approved to operate in the 9 states in and after 2010. The origination volume is the average loan issuance in million of dollars 1-4 years prior to and 1-4 years after the approval years of the 9 states.



# Figure A4. Lending Club Loan Origination Around Approval of Entrance: by state

This figure shows the loan origination volume of Lending Club around the year that Lending Club was approved to operate in the 9 states in and after 2010. The origination volume is the total loan issuance in million of dollars for each of the 9 states during 2010 and 2018. The red dashed line shows the approval year for each of the nine states.

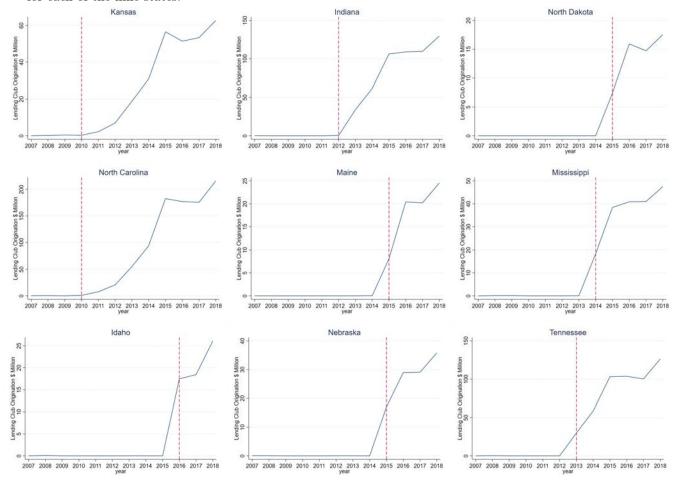


Figure A5. McKinsey (2012) "Breakthrough IT Banking"



BUSINESS TECHNOLOGY OFFICE

## **Breakthrough IT banking**

Some Asian banks achieve superior returns despite relatively low IT expenditures. What's their secret?

Sai Gopalan, Gaurav Jain, Gaurav Kalani, and Jessica Tan Banks have long relied on technology to introduce products such as online banking, ATMs, and mobile payments, and to improve back-office

efficiency. But that reliance comes with a price. Globally, the banking sector spends an average of 4.7 percent to 9.4 percent of operating income on IT, while other sectors spend less: insurance companies and airlines, for example, spend 3.3 percent and 2.6 percent of income, respectively.

Our Asian Banking IT Benchmarking Study¹ finds, however, that a bank's high IT expenditures do not always correlate with superior performance. Some banks with large IT budgets often have trouble leveraging investments to generate commensurately high revenue growth and operational efficiency. Survey data show that 66 percent of banks with higher-than-average IT spending relative to income generated lackluster results, with revenue growth 0.4 percentage points lower than the industry standard and a cost-income (C/I) ratio 2.5 percentage points higher.

By contrast, 23 percent of the 44 banks surveyed outperformed the market on both revenue growth (up 10.9 percentage points) and C/I ratio (down 4.6 percentage points) while spending 29 percent less on IT than other banks in our study. These outperforming banks are more likely to view IT as a strategic enabler, and their investments mirror this outlook. Outperformers direct a higher share of spending toward technologies designed to create new business value and a lower share of spending on support operations, such as finance and human resources. These banks are also more likely than the lower performers to promote efficiency through a consolidated IT footprint as well as formal vendor- and demand-management practices.

The common denominator linking highperforming Asian banks is a commitment to strong governance and spending alignment with the needs of the business. This finding supports our experience with bank clients in Europe and the Americas, and prompted us

<sup>&</sup>lt;sup>1</sup>The 2010 biennial McKinsey Asian Banking IT Benchmarking survey comprised 44 banks across 11 Asia-Pacific countries, with the results tracked against prior year benchmarks from 2006 onward.

## Figure A6. Definition of Different Types of IT Spending

The modeled IT budget for communication services at this site.

It is defined as the network equipment that companies operate to support their communications needs.

### routers

- carrier line equipment fiber optic equipment
- switches
- private branch exchanges (PBXes)
- radio and TV transmitters
- Wi-Fi transmitters
- desktop telephone sets; wide-area network (WAN) and local-area network (LAN) equipment
- videoconferencing and telepresence equipment
- cable boxes
- other network equipment.
- end-client mobile devices like cell phones/iPhones that are bought by individuals

### (a) Figure A

The modeled IT budget for software at this site.

It is defined as software from third parties, whether that software is packaged or semipackaged software delivered on CD and installed within the company, hosted by a third party, offered on a SaaS basis from a multitenant shared-instance server accessible by a browser, or custom-created for a company by third-party

### It includes:

- o license, maintenance, subscriptions and software vendor-provided services revenues for all categories of middleware software (including storage management systems, database management systems, IT management systems, security software, application servers and application development software)
- - o desktop applications
  - o information management software (like business intelligence and enterprise content
  - o process applications (like ERP, CRM, SCM or PLM)

  - o ePurchasing software orisk and payment management software
  - We also include vertical industry applications (like banking management systems, security trading systems, insurance underwriting or claims management software, retail management software, or hospital information systems). Finally, we include computer operating systems software, even though that cost is often bundled
- vertical industry applications (like banking management systems, security trading systems, insurance underwriting or claims management software, retail management software, hospital information systems)
- computer operating systems software (even though that cost is often bundled)

### (b) Figure B

### SERVICES BUDGET

The modeled IT budget for IT-related services at this site.

It is defined as project-based consulting or systems integration services that vendors provide to businesses and Governments, whether on or off-site.

### It includes:

- contractors, consulting services for IT strategy, security assessments and process change
- systems integration 0
- o project services
- mainframe outsourcing, desktop support outsourcing, distributed systems outsourcing, network outsourcing, application hosting, application management outsourcing and application testing. These applications are single-instance software deployments, generally owned rather than subscribed to, and thus are different from SaaS.
- o computer hardware support and maintenance services.

### (c) Figure C

### HARDWARE\_BUDGET

The modeled IT budget for hardware at this site.

It includes the classic computer hardware that IT departments buy and support, regardless of whether the IT department itself operates that equipment (such as servers) or oversees the use of this equipment by employees (such as PCs):

- o PCs: personal computers (laptops, desktops, and tablets)
- o Servers/Mainframes
- Peripherals: monitors, terminals, printers, keyboards, mice, USB devices, etc...
- Storage: storage devices (NAS, DAS, tape)
- Other hardware: hardware specific to the industry (point-of-sales equipment based on PCs, smart cards, embedded computer chips, etc...)

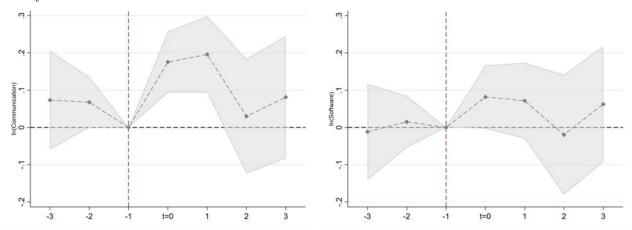
### (d) Figure D

## Figure A7. IT Spending Around Small Business Tax Credit Policy: Alternative IV

This figure reports the event studies of IT spending around the small business tax credit event. The specification is

$$\ln IT_{i,c,t} = \alpha_{i,t} + \alpha_{i,c} + \sum_{s=\in[-3,3],s\neq-1} \beta_s \times \mathbb{1}_{\{\text{t-2014}=s\}} \times \text{High exposure}_{pre} + \Pi_t \times \mathbf{X}_{i,c,t} + \epsilon_{i,c,t}$$

where for bank i at county c in year t,  $\alpha_{i,t}$  are the bank-year fixed effects,  $\alpha_{i,c}$  are the bank-county fixed effects.  $1_{\{t-2014=s\}}$  is a dummy variable that is equal to one if the distance between year t and the event year (2014) is s. "High QSB exposure pre" is equal to one if the average  $\frac{\text{Employees of qualified small businesses}}{\text{Total employees}} \text{ among the top tercile between 2011-2013}, and is equal to zero if the average <math display="block">\frac{\text{Employees of qualified small businesses}}{\text{Total employees}} \text{ in the bottom tercile between 2011-2013}. Bank control variables include banks' revenue per employee. County level control variables include pre-shock unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of small business establishments in nontradables sectors, logarithmic of total small business loan, and GDP per capita. Standard errors are clustered at county level.$ 



## Table A1. Summary statistics of the key variables in regressions

This table presents the summary statistics of the key variables showing up in regression. Revenue is the banks' total income at bank level. "Software/Revenue"/"Communication/Revenue" are the software spending/communication spending scaled by total revenue at bank-year level. "Software" and "Communication" are the dollar amount of software spending or communication spending at bank-county level. " $\ln(\text{Software})$ " and " $\ln(\text{Communication})$ " are the natural log of software spending and communication spending at bank-county level.  $\Delta \ln(\text{Software})_{\text{post}}$ , ( $\Delta \ln(\text{Communication})_{\text{post}}$  and  $\Delta \ln(\text{CRA})_{\text{post}}$ ) are the difference in the natural log of software spending (communication spending and CRA lending) of a bank during 2014-2017 compared to 2011-2013. "CRA amount" is the total dollar volume of CRA loans issued to small businesses with annual revenue below \$1 million at bank-county level during 2011-2017, and the unit is in thousands of dollars. "Refinance/Origination" is the total dollar amount of mortgage refinance scaled by total dollar amount of mortgage origination during 2011-2016. "Mortgage Origination" is the total volume of mortgage origination during 2011-2016 in thousands of dollars. "Establishments" is the total number of establishments in thousands. "Revenue/employee" is in million of dollars per employee. "Non-tradable share" is the share of non-tradable sector establishments among all small business establishments. Non-tradable sector classification follows Mian and Sufi (2014) by using retail related industries and restaurant industries. "Workforce participation" is the total number of employees working in small business establishments scaled by total population.

	Mean	S.d.	25-th	Median	75-th
Revenue (bank-level)	430828.501	4087000.000	8025.000	26120.000	73501.000
Software/Revenue (bank level)	0.0067	0.0580	0.0006	0.0021	0.0055
Communication/Revenue (bank level)	0.0011	0.0070	0.0001	0.0007	0.0014
Software	85192.901	203395.500	2404.276	24695.051	92180.265
Communication	25380.482	424762.900	504.729	12643.667	30798.941
$\ln(\text{Software})$	9.510	2.810	7.785	9.595	11.186
ln(Communication)	7.411	2.895	6.291	7.921	9.277
$\Delta \ln(\text{Software})_{\text{post}}$	0.059	1.106	-0.335	0.139	0.467
$\Delta \ln(\text{Communication})_{\text{post}}$	0.332	1.135	-0.158	0.086	0.548
CRA amount	4523.008	12244.021	353.000	1425.000	4390.000
$\Delta \ln(\mathrm{CRA})_{\mathrm{post}}$	0.258	0.806	-0.436	0.109	0.406
Refinance/Origination	3.012	5.226	0.709	1.354	2.464
ln(Refinance/Origination)	0.353	0.901	-0.257	0.341	0.928
Mortgage Origination	16571.100	42038.850	773.000	2736.500	10537.000
Revenue/Employee	0.287	0.380	0.165	0.227	0.325
Real GDP per capita	48.731	43.453	31.672	42.273	56.157
Non-tradable share	0.641	0.113	0.569	0.640	0.711
Establishment	10226.778	25170.742	653.000	2176.000	9011.000
Workforce participation	0.316	0.149	0.216	0.301	0.396
Unemployment rate	6.130	2.693	4.100	5.600	7.700
Equity/Assets	0.113 - 6	9  0.027	0.096	0.109	0.126
Deposits/Total assets	0.799	0.089	0.758	0.815	0.863

## Table A2. Summary Statistics of Banks' IT Spending by Bank Size Group

This table presents the summary statistics of banks' IT spending by banks' size groups. Banks in the sample are split into five groups. Total IT Spending is the sum of all types of IT spending in millions of dollars. No. of IT employees is the total amount of IT-related employees. IT Spending/Revenue is total IT Spending scaled by banks' total income, IT Spending/Non-interest expense (IT/NIE) is total IT spending scaled by non-interest expenses. The different categories of IT spending are the four categories of IT spending scaled by total IT spending.

	Mean	S.d.	Median		Mean	S.d.	Median
< \$100 Million				\$100 Million-\$1 Billion			
IT Spending/Total Assets	0.001	0.007	0.000	IT Spending/Total Assets	0.001	0.001	0.000
IT Spending/Income	0.015	0.023	0.007	IT Spending/Income	0.015	0.035	0.007
IT Spending/NIE	0.019	0.035	0.010	IT Spending/NIE	0.021	0.032	0.010
Communication/Total	0.126	0.138	0.079	Communication/Total	0.096	0.113	0.055
Software/Total	0.317	0.136	0.342	Software/Total	0.314	0.161	0.347
Services/Total	0.311	0.121	0.333	Hardware/Total	0.173	0.108	0.167
Hardware/Total	0.188	0.113	0.191	Services/Total	0.313	0.131	0.326
Other/Total	0.091	0.148	0.050	Other/Total	0.055	0.060	0.020
	Mean	S.d.	Median		Mean	S.d.	Median
\$1 Billion-\$10 Billion	Mean	S.d.	Median	\$10 Billion-\$250 Billion	Mean	S.d.	Median
\$1 Billion-\$10 Billion IT Spending/Total Assets	Mean 0.001	S.d. 0.002	Median 0.000	\$10 Billion-\$250 Billion IT Spending/Total Assets	Mean 0.001	S.d. 0.001	Median 0.000
				· ·			
IT Spending/Total Assets	0.001	0.002	0.000	IT Spending/Total Assets IT Spending/Income IT Spending/NIE	0.001	0.001	0.000
IT Spending/Total Assets IT Spending/Income	0.001 0.016	0.002 0.028	0.000 0.008	IT Spending/Total Assets IT Spending/Income	0.001 0.023	0.001 0.092	0.000 0.008
IT Spending/Total Assets IT Spending/Income IT Spending/NIE	0.001 0.016 0.022	0.002 0.028 0.027	0.000 0.008 0.012	IT Spending/Total Assets IT Spending/Income IT Spending/NIE	0.001 0.023 0.029	0.001 0.092 0.060	0.000 0.008 0.013
IT Spending/Total Assets IT Spending/Income IT Spending/NIE Communication/Total	0.001 0.016 0.022 0.068	0.002 0.028 0.027 0.082	0.000 0.008 0.012 0.042	IT Spending/Total Assets IT Spending/Income IT Spending/NIE Hardware/Total Communication/Total Software/Total	0.001 0.023 0.029 0.159	0.001 0.092 0.060 0.124	0.000 0.008 0.013 0.125
IT Spending/Total Assets IT Spending/Income IT Spending/NIE Communication/Total Software/Total	0.001 0.016 0.022 0.068 0.308	0.002 0.028 0.027 0.082 0.157	0.000 0.008 0.012 0.042 0.237	IT Spending/Total Assets IT Spending/Income IT Spending/NIE Hardware/Total Communication/Total	0.001 0.023 0.029 0.159 0.065	0.001 0.092 0.060 0.124 0.081	0.000 0.008 0.013 0.125 0.043

	Mean	S.d.	Median
> \$250 Billion			
IT Spending/Total Assets	0.002	0.002	0.001
IT Spending/Income	0.048	0.080	0.022
IT Spending/NIE	0.079	0.128	0.038
Communication/Total	0.062	0.038	0.056
Software/Total	0.310	0.217	0.265
Hardware/Total	0.154	0.094	0.134
Services/Total	0.276	0.412	0.134
Other/Total	0.036	0.061	0.012

# Table A3. Bank Characteristics and Banks' IT Spending: IT scaled by total deposits

This table presents the results of correlation between banks' IT spending and banks' characteristics. The regression specification is as follows.

$$\frac{\text{Type S IT spending}}{\text{Deposits}}_{i.10\text{-}19} = \alpha + \beta \frac{\text{Type L loan}}{\text{Total loan}}_{i,10\text{-}19} \text{ or (Bank Char)} + \gamma \mathbf{X} + \epsilon_i$$

Panel A shows how the banks' loan specialization correlates with banks' IT spending. Type L loan/Total Loan is the average of a specific type of loan scaled by total loan. Among them, Personal loan/Total Loan is the sum of personal loans and real estate loans to 1-4 family units scaled by total loan; Agriculture/Total loan is the agricultural loan scaled by total loan; CRA/Total loan is the sum of small business loans reported in CRA scaled by total loan; "Other C&I/Total loan" is the total C&I loan minus small business loans reported in CRA, scaled by total loan; "Mortgage refinance" is the total amount of mortgage refinance reported in HMDA scaled by the bank's total loan; "Other personal loans" is the deduction of "Mortgage refinance" from "Personal and mortgage loans." %Refinance is the frequency of refinance as a percent of total number of mortgage issuance that are reported in HMDA. "XX IT/Deposits" are total dollar amount of IT spending scaled by total deposits. Panel B shows how a bank's hierarchical structure correlates with its IT spending. "Hierarchical layer" is the number of types of its locations. "In(num offices)" is the logarithmic of total number of offices. Control variables include net income scaled by total assets, deposits scaled by total assets, revenue per employee, salaries scaled by total assets and equity scaled by total assets. Fixed effects include bank size group, and banks' headquarter state fixed effects. Panel C shows how a bank's role in the syndicated loan market correlates with its IT spending. "Lead bank is the frequency of a bank's showing up as a lead bank in the syndicated loan market as a share of total number of syndicated loans lent out. All of the loan profile variables are calculated as the average of the loan profile of a bank between 2010 and 2019.

\*\*\*\*, \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

Panel A: Loan Specialization						
	Software/Deposits	Communication/Deposits	Hardware/Deposits	Services/Deposits		
	(1)	(2)	(3)	(4)		
C& I loan/Total loan	0.0121	0.0203	0.0329	-0.0212		
,	(0.0279)	(0.0275)	(0.0278)	(0.0278)		
CRA/Total loan	-0.226***	0.101***	0.0331	-0.000236		
	(0.0334)	(0.0332)	(0.0338)	(0.0338)		
Other C&I loan/Total loan	0.0392	0.00988	0.0306	-0.0199		
	(0.0278)	(0.0274)	(0.0277)	(0.0277)		
Personal loan/Total loan	0.0463***	0.0194	0.0249	0.00213		
7	(0.0173)	(0.0170)	(0.0173)	(0.0173)		
Refinance/Total loan	0.0719**	0.0340	$0.0499^*$	-0.00574		
	(0.0293)	(0.0285)	(0.0292)	(0.0292)		
Other personal loans/Total loan	-0.00296	-0.00728	-0.0228	0.0209		
	(0.0321)	(0.0311)	(0.0320)	(0.0320)		
Refinance frequency	0.0774***	0.0428	0.0706**	$0.0562^*$		
	(0.0299)	(0.0285)	(0.0288)	(0.0296)		
Agriculture loan/Total loan	-0.236	-0.0808	-0.345	-0.0828		
<u> </u>	(0.330)	(0.325)	(0.329)	(0.329)		
Pa	anel B: Hierarchica	al Complexity and IT Sp	ending			
Hierarchical layer	-0.00409	$0.0475^{*}$	-0.00377	-0.00331		
·	(0.0257)	(0.0254)	(0.0257)	(0.0258)		
ln(num of offices)	$0.0586^{'}$	$0.104^{**}$	0.151***	0.0690		
,	(0.0521)	(0.0517)	(0.0521)	(0.0517)		
	Panel C: Banks'	Role in Syndicated Lend	ing			
% Lead bank/Total syndicate	0.0310*	0.0408**	0.0222	0.00746		
, ,	(0.0181)	(0.0180)	(0.0183)	(0.0181)		

### Table A4. Small Business Loan, Mortgage Refinance and Bank IT Spending: II

This table presents regression results of banks' IT spending on their mortgage refinance (or small business lending) and relevant control variables at bank-county-year level. The sample period is 2010 to 2019.

$$\frac{\text{Type S IT spending}}{\text{Revenue}}_{i,c,t} = \alpha_i + \mu_{c,t} + \beta \ln(\text{refinance})_{i,c,t} \text{ or } \ln(\text{CRA})_{i,c,t} + \gamma X + \epsilon_{i,c,t}$$

In refinance i,c,t is the natural logarithm of newly issued mortgage refinance of bank i at county c in year t in reported in HMDA, ln CRAi,c,t is the natural logarithm of small business loans issued by bank i in county c and in year t. Software (Communication)/Revenue is software (communication) spending scaled by total revenue, measured at bank-county-year level. Control variables include net income scaled by total assets, deposits scaled by total assets, salaries scaled by total assets and equity scaled by total assets, control variables are measured at bank-year level. Fixed effects include bank fixed effects, county fixed effects and year fixed effects (or county  $\times$  year fixed effects). Standard errors are clustered at county and bank level. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

		Software/Re	Panel A	Cov	mmunication	/Povonuo	
				Communication/Revenue			
	(1)	(2)	(3)	(4)	(5)	(6)	
Ln(Mortgage refinance)	0.0325****	0.0255***	0.0222**	0.0148**	0.00569	0.0104	
	(0.00765)	(0.00935)	(0.00894)	(0.00695)	(0.00665)	(0.00679)	
Revenue per Employee		-0.121***	-0.0928***		-0.123***	-0.0951***	
		(0.00763)	(0.00733)		(0.00663)	(0.00652)	
Net income/Assets		-0.0115	-0.00791		0.00226	0.00474	
		(0.00829)	(0.00755)		(0.00747)	(0.00735)	
Equity/Assets		-0.0280**	-0.0282**		-0.0212*	-0.0161*	
,		(0.0135)	(0.0120)		(0.0112)	(0.00946)	
Deposits/Assets		0.00645	0.000818		-0.0103	-0.00686	
-		(0.0189)	(0.0157)		(0.0149)	(0.0117)	
Salaries/Assets		-0.926	-1.036		-0.493	-0.655	
		(0.809)	(0.794)		(0.756)	(0.824)	
Bank FE	Y	Y	Y	Y	Y	Y	
Controls	N	Y	Y	N	Y	Y	
Other fixed effects	Count	y, Year	$\mathbf{County}{\times}\mathbf{Year}$	Count	y, Year	$\mathbf{County}{\times}\mathbf{Year}$	
AdR-squared	0.481	0.473	0.499	0.508	0.501	0.523	
N	163934	149899	145848	163965	149940	145879	

	Panel B Software/Revenue			Communication/Revenue			
	(1)	(2)	(3)	(4)	(5)	(6)	
Ln(CRA)	-0.0290**	-0.0255	-0.0139*	0.0174**	0.0180***	0.0214***	
	(0.0136)	(0.0163)	(0.00805)	(0.00779)	(0.00683)	(0.00678)	
Revenue per Employee		-0.115***	-0.0929***		0.00337	0.0180	
		(0.00727)	(0.00722)		(0.0196)	(0.0187)	
Net income/Assets		0.0141	-0.00679		-0.00331	-0.000964	
		(0.00981)	(0.00707)		(0.00986)	(0.00941)	
Equity/Assets		-0.0331***	-0.0296**		-0.0183	-0.0139	
		(0.0125)	(0.0116)		(0.0112)	(0.00958)	
Deposits/Assets		-0.0542**	-0.00449		-0.00919	-0.00548	
		(0.0219)	(0.0152)		(0.0146)	(0.0115)	
Salaries/Assets		-0.318	-0.567		-0.179	-0.294	
		(0.688)	(0.652)		(0.769)	(0.818)	
Bank FE	Y	Y	Y	Y	Y	Y	
Controls	N	Y	Y	N	Y	Y	
Other fixed effects	Count	y, Year	$\mathbf{County}{\times}\mathbf{Year}$	Count	y, Year	$\mathbf{County}{\times}\mathbf{Year}$	
AdR-squared	0.489	0.467	72  0.507	0.530	0.518	0.539	
N	167452	169177	148936	167483	152888	148967	

### Table A5. Soft Information and Banks' IT Spending: Alternative IV

$$\begin{split} \Delta \ln(\text{CRA})_{i,c,post} &= \tilde{\alpha_i} + \mu_1 \times \frac{\text{Emp of qualified small businesses establishments}}{\text{Total employee}} \\ &\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \widehat{\Delta \ln(\text{CRA})}_{i,c,post} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c} \end{split}$$

The last two columns show the following OLS specification:

$$\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \Delta \ln(\text{CRA})_{i,c,post} + \mu_c + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln (CRA)_{i,c,post}$  is the change in average natural log of small business loans reported in CRA of bank i at county c during the years 2014-2017 compared with 2011-2013. Bank control variables include pre-shock revenue per employee. County level control variables include the pre-shock unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of small business establishments in non-tradable sector and GDP per capita. Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	First stage	$\ln(\text{Software})$	ln(Communication)	$\ln(\mathrm{Software})(\mathrm{OLS})$	$\ln({\rm Communication})({\rm OLS})$
	(1)	(2)	(3)	(4)	
Emp of qualified small businesses establishments Total employee	0.053***				
rotar employee c,pre	(0.012)				
$\Delta \widehat{\ln(\mathrm{CRA})}$		0.086	$0.454^{*}$		
,		(0.248)	(0.270)		
$\Delta \ln(\mathrm{CRA})$		, ,	, ,	0.004	$0.019^{*}$
				(0.010)	(0.011)
Bank FE	Y	Y	Y	Y	Y
Clustered	Y	Y	Y	Y	Y
AdR-squared	0.435	-0.183	-0.332	0.119	0.103
N	17,941	17,937	17,936	19,548	19,547

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

Table A6. Economic magnitude of the  $\mathrm{Did/IV}$  regressions

Analysis and Table	Economic magnitude estimation
IT spending and banks' loan profile (Column (2) of Table 3)	As is shown in Table A1 the standard deviation of communication/revenue is 0.007, and the average revenue is \$429 million. Given a coefficient of 0.08, the implied increase of communication spending is $0.007 \times 0.08 \times $429$ million=\$0.24 million. In terms of the communication as a share of total revenue, a one standard deviation increase in the C&I loan as a share of total personal loan (13 percentage points) is associated with $0.007 \times 0.08 = 5.6$ basis points higher communication spending as a share of total revenue. Compared with the average "C&I loan/total loan" and the average size of "communication/revenue", increasing "C&I loan/total loan" from 0.19 to 0.32 (68% increase compared to mean 0.19) leads to "communication/revenue" to increase from 0.0011 to 0.00166 (0.51% increase compared to mean 0.0011).
IT spending responses to SHOP tax credit jump (Table 6)	For the first stage, we estimate a one standard deviation of increase in qualified business share leads to 1.032 standard deviation in $\Delta \ln(\text{CRA})$ (standard deviation 0.86, mean 0.258), or $\Delta \ln(\text{CRA})$ increases by $0.86 \times 1.032 = 0.88$ from the mean of 0.258 to 1.145. This means CRA loan (mean of CRA loans at bank-county level is \$4.523 million) increases by an extra $(e^{1.145} - e^{0.258}) \times $4.523$ million, or \$8.366 million. With interest rate spread at 2%, this amounts to \$167K. Similarly, in the second stage of the regression, $\Delta \ln(\text{communication})$ (standard deviation 1.135, mean 0.317) increases by 0.67 standard deviation, which means communication spending increases by an extra $(e^{1.135 \times 0.67 + 0.317} - e^{0.317}) \times $25,380 = $28,183$ . Compared to the extra increase in CRA loans, this means \$40,298.79/\$167,000 = 24.1%
IT spending responses mortgage refinancing increases (Table 8)	In our data $\ln(\text{Refinance}/\text{Origination})$ has a standard deviation of 0.650, which means "Refinance/Origination" (mean 3.01) increase by $(e^{0.650}-1)\times 3.01=2.75$ from mean. Similarly, as the standard deviation of $\ln(\text{software})$ is 2.81, a 0.313-standard-deviation increase means software spending (mean \$85,192.90) increases by $((e^{2.81\times0.315}-1)\times \$85,192.90=\$133.26K)$ from the mean. The average size of mortgage origination at bank-county level is \$16,571 K, therefore, "Refinance/Origination" increasing by 2.75 means refinancing increases by \$16,571K \times 2.75=\\$455,570K. Banks typically charge around 1%-2% of the loan amount as fees (see this report for reference). Consequently, from the \$455,570K banks will be able to generate extra \$455.57K in revenue assuming they charge 1% of fee. The increased software spending of \$133.26K squares with the increased revenue.

### Table A7. Impact of SHOP program on small businesses: by business sizes

This table presents the estimated impact of small business health insurance tax credit on small businesses with different employee sizes. The estimated annual revenue and net profit come from the research summary by Zippia. "Ave num of emp" is the average number of employees for each firm size bracket, for example, "ave num of emp" for firms of size "2-4 employees" is 3. Net profit is calculated as the profit margin multiplied by the "Ave revenue." As is estimated the research summary by Zippia, the estimated profit margin for small businesses are around 7%-10%, and we use 8.5% as the average profit margin in our calculation. The "Ave health premium" is the average health insurance package per employee contributed by employers, and the estimates comes from the survey by Kaiser Family Foundation Report.

Panel A:Firm size	Ave num of emp	Ave Revenue	Net profit	Ave health premium	Total health	50% credit	35% credit	\$ Change	Saving
					premium costs				of tota
1 employee	1	44000	3740	6485	6485	3242.5	2269.75	972.75	0.260
2-4 employee	3	387000	32895	6485	19455	9727.5	6809.25	2918.25	0.088
5-9 employee	7	1080000	91800	6485	45395	22697.5	15888.25	6809.25	0.074
10-19 employee	15	2164000	183940	6485	97275	48637.5	34046.25	14591.25	0.080
20-99 employee	60	7124000	605540	6485	389100	194550	136185	58365	0.096
100-499 employee	300	40775000	3465875	6485	1945500	972750	680925	291825	0.084

Panel B:		
small business employee size	average \$ saving/net revenue	share among qualified
1-4 workers	0.17	0.64
5-9 workers	0.08	0.22
10-19 workers	0.07	0.14
Weighted average \$ savings/net revenue for QSB	0.14	

### Table A8. Growth of small business establishments by size around the event

This table presents the growth of total number of small businesses of different sizes around in 2014, with the following regression specification:

$$\Delta \ln(\text{Num small business of size S})_{c,post} = \alpha + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln(\text{Number of small business establishments of size S})_{c,post}$  is the change in natural log of small business establishments of size "S" in county c during the years 2014-2017 compared with 2011-2013. County level control variables include the pre-shock unemployment growth, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of non-tradable sector business establishments, real GDP per capita, and average of local banks' revenue per employee. Robust standard errors are reported in the parentheses. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

Panel A							
	% growth est 1-4	$\frac{\%}{\%}$ growth est 5-9	% growth est 10-19	% growth est 20-49	% growth est 50-99	% growth est 100-249	% growth est> 249
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
$ \frac{\text{\# Qualified small business est}}{\text{Total \# of establishments}} \right)_{c,pre} $	-0.079***	0.133***	0.051**	0.046	-0.048**	-0.086***	-0.055
· /K · ·	(0.018)	(0.019)	(0.020)	(0.021)	(0.023)	(0.045)	(0.042)
Controls	N	N	N	N	N	N	N
R-squared	0.006	0.016	0.002	0.001	0.001	0.004	0.001
N	3,060	3,064	3,057	3,050	2,977	2,802	$2,\!374$
Panel B							
	% growth 1-4	% growth 5-9	% growth 10-19	% growth 20-49	% growth 50-99	% growth 100-249	% growth $> 249$
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	-0.037	0.163***	0.160***	0.079	-0.020	-0.052	0.055
-	(0.038)	(0.037)	(0.050)	(0.059)	(0.050)	(0.038)	(0.056)
L.Revenue/Employee	-0.039**	0.006	-0.009	-0.037*	-0.028	0.028	0.015
	(0.017)	(0.017)	(0.020)	(0.021)	(0.023)	(0.025)	(0.027)
L.Population growth	$0.271^{***}$	$0.224^{***}$	$0.169^{***}$	$0.150^{***}$	0.098***	0.098***	0.103***
	(0.026)	(0.024)	(0.031)	(0.023)	(0.031)	(0.033)	(0.033)
L.Work force participation	-0.098***	-0.091***	-0.039	-0.095***	-0.066**	-0.082***	-0.001
	(0.032)	(0.025)	(0.027)	(0.024)	(0.026)	(0.029)	(0.029)
L.ln(Total establishments)	$0.160^{***}$	$0.107^{***}$	$0.145^{***}$	$0.113^{***}$	$0.087^{***}$	0.095***	0.109***
	(0.030)	(0.029)	(0.032)	(0.029)	(0.031)	(0.035)	(0.035)
L.Non-tradable share	0.002	-0.021	-0.017	0.075**	0.018	0.013	-0.037
	(0.024)	(0.023)	(0.028)	(0.030)	(0.027)	(0.026)	(0.031)
L.Unemployment growth	-0.026	-0.064***	-0.051*	-0.124***	-0.037	-0.087***	-0.097***
	(0.023)	(0.022)	(0.028)	(0.023)	(0.029)	(0.024)	(0.028)
L.GDP per capita	1.013***	$0.535^{**}$	$0.868^{**}$	1.041***	0.051	$0.836^{**}$	$0.906^{***}$
	(0.280)	(0.272)	(0.356)	(0.352)	(0.348)	(0.365)	(0.301)
Controls	Y	Y	Y	Y	Y	Y	Y
R-squared	0.127	0.083	0.059	0.072	0.016	0.030	0.039
N	2,942	2,945	2,945	2,938	2,888	2,724	2,302

### Table A9. Employment growth by small business sizes around the event

This table presents the growth of total number of small business establishments of different sizes around the activation of SHOP in 2014, the regression specification is as follows:

$$\Delta \ln(\text{Employees of establishments of size S})_{c,post} = \alpha + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln(\text{Employees of establishments of size S})_{c,post}$  is the change in natural log of employees of business establishments with size "S" in county c during the years 2014-2017 compared with 2011-2013. County level control variables include the pre-shock unemployment growth, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of non-tradable sector business establishments, real GDP per capita, and average of local banks' revenue per employee. Robust standard errors are reported in the parentheses. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	$\Delta \mathrm{Em}$	np 1-19	$\Delta \text{Emp} > 20$		$\Delta { m Em}$	p500+
	(1)	(2)	$\overline{(3)}$	(4)	$\overline{(5)}$	(6)
	0.109**	0.243***	-0.012	0.067	0.041	0.032
c,pre	(0.044)	(0.032)	(0.020)	(0.057)	(0.033)	(0.045)
L.Revenue/Employee		-0.011		-0.039*		-0.003
		(0.020)		(0.021)		(0.023)
L.Population growth		0.188***		0.079***		0.104***
		(0.025)		(0.030)		(0.029)
L. Work force participation		-0.100***		-0.054**		$-0.049^*$
		(0.023)		(0.027)		(0.026)
L.ln(Total establishments)		$0.286^{***}$		$0.216^{***}$		-0.019
		(0.025)		(0.026)		(0.029)
L.Non-tradable share		0.001		-0.000		-0.026
		(0.023)		(0.022)		(0.026)
L.Unemployment growth		-0.153***		-0.101***		-0.171***
		(0.025)		(0.027)		(0.028)
L.GDP per capita		-0.102		0.642***		0.954**
		(0.252)		(0.223)		(0.381)
Clustered	Y	Y	Y	Y	Y	Y
AdR-squared	0.012	0.151	-0.000	0.055	0.001	0.045
N	3,043	2,934	3,024	2,927	2,972	2,883

### Table A10. Growth of small businesses with different employee sizes

This table presents the correlation between the change in growth of small businesses with 1-4 employees and change in growth of small businesses with 5-9 employees or 10-19 employees.

 $\Delta \ln(\text{Num small business of size } 1-4)_{c,post} = \alpha + \mu_1 \times \Delta \ln(\text{Num small business of size } S)_{c,post} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$ 

 $\Delta \ln(\text{Number of small business establishments of size S})_{c,post}$  is the change in natural log of small business establishments of size "S" in county c during the years 2014-2017 compared with 2011-2013. County level control variables include pre-shock average of revenue per employee banks in a county, unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, logarithmic of total small business loan, and GDP per capita. Robust standard errors are reported in the parentheses. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

		% growth		
	$\overline{(1)}$	$\overline{(2)}$	$\overline{(3)}$	$\overline{(4)}$
% growth est 5-9	-0.085***	-0.132***		
	(0.028)	(0.031)		
% growth est 10-19			-0.033	-0.033
			(0.044)	(0.034)
L.Revenu/Employee		-0.039**		-0.041**
		(0.019)		(0.019)
L.Population growth		$0.210^{***}$		$0.193^{***}$
		(0.027)		(0.026)
L.Work force participation		-0.077**		-0.062**
		(0.031)		(0.030)
L.ln(Total establishments)		$0.194^{***}$		0.184***
		(0.031)		(0.032)
L.Non-tradable share		0.012		0.014
		(0.026)		(0.027)
L.Unemployment growth		-0.059***		-0.054**
		(0.022)		(0.022)
L.GDP per capita		$0.659^{**}$		$0.635^{**}$
		(0.263)		(0.258)
Controls	N	Y	N	Y
Clustered	Y	Y	Y	Y
AdR-squared	0.007	0.095	0.001	0.081
N	3,044	2,934	3,035	2,933

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

### Table A11. Placebo: Small Business Loan Event using 2018 as event year

This table displays the placebo test outcomes for the small business health care tax credit event, assuming the year 2018 as the the event year of analysis. The regression specifications are the following specification:

$$\begin{split} &\Delta \ln(\text{CRA})_{i,c,post} = \tilde{\alpha_i} + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c} \\ &\Delta \ln(\text{IT})_{i,c,post} = \alpha_i + \beta \times \widehat{\Delta \ln(\text{CRA})}_{i,c,post} + \gamma \mathbf{X}_{i,c} + \epsilon_{i,c} \end{split}$$

Aln (CRA)<sub>i,c,post</sub> is the change in average natural log of small business loans reported in CRA of bank i at county c during the years 2018-2019 compared with 2016-2017. Bank control variables include pre-shock revenue per employee of the bank in a county. County level control variables include the pre-shock unemployment growth rate, labor force participation rate, population growth rate, logarithmic of total number of establishments, logarithmic of total small business loan, share of nontradable sector establishments, and GDP per capita. Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	First stage	$\ln(\text{Software})$	ln(Communication)
	(1)	$\overline{(2)}$	$\overline{\qquad \qquad }(3)$
Qualified small businesses establishments Total establishments c.pre	-0.852	, ,	· /
Total establishments $c,pre$	(0.521)		
$\Delta \ln(\text{CRA})$	,	-0.648	-0.503
,		(0.553)	(0.614)
L.Revenue/Emp	-0.027***	-0.078***	-0.084***
, -	(0.005)	(0.018)	(0.020)
L.Work force participation	0.173***	0.147	0.144
	(0.054)	(0.140)	(0.147)
L.Unemployment growth	-0.198*	-0.280	-0.104
	(0.111)	(0.206)	(0.211)
L.ln(Total establishments)	-0.319***	-0.201	-0.165
	(0.012)	(0.171)	(0.189)
L.GDP per capita	-0.185***	-0.013	-0.017
	(0.020)	(0.112)	(0.123)
L.Population growth	0.016**	0.004	0.008
	(0.007)	(0.014)	(0.015)
L.Non-tradable share	-0.175**	0.013	-0.002
	(0.078)	(0.161)	(0.177)
Bank FE	Y	Y	Y
Clustered	Y	Y	Y
F-stat	2.623		
AdR-squared	0.687	-0.381	-0.303
N	14,593	14,585	14,576

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

### Table A12. Responses of mortgage loans around the activation of SHOP 2014

This table presents the placebo tests for the impact of "QSB share" on the change in other types of credit demand in the local economy around the activation of SHOP in 2014. The regression specification is as follows:

$$\Delta \ln(\text{Mortgage loans})_{i,c,post} = \alpha_i + \mu_1 \times \left(\frac{\# \text{ Qualified small business est}}{\text{Total } \# \text{ of establishments}}\right)_{c,pre} + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$$

 $\Delta \ln(\text{Mortgage loans})_{i,c,post}$  is the change in natural log of total mortgage/mortgage origination/mortgage refinance by bank i in county c during the years 2014-2017 compared with 2011-2013. Bank control variables include pre-shock revenue per employee of the bank in a county. County level control variables include the pre-shock unemployment growth, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of non-tradable sector business establishments, and real GDP per capita. Fixed effects include bank fixed effects. Robust standard errors are reported in the parentheses. \*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	$\Delta$ ln(total mortgage)	$\Delta$ ln(origination)	$\Delta  \ln({\rm refinance})$
	(1)	$\phantom{aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa$	$\overline{\qquad \qquad }(3)$
	0.057	0.134	0.049
e,pre	(0.200)	(0.147)	(0.241)
L.Revenue/Employee	-0.033***	0.002	-0.037***
, -	(0.010)	(0.008)	(0.012)
L.Unemployment growth	-0.493***	$0.502^{***}$	-0.969***
	(0.137)	(0.120)	(0.180)
L.GDP per capita	0.001	$0.010^{*}$	-0.005
	(0.007)	(0.006)	(0.009)
L.Work force participation	-0.496***	-0.066	-0.513***
	(0.064)	(0.053)	(0.083)
L.Non-tradable share	-0.022	-0.185***	0.137
	(0.069)	(0.060)	(0.088)
L.ln(Total establishments)	0.052***	-0.029***	0.080***
	(0.010)	(0.008)	(0.013)
L.Population growth	$0.067^{***}$	-0.019***	0.084***
	(0.009)	(0.006)	(0.011)
Bank FE	Y	Y	Y
Clustered	Y	Y	Y
AdR-squared	0.429	0.232	0.414
N	15,788	14,994	14,994

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

### Table A13. Hard information and IT: with house price in control

The table presents the response of IT spending to changes in mortgage refinance intensity during the low interest rate episode, with house price index (HPI) in 2010 in as control variables. The 2SLS specification is as below:

ln(Refinance/Origination)<sub>i,c</sub> = 
$$\tilde{\alpha}_i + \mu_1 \times \Delta \text{Payments}_c + \mu_2 \ \mathbf{X}_{i,c} + \epsilon_{i,c}$$
  
ln(Type S Spending)<sub>i,c</sub> =  $\alpha_i + \beta \times \ln(\text{Refinance/Origination})_{i,c} + \gamma \ \mathbf{X}_{i,c} + \epsilon_{i,c}$ 

 $ln(Type~S~Spending)_{i,c}, ln(Refinance/Origination)_{i,c}$ , Payments gap, and Mortgage rates are as defined in Table 8. Bank control variables include banks' revenue per end of the bank in a county. County level control variables include the unemployment growth, labor force participation rate, population growth rate, logarithmic of total modes of establishments, share of non-tradable sector business establishments, real GDP per capita, and specifically, we control for the county's house price index in 2010 Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*\*, \*\*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	First stage	$\ln(\mathrm{Software})$	ln(Communication)	First stage	$\ln(\mathrm{Software})$	$\ln({\rm Communication})$
	(1)	(2)	(3)	(4)	(5)	(6)
$\Delta$ Payment <sub>c</sub>	0.994***					
	(0.233)					
$\Delta { m Mortgage\ rate}_c$				1.875***		
				(0.624)		
Ln(Refinance/Origination)		$0.436^{**}$	0.207		$0.436^{*}$	0.337
, - ,		(0.195)	(0.212)		(0.261)	(0.337)
HPI	-0.000	0.002***	0.001**	0.001	0.002***	0.001
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
Unemployment growth	-0.670**	$0.687^{**}$	$1.353^{*}$	-0.811***	$0.687^{**}$	1.705
	(0.296)	(0.267)	(0.788)	(0.299)	(0.308)	(1.087)
Real GDP per Capita	-0.036**	-0.003	-0.001	-0.036**	-0.006	-0.004
	(0.014)	(0.017)	(0.016)	(0.014)	(0.016)	(0.015)
Revenue/Emp	0.054***	$0.046^{***}$	$0.064^{***}$	0.055****	$0.046^{**}$	$0.057^{**}$
	(0.018)	(0.017)	(0.016)	(0.018)	(0.019)	(0.023)
Workforce Participation	-0.096	0.876***	1.064***	-0.113	0.876***	1.085***
	(0.115)	(0.126)	(0.088)	(0.116)	(0.126)	(0.104)
Population growth	-0.140***	$0.073^{**}$	$0.054^{*}$	-0.133***	$0.073^{**}$	0.070
	(0.018)	(0.030)	(0.031)	(0.018)	(0.036)	(0.044)
$ln(Total\ establishments)$	0.892**	-2.204***	-2.262***	0.900**	-2.175***	-2.227***
	(0.374)	(0.552)	(0.543)	(0.375)	(0.536)	(0.527)
Non-tradable share	-0.025*	0.081***	0.081***	-0.025*	0.078***	0.078***
	(0.014)	(0.016)	(0.016)	(0.014)	(0.016)	(0.015)
Bank FE	Y	Y	Y	Y	Y	Y
Clustered	Y	Y	Y	Y	Y	Y
F-stat	11.716			11.538		
AdR-squared	0.426	-0.239	0.117	0.425	-0.239	-0.049
N	14,263	14,263	154,263	14,263	14,263	14,263

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

### Table A14. Correlation between mortgage payment savings and household wealth

This table presents the correlation between  $\Delta$ Payments or local households' mortgage repayment savings to local household income and the banks' expansion responses to local household income. The regression equations are as follows

$$\begin{array}{rcl} \text{Real GDP per capita}_c & = & \alpha + \beta \times \Delta \text{Payments}_c + \epsilon \\ \Delta \text{Branches}_{i,c}/\Delta \text{Employees}_{i,c} & = & \alpha_i + \beta \times \text{Real GDP per capita}_c/\Delta \text{Payments}_c + \epsilon_{i,c} \end{array}$$

 $\Delta$ Payments<sub>c</sub> is local households' mortgage repayment savings during the low interest rate episodes in a county c. The detailed explanation of is provided in Section 4.3 of the main text.  $\Delta$ Branches<sub>i,c</sub> is the log difference of bank i's average number of branches in 2011-2016 compared to the number of branches in 2010 in county c.  $\Delta$ Employees<sub>i,c</sub> is the log difference of bank i's average number of employees in 2011-2016 compared to the number of employees in 2010 in county c \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	GDP per capita (2011-2013)	GDP per capita (2011-2016)
	(1)	$\frac{}{(2)}$
$\Delta$ Payment <sub>c</sub>	0.007	0.025
	(0.018)	(0.018)
AdR-squared	-0.000	0.000
N	3,001	3,035

	$\Delta$ Branches	$\Delta$ Employees	$\Delta$ Branches	$\Delta$ Employees
	(1)	$\frac{}{(2)}$	$\overline{(3)}$	(4)
Real GDP per capita (2011-2016)	0.059	0.076		
	(0.047)	(0.052)		
$\Delta$ Payment <sub>c</sub>			0.022	$0.025^{*}$
			(0.015)	(0.014)
Bank FE	Y	Y	Y	Y
AdR-squared	0.057	0.022	0.009	0.012
N	15,615	15,610	16,537	16,531

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

Table A15. Correlation between mortgage repayment savings (mortgage refinance rate savings) and county-level characteristics.

This table presents the correlation between the  $\Delta Payment_c$  ( $\Delta Mortgage rate_c$ ) and the major county-level characteristics variables.  $\Delta Payment$  is the county level average mortgage repayment savings during the low interest rate episode.  $\Delta Mortgage$  rate is the county level average mortgage rate savings during the low interest rate episode.  $\Delta HPI$  is the average annual change in house price during 2011-2016.  $HPI_{2010}$  is the county level house price index in 2010. "Non-tradable share" is the share of small businesses in non-tradable sector. "Workforce participation rate" is the total number of employees working in small business establishments as a share of total population. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	$\Delta { m Payment}$	$\Delta$ Mortgage rate
Δ Payment	1	
$\Delta$ Mortgage rate	0.43***	1
$\Delta \mathrm{Delta}\ \mathrm{HPI}$	0.12***	0.07
$\mathrm{HPI}_{2010}$	0.26***	0.05
$\log(\text{Establishment})$	0.01	-0.02
Unemployment rate	0.11***	0.08*
Real GDP per Capita	0.06	-0.05
Non-tradable Share	0.20***	0.08*
Workforce Participation Rate	0.05	-0.02

### Table A16. Payment gap and small business loan during low interest rate period

This table presents the correlation of issuance of mortgage refinance loans (small business loans) and the local households' mortgage repayment savings during the low interest rate episode.

(Refinance (CRA)/Deposits)<sub>i,c</sub>or (Refinance (CRA)/Revenue)<sub>i,c</sub> =  $\alpha_i + \mu_1 \times \Delta Payments_c + \mu_2 \mathbf{X}_{i,c} + \epsilon_{i,c}$ 

(Refinance (CRA)/Deposits) $_{i,c}$  or (Refinance (CRA)/Revenue) $_{i,c}$  are the total mortgage refinance (small business loans) scaled by total deposits (or total revenue) of bank i in county c during 2011 and 2013. Bank control variables include pre-shock revenue per employee of the bank in a county. Payments gap is the hypothetical amount of interest payments that could be saved due to the interest rate decrease, if local households chose to refinance their mortgages during the year of 2011 and 2013. Bank control variables include banks' revenue per employee of the bank in a county. County level control variables include the pre-shock unemployment growth, labor force participation rate, population growth rate, logarithmic of total number of establishments, share of non-tradable sector business establishments, real GDP per capita. Fixed effects include bank fixed effects. Standard errors are clustered at county level. \*\*\*, \*\*, and \* denote significance at the 1%, 5%, and 10% levels, respectively.

	Refinance/Deposits	Refinance/Revenue	CRA/Deposits	CRA/Revenue
	(1)	$\frac{}{(2)}$	(3)	$\phantom{aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa$
$\Delta$ Payment <sub>c</sub>	0.363***	0.389**	0.008	-0.244
	(0.105)	(0.185)	(0.137)	(0.210)
Revenue/Emp	$0.037^{***}$	-0.174***	0.028***	-0.197***
	(0.007)	(0.017)	(0.008)	(0.023)
Unemployment growth	0.600***	-0.960***	0.611***	-0.661***
	(0.099)	(0.200)	(0.121)	(0.237)
Population growth	$0.046^{***}$	-0.023*	$0.037^{***}$	-0.018
	(0.007)	(0.012)	(0.008)	(0.016)
% workers in small business	$0.964^{***}$	-1.122***	1.458***	-0.926***
	(0.065)	(0.087)	(0.078)	(0.101)
$ln(Total\ establishments)$	-0.172**	0.043	-0.127**	0.015
	(0.069)	(0.104)	(0.059)	(0.036)
Non-tradable share	0.114**	-0.491***	$0.242^{***}$	-0.252**
	(0.049)	(0.093)	(0.059)	(0.121)
GDP per capita	-0.320*	0.691**	-0.080	0.432
	(0.164)	(0.286)	(0.189)	(0.327)
Bank FE	Y	Y	Y	Y
Clustered	Y	Y	Y	Y
AdR-squared	0.538	0.033	0.568	0.062
N	16,953	15,663	20,197	18,405

# Table A17. Comparing Lending Club Magnitude with Commercial Banks' Personal Loan

This table presents the summary statistics of Lending Club's loan outstanding balance and loan origination compared to the personal loan balance and origination of banks during 2015-2019. "personal loan" is the aggregate of total personal loan on banks' balance sheets in Call Re-"Apersonal loan" is the change of aggregate total personal loans of banks from last year. Mathematically, Total personal loan  $t = \sum_b Personal loan balance_{b,t}$ , and  $\Delta Total personal loan_t = \Delta Total p$  $\sum_{b}$  Personal loan balance<sub>b,t</sub> -  $\sum_{b}$  Personal loan balance<sub>b,t-t</sub>. "LC Origination" is the total amount of newly loans Lending Club originated in a year, and "LC outstanding" is the total estimated dollar amount of Lending Club's loan outstanding. Mathematically, LC Origination<sub>t</sub> =  $\sum_{l}$  LC loan<sub>l,t</sub>. "LC outstanding" is estimated using the loan issuance information by aggregating dollar amount of loans issued in past years that have not matured. Mathematically, LC Outstanding  $=\sum_{m\in\{1,2,\dots,5\}}\sum_{k\in\{0,1,2,\dots,m-1\}} \text{LC origination}_{t-k}^m \frac{m-k}{m}$ , where m are the loan maturities of the loans offered by Lending Club in previous years, and LC origination m = m = 1stands for total loan amounts with maturity m that were issued k years ago, where k is less than m. To estimate the personal loans of the banking sector at the state level, we first estimate each banks' personal loan at the state level. Specifically, for each bank b, we calculate Total personal loan b s tTotal personal  $loan_{b,t} \times \frac{Deposits_{b,s,t}}{Deposits_b}$ . The state level personal loan balance of banking sector is calculated by summing all the banks' state-level balances for banks with operation in that state: Total personal  $loan_{s,t} = loan_{s,t} = lo$  $\sum_{b \in s}$  Total personal loan b,s,t. And flow of total personal loan by banking sector at state level is estimated by  $\Delta \text{Total personal loan}_{t,s} = \text{Total personal loan}_{t,s} - \text{Total personal loan}_{t-1,s}$ 

_	٠,٠	_	0 1,0					
LC origination/ $\Delta$ personal loan								
(5	state-leve	el)	(nation-level)					
25th	50th	75th	Mean					
1.18%	2.72%	9.37%	6.20%					
LC bal	lance/pe	rsonal loa	n outstanding					
(	state-leve	el)	(nation-level)					
25th	50th	75th	Mean					
1.74%	3.88%	12.36%	1.76%					

Table A18. Importance of Loan types in Profit Making for Banks across Size Groups

This table presents the interest income as a share of total interest income for different types of loans using banks' income statement information from Call Report. Banks are categorized into 5 groups according to their average asset size during 2010-2019.

	Credit card	Personal loan (all)	C& I	Agriculture	Mortgage
< \$100 Million	0.026	0.131	0.176	0.255	0.296
\$100 Million-\$1 Billion	0.024	0.087	0.171	0.122	0.315
\$1 Billion-\$10 Billion	0.061	0.095	0.190	0.028	0.274
\$10 Billion-\$250 Billion	0.211	0.222	0.223	0.008	0.284
> \$250 Billion	0.142	0.219	0.165	0.002	0.239

# B Data Appendix: Construction of Bank IT Spending

In this section, we provide detailed description of our data matching process and conduct some supplemental cross-validation checking analysis that serve as a quality examination of the bank IT spending data we leverage.

We start with an elaboration of our sample construction procedure, which involves the matching between the Harte Hanks data set with other data sets such as the Call report and the Summary of Deposits (SOD) data set.

### **B.1** Details of Sample Construction

#### **B.1.1** Elaboration of the Matching Procedure

Our overall matching procedure follows existing works conducting name-based matching in the banking sector (Lerner et al. (2021)). Below we elaborate the whole matching process step by step.

Step I: Initial construction of matched pairs. Our matching algorithm starts by dropping the sites with names that are *very* different from the official bank names. Similar to the old algorithm, we first take the bank names from SOD data set and obtain the banks' homepages from Google. The first step is to extract a smaller set of site names in the site-level bank IT data that are similar to the names of the banks in SOD. We drop the suffixes ", national association", "national association", "fsb", "fsb", "na.", "na.", "f.s.b.", "f. s. b.", ", f. s. b.", ", s.b.", ", s.b.", ", s.b.", ", s.b.", ", fsb", ", fsb", ", fsb", ", a fsb", and ", a federal savings and loan association", "bank", and "national bank", etc in the SOD data. This serves as the basic cleaning of the naming convention, this step is the same as that in the old matching algorithm.

Afterwards, we split the names into at most two keywords by spaces. For example, Wells Fargo Bank is labeled as "Wells" and "Fargo." This is because many site names in the bank IT data set, which are going to be merged with official bank names in SOD later, are written without spaces. In the Wells Fargo Bank case, the site names could be written as "wellsfargo bank" or "thewellsfargobank." Given that most sites in the IT data set also include a url variable that label the website address of the bank's homepage, we conduct the matching using url first. When matching with url does not work, we then conduct matching based on keywords in names constructed above. Specifically, we drop sites that do not contain the keywords. This serves as the first step of dropping sites that are *very* different from the official bank names in SOD.

For the matched pairs obtained after this initial step in our sample construction, we then calculate the Levenshtein distance of the string names between each site's name and all the names that showed up in SOD. Similar to our old algorithm, for each matched site name (in the bank IT data set) after the first step, we keep the bank name (in SOD) that is closest to the site name. This step guarantees that the same site in the bank IT dataset will not be matched with multiple

banks with distinct names, which thus reduces the false positive cases in our matching process and ensures that there is no double counting in the later calculation of the bank-level IT spending.

The matching procedure in the first step serve as the initial filtering of the bank IT data set, which targets at dropping the sites that are significantly different from any of the bank names in the SOD dataset. Yet, there might still exist many false positive cases after these two steps in the matching process. For instance, in Table B1, we list several concrete examples of false positive matches after the implementation of the first step in our matching procedure as described above.

Step II: Further data cleaning for accurate string matching. Before we drop the matched pairs by applying a threshold in Levenshtein distance of 0.1 (which is implemented in Step III), we conduct two further steps (filters) in cleaning the matched pairs obtained after Step I. These two additional filters serve to make the calculated Levenshtein distance more informative about the discrepancy between the names of establishments in Harte Hanks and banks in the Call Reports.

First filter (cleaning site names in Harte Hanks data). The first filter aims at cleaning names in Harte Hanks by manually identifying systematic differences in the names between the official names in Call Reports and the sites names (especially of non-financial ones) in Harte Hanks. For example, the names of Wells Fargo's sites that engage in non-deposit-related financial services would in general have suffixes like "services LLC." Therefore, we manually drop these common substrings that often exists in the names of non-financial sites to ensure that the string distances between the official names and the non-financial sites are inaccurately increased due to different suffix. In other words, this filter guarantees that the calculated Levenshtein distances accurately reflect the dissimilarities between the official names and the non-financial sites, without any systematic upward bias caused by these common substrings.

Specifically, in this first adjustment we clean site names by dropping strings that consistently show up in the site names but are not in the official bank name. For example, we drop strings like "financial inc.," "clearing servicess llc," "& company," "home mortgage inc.," "advisors llc," " usa," "rail corporation," "financial acceptance inc.," and so on. These cases tend to be more common in larger banks, which often have numerous associated non-financial establishments. Concrete cases of such instances can be seen in Figure B1, which showcases examples from Wells Fargo. For example, in Figure B1 "wells fargo clearing services llc" has a very large Levenshtein distance of 0.5 from "wells fargo". After dropping the substring " clearing services llc", the Levenshtein distance becomes zero.

This step of data cleaning aims to reduce the occurrence of mistakenly dropping matches due to an exaggerated Levenshtein distance. For example, in Figure B1, "wells fargo clearing services llc" has a very large Levenshtein distance of 0.5 from "wells fargo", despite the fact that both name strings are referring to the Wells Fargo bank. Furthermore, we expect this effect to be particularly pronounced for big banks, since the problem described above is more severe for big banks due to

their more complicated within-bank establishment structure. By implementing this adjustment, we aim to mitigate the impact of these problem due to name complexity and enhance the accuracy of our matching process.

Second filter (more stringent matching). The second filter aims to reduce false positive rates in our matching process. The false positive cases are mainly due to the fact that many small but different banks could have very similar names. The errors in this stage are more miscellaneous. For example, People's Bank of Alabama will have a small Levenshtein distance from People's Bank. Despite the similarity in names, these two banks are distinct entities.

To tackle this specific issue, we have implemented a manual filtering process to identify and exclude bank-site pairs that may result in such mismatches. For instance, if the site name (in Harte Hanks) does not include any state names while the bank name (in Call Report) does, we discard the matching. In another example, our previous procedure mistakenly matched Citizens Trust Bank from Harte Hanks with Citizen Bank from Call Report, or West Bank from Harte Hanks with Bank of the West from Call Report. We manually drop such false positive cases. This step is expected to be particularly effective for dropping false-positive matches for smaller-sized banks, since this common name issue is more prevelant among small banks.

In the end, the Harte Hanks dataset often uses "&" for abbreviation of "and." We therefore change the symbol to texts and drop all spaces in the names from both datasets, and recalculate the Levenshtein distance. A summary of the manually adjusted algorithm:

- require title and "bank" to have no other words in between, and require the keyword in the title right before "bank" to be exactly the same
- require "XX Bank" to be different as "Bank of the XX"
- require entity names with and without ending location name to be different
- disallow abbreviation to be matched

Step III: Dropping matched pairs based on Levenshtein distance. As the final step of the sample construction procedure, we drop matched site-bank pairs with Levenshtein distances higher than 10%, where the calculation of the Levenshtein have incorporated the two adjustments as described above.

#### **B.1.2** Evaluating Matching Quality

In this section, we conduct a quality examination of our matching algorithm by comparing our matched sample to other alternative sources on some main bank-level variables. In particular, these variables include the number of branches, total revenue, and total number of employees at the bank level.

<sup>&</sup>lt;sup>1</sup>Table B1 provides more concrete examples for such cases.

Comparison with SOD: Number of Sites/Branches. To examine the coverage performance (at the extensive margin) of our matched sample, for each matched bank we examine the difference between the number of financial establishments implied by our matched sample and the number of branches from SOD. This exercise is essential to our later analysis regarding the impact of different types of credit demand shock on local banks' IT spending. Since our analysis is conducted at the bank-county level that involves the aggregation of a banks' sites operating in a specific county, it is important to ensure that our matched sample has a relatively stable (and ideally full) site coverage ratio for all banks.

Note that to make this comparison meaningful, we restrict our matched sample to financial sites whose SIC codes start with 6. Figure B2 plots the distribution of the log difference between the number of financial sites and the number of branches in SOD. The four panels are the distributions by four revenue quartiles. As shown in Figure B2, the number of sites for matched banks is very close to the number of branches for the corresponding bank from SOD. Over the whole matched sample, the average difference in terms of site coverage is -2.7% and the median is zero. The IQR is -16.0% and 8.0%.

Furthermore, a bank-level regression of number of bank branches from SOD on the number of financial sites from matched sample with Harte Hanks, shown in Table B2, also suggests a decent sample coverage rate for our matched sample. Specifically, the regression coefficient and intercept are 0.880 and 0.055 respectively.

Comparison with Call Report: Bank-level revenue and employment. The comparison with SOD on matched banks' number of financial sites discussed above confirms that our matched sample has high qualities in terms of coverage and matching rates for financial sites of banks. In this section, we compare the bank-level total number of employees and total revenues aggregated from matched sites in the Harte Hanks dataset, to those from Call Report. Note that in this comparison we include all matched sites (rather than only restricting to financial sites as in the above comparison regarding number of sites/branches) in the calculation of the bank-level total employee and total revenue, since these aggregate bank-level variables from Call Report also account for the employment and revenue of non-financial sites. This comparison thus allows us to evaluate the coverage performance of our matched sample at the bank-level across all sectors. Furthermore, as the Harte Hanks data is collected from surveys, cross-checking with Call Report on important bank-level variables such as total revenue and number of employees also helps us to evaluate the quality of this survey-based data set.

The left panels of Figure B4 display the binned-scatter plots comparing log number of employees. As shown in panel A, the aggregate number of employees across all matched sites following our matching algorithm exhibit a high proximity to the bank-level counterparts from Call report. In a more formal comparison, columns (3) of Table B2 display the OLS regression estimates of this correlation. For our matched sample, the slope coefficient is 0.907 and the intercept is insignificant

from zero.

The right panel of Figure B4 plot the *total revenue* from Harte Hanks and *total income* from Call Report, which shows a decent correlations between the two variables. In columns (2) of Table B2, we also find regression estimates confirming this graphical alignment, reflected by a slope coefficient close to 1 and an insignificant intercept. It is worthy noting that there is a difference between revenue and income in accounting terms. But Call Report does not have revenue and only has income, which could be a reason why the correlation is lower for revenue (R-squared 0.71) than employees (R-squared 0.83).

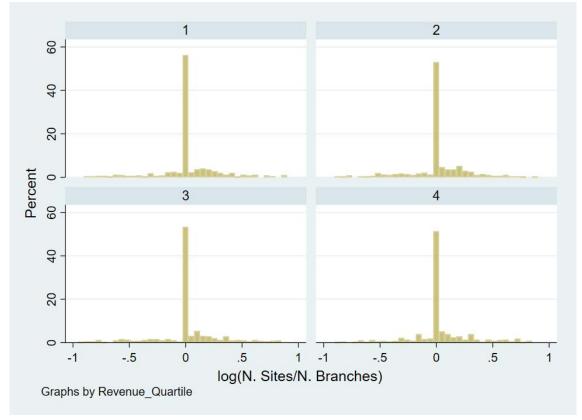
# Figure B1. Examples of Non-Financial Sites of Wells Fargo

This figure presents some examples that some sites of banks are not labeled as in the financial sector (SIC code starting with 6) for the Wells Fargo bank.

company	sic3_code	sic3_desc
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo clearing services llc	874	Management & public relations
wells fargo financial inc.	874	Management & public relations
wells fargo financial retail service inc	738	Miscellaneous business services
wells fargo equipment finance inc.	738	Miscellaneous business services
wells fargo financial security services inc.	738	Miscellaneous business services
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo & company	738	Miscellaneous business services
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo clearing services llc	738	Miscellaneous business services
wells fargo home mortgage inc.	874	Management & public relations
wells fargo & company	823	Libraries
wells fargo home mortgage inc	737	Computer & data processing services
wells fargo home mortgage inc	874	Management & public relations
wells fargo clearing services llc	874	Management & public relations
wells fargo home mortgage inc.	874	Management & public relations
wells fargo clearing services 11c	738	Miscellaneous business services

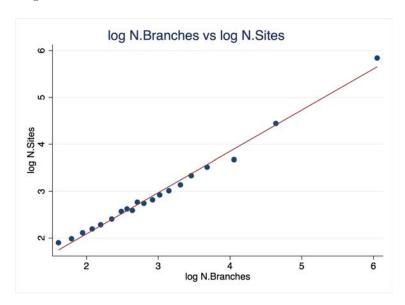
Figure B2. Comparison with SOD: Number of Financial Sites in Matched Sample and Number of Branches in SOD

This figure compares the number of sites of banks in the sample compared to the number of branches in Summary of Deposits (SOD). Each subplot presents the histogram of the log ratio between the total number of sites belonging to financial sector (SIC code starting with "6") in our sample and the total number of branches in SOD. Subplots are grouped by bank revenue quartiles.



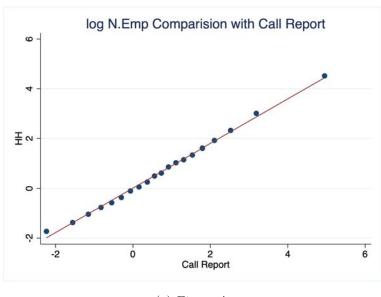
# Figure B3. Comparison with SOD: Number of Financial Sites and Branches

This figure gives the binned-scatter plots between the total number of sites belonging to financial sector of the matched banks in Harte Hanks, and the total number of branches in SOD. Financial sector sites are defined as sites with SIC sector code starting with "6." The plot is characterized by dividing the sample into 20 bins of the natural log of total number of branches of banks in SOD.

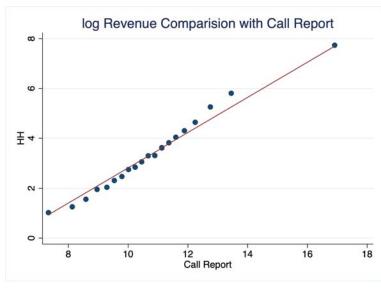


# Figure B4. Comparison with Call Report: Number of employees and total revenue

This figure compares total revenue and total number of employees of banks in matched sample in Harte Hanks dataset and Call Report. In particular, the figures examine the association between log number of employees and log total revenue in the matched sample and those in Call Report. Panels A plots the log total number of employees from Harte Hanks and that from Call Report. Panels B plots the log total revenue from Harte Hanks and log total income from Call Report.



(a) Figure A



(b) Figure B

# Table B1. Examples of False Positive Matches

This table gives some examples of the name discrepancy between Call Reports and Harte Hanks after Step I in the matching procedure.

(1)	(2)
Harte Hanks	Call Report
citizens trust bank	citizens bank
people's bank of alabama	people's bank sb
signature bank of arkansas	signature bank
central state bank of calera	central state bank
first national bank of jasper	first national bank
community bank and trust alabama	community bank and trust
the citizens bank corporation	the citizens bank
first arkansas bank & trust	first arkansas bank and trust
merchants & planters bank	merchants and planters bank
cb&s bank	the citizens bank
mnb Bank	the malvern national bank
traders and farmers bank	farmer's bank
integrity first bank	first bank & trust company
west bank	bank of the west

### Table B2. Comparing other variables of Harte Hanks with Call Report

This table shows the comparison of other bank-level variables between the matched sample in Harte Hanks and those in Call Report. "log N.Sites" is the natural log of total number of financial-sector sites. "log Reven (HH)" is the natural log of total revenue of the matched banks in Harte Hanks. "log Emp (HH)" is the natural log of total number of employees of the matched banks in Harte Hanks. "log N.Branches (SOD)" is the natural log of total number of branches of the matched banks in SOD. "log Income (CR)" and "log Emp (CR)" are the natural log of total income and natural log of total number of employees of the matched banks in Call Report.

	log N.Branches (SOD) (1)	log Income (CR) (2)	log Emp (CR) (3)
log N.Sites	0.880*** (0.009)		
log Reven (HH)		$0.753*** \\ (0.005)$	
log Emp (HH)			$0.907*** \\ (0.005)$
cons	0.055*** (0.011)	-0.191 $(0.136)$	$0.005 \\ (0.007)$

### B.2 IT Data Quality Checking: Comparison with Y-9C

#### B.2.1 Y-9C Data on Bank IT Investment

To evaluate the quality of the Harte Hanks bank IT data and its reliability as a measure of banks' IT investment, we conduct a thorough comparison with the relevant information of banks' investment in information technology as reported in "FR Y-9C Consolidated Financial Statements for Holding Companies."

The Y-9C filings include detailed balance sheet and income data, as well as information about loan performance, derivatives, off-balance-sheet activities, and other aspects of bank holding companies (BHC) operations. Data are reported on a consolidated basis, incorporating both bank and non-bank subsidiaries controlled by the BHC.<sup>2</sup> Per regulatory requirements, bank Holding Companies with consolidated total asset sizes above certain threshold are required to report. The reporting threshold went through two adjustment, once in 2015, when the reporting threshold was adjusted from \$500 million to \$1 billion, and in 2018, the reporting threshold was adjusted from \$1 billion. We utilize the 2018 threshold and take the sub-sample of banks with total assets \$1 billion in the Harte Hanks database and compare with BHCs in Y-9C with total assets \$1 billion or more.

In the consolidated income statement of Y-9C, BHCs report major categories of "Total noninterest expenses (BHCK4093)" including "Salaries and employee benefits (BHCK4135)," "Expenses of premises and fixed assets (BHCK4217)," "Goodwill and impairment losses (BHCKC216)," "Amortization and impairment losses for other intangible assets (BHCKC332)," and "Other noninterest expense (BHCK4092)." The item "Other noninterest expense (BHCK4092)" is further expanded into several sub-categories of expenses, such as advertising and marketing, postage, legal fees, etc. Following Kovner et al. (2014), who define term IT expenses of BHC as "IT and Data processing expenses", we define the total IT of a BHC in Y9C as the sum of "Data processing expenses (BHCKC017)," "Telecommunication (BHCKF559)," and other write-in expenses.<sup>3</sup> For the rest of this letter, we will label the sum of "BHCKC017," "BHCKF559," and text items "text8565"-"text8567" as the "Total IT spending in Y-9C."

Figure B5 displays a snapshot of the summary statistics in Kovner et al. (2014) of the IT expenses using the above definition from 2008-2012. One issue worthy noting is that the IT expense reported by Kovner et al. (2014) used "other noninterest expenses (BHCK4092)" as denominator, which are much smaller than "total noninterest expenses (BHCK4093)." One critique from R7 of our data's underestimation of IT is based on "IT/other noninterest expense," which is mechanically greater than "IT/total noninterest expenses." To make consistent and fair comparisons, in the latter part of this section, we directly compare the *dollar amount* of IT spending in the new version of our matched sample with that in Y-9C. In constructing measures for banks' scaled IT spending

<sup>&</sup>lt;sup>2</sup>Avraham et al. (2012) provides more detailed introduction on this.

<sup>&</sup>lt;sup>3</sup>These write-in expenses include "information technology," "software," and "internet banking" filled in text items "text8565," "text8566," and "text 8567."

intensity, we plot time-varying evolution consistently using "IT/total noninterest expenses."

Before proceeding to the comparison, we would like to highlight the features of BHCs and the reporting guideline regarding Y-9C, which will determine the empirical measures and methodology we employ to compare Harte Hanks data with Y-9C in proper manners.

#### Issue 1: Y-9C data is consolidated at the bank holding company (BHC) level

Bank holding companies, especially the larger ones, often have complicated subsidiaries, which are not necessarily deposit taking or lending entities. For instance, as pointed out by Avraham et al. (2012), big bank holding companies in U.S. have subsidiaries organized as "funds, trusts and other financial vehicles," "insurance carriers," "management of companies," "professional, scientific and technical services," etc. Total number of subsidiaries of large bank holding companies such as JP Morgan was more than 3000 in 2012.

Importantly for our later comparison, many of these subsidiaries often have substantially different names compared to its major commercial banking subsidiary. For instance, the Bank of America BHC has subsidiaries sharing similar names of "Bank of America," such as "Banc of America Development, Inc.", but it also has many subsidiaries whose names are drastically different from "Bank of America", for example, "Boston Advisors, Inc.", "BANA GA Mortgage Company", "First Capital Corporation of Boston," etc. The richness of industries coverage of subsidiaries of BHC's lead to the complicated name structures of these subsidiaries. The larger the bank size, the more complicated the subsidiaries' names. For smaller banks, the number of subsidiaries is relatively small and the name differences also tend to be limited. For example, Compass Bank has 66 subsidiaries compared to Bank of America, which had more than two thousand subsidiaries.<sup>4</sup>

Since Y-9C is a consolidated report of BHCs encompassing all the associated subsidiaries, we will need to take into consideration of these subsidiaries in the Harte Hanks database in order to properly compare with the IT expenses as reported in Y-9C. Furthermore, it is worthy noting that BHCs frequently operate overseas subsidiaries. For example, "Asian American Merchant Bank Ltd." is a subsidiary of Bank of America based in Singapore. These overseas subsidiaries' financial statements are also consolidated into the Y-9C filings. Since the Harte Hanks database only contains establishments in the U.S., we are aware of this difference and acknowledge that Harte Hanks data is not able to fully capture the IT spending for those banks with substantial amount of assets held by foreign subsidiaries. For example, as of 2012, Citigroup Incorporate had only 935 out of 1645 subsidiaries that are operating in U.S., and assets held by its domestic subsidiaries were 68.8% among its total asset. Augmenting all of the above features regarding the features of BHCs, we supplement the manual inclusion of subsidiaries to our existing algorithm-based matching using the names of all the subsidiaries disclosed by BHC's SEC reports. Without incorporating the IT spending of these subsidiaries, the IT expenses will be underestimated compared to that in Y-9C, which is consolidated at the BHC level.

<sup>&</sup>lt;sup>4</sup>The full list of Bank of America's subsidiaries as of 2009 can be found here. The full list of Compass Bank's subsidiaries can be found here.

#### Issue 2: reporting in Y-9C data is censored

Another notable feature of the Y-9C data is related to the reporting guidelines for the items in the other noninterest expense. Banks whose data processing expenditure not exceeding 7% of "other noninterest expenses" are not mandated to report, and when no reporting is made by a bank in a certain year, in the data set the entry will appear as a "0" or a missing value. The paragraph below is extracted from the reporting guidelines. The figure B6 shows a screenshot of the Y9C report of JP Morgan Chase in 2019 that reflects this reporting rule. As is shown, the entry for the item BHCKC017 "Data Processing Expenses" of JP Morgan Chase the September report of 2019 was zero. The paragraph below is extracted from the Y-9C reporting guidelines:

"Report all operating expenses of the holding company for the calendar year-to-date not required to be reported elsewhere in Schedule HI. Disclose in Schedule HI, Memoranda items 7(a) through 7(p), each component of other noninterest expense, and the dollar amount of such component, that is greater than \$100,000 and exceeds 7 percent of the other noninterest expense reported in this item. If net gains have been reported in this item for a component of "Other noninterest expense," use the absolute value of such net gains to determine whether the amount of the net gains is greater than \$100,000 and exceeds 7 percent of "Other noninterest expense" and should be reported in Schedule HI, Memoranda item 7."

Figure B7 shows a concrete example of Y-9C report made by First ST BK KIOWA KS. The orange line shows the IT spending scaled by "other noninterest expenses" (ONIE), and the blue line shows the dollar amount of spending in thousands of dollars. It is apparent that for quarters during which the bank reported "0," the positive spending it reported in adjacent periods was also very close to at least one of the threshold for mandated reporting. This pattern further indicates that the "0" entries can be interpreted as missing entries. Among all the Y9C quarterly reports of IT spending after 2010 made by BHC's with greater than \$ 1 billion asset size, 70.99% of the observations have positive reporting on the IT spending, 4.91% are missing entries and 24.10% are zero entries.

Given this reporting rule, the missing observations or "0" entries in Y-9C do not provide informative measures for us to make comparison, as zero or missing entrance does not necessarily mean the BHC spent zero in a given quarter, and it will also introduce estimation error if we simply impute an estimated spending to those missing entries. To get around this reporting issue and make informative and proper comparison with Y-9C, we will focus on the year-end quarter (4-th quarter) reporting of BHC's for which both BHCKC017 "Data Processing Expenses" and BHCKF559 "Telecommunication" are non-missing. Since the BHCKC017 "Data Processing Expenses" (more than 33,000 bank-quarter frequencies) and BHCKF559 "Telecommunication" are reported of much more higher frequencies (around 24%) than the write-in text items (around 2.1%), and all of these items are subject to the reporting criteria, we focus on comparing the summation of the IT-related variables when both BHCKC017 "Data Processing Expenses" and BHCKF559 "Telecommunication" are reported.

#### B.2.2 Sample Adjustment for Comparison with Y-9C

Based on the characteristics of BHCs and the reporting rules in Y-9C filings, we apply two major adjustments to our constructed sample (following procedure as described in Appendix B.1.1) in order to make proper comparisons with Y-9C.

In the first part of the constructions, we extract the BHCs in Y9C with total consolidated assets of more than \$1 billion dollars, and for each of the BHCs satisfying the asset threshold, we extract the reporting of item BHCKC017 "Data processing expenses" and BHCKF559 "Telecommunication" of that BHC during 2010 and 2019 requiring that there are no missing observations or zero values for both of these two items for the fourth quarter within a year. This part of construction give us the sub-sample of BHCs in Y-9C filings with reliable and informative report on their IT spending for us to make comparison with.

In the second part of the construction, we extend our existing algorithm described in the previous section by manually incorporating the subsidiaries of BHCs. Since the number of subsidiaries and the degree of complexity of the subsidiaries' names tend to increase with the BHCs' size, we first sort the BHCs by their asset sizes. We then manually input subsidiaries' names and supplement to the algorithm-based matching to get the estimate of BHCs' IT spending in the Harte Hanks database. As BHCs' sizes get smaller, the total number of subsidiaries gets smaller and the number of subsidiaries with substantially different names from the BHCs' also gets smaller. Therefore, as the BHCs' size gets smaller, the manual adjustments of subsidiaries should make fewer effective supplements to the IT spending of BHC's implied by the Harte Hanks data based on "algorithm-only" matching. The manual supplementation of subsidiaries to the existing algorithm (refer to as "matching iteration adjusting for BHC subsidiaries" in later analysis) is formally described as follows:

- 1. Step 1: Take the sub-sample of BHCs in Y-9C such that at each bank-year, all relevant IT expenditures are reported.
- 2. Step 2: Sort the sub-sample of the BHCs in step 1 by their average asset size between 2010-2019.
- 3. Step 3: Starting from the largest BHCs in the Y-9C, for each BHC, manually collect the subsidiaries' information from that BHC's SEC disclosure, and incorporate the subsidiaries' names in the current algorithm to calculate the total IT spending for that BHC.
- 4. Step 4: Assess the matching in the updated algorithm by comparing the difference in IT spending under old algorithm and the newly manually supplemented BHC-version algorithm adjusted for subsidiaries.
- 5. Step 5: Proceed to the next BHC on the list sorted by bank asset size until the assessment in step 4 is small.

6. Step 6: For the remaining BHCs treat the IT spending in Harte Hanks as the including all the subsidiaries.

Table B3 illustrates the iterations of the algorithm comparisons as described above. For each BHC, we label the average dollar amount of IT spending reported by item "BHCKC017," "BHCKF559," and other write-in items in "text8565"-"text8567" in Y-9C filings as "Y9C", which is reported in Column (2). We label the average dollar amount of IT spending under the algorithm supplemented by manual inclusion of subsidiaries as "Manual adjustment + Algorithm," which is reported in Column (3), and we label the average dollar amount of IT spending under only the algorithm as "Algorithm-only", which is reported in Column (4). We further compare the difference between the IT expenses reported in Y-9C and IT spending in Harte Hanks based on "algorithm-only," which is reported in column (5), as well as the difference between the IT expenses reported in Y9C and the IT spending in Harte Hanks based on "manual adjustment + algorithm" and "algorithm-only," which is reported in Column (6). Specifically, in Column (5), we calculate  $|\log(\mathrm{IT}_{\mathrm{Harte Hanks}}^{\mathrm{algorithm-only}}|$ , and in Column (6), we calculate  $|\log(\mathrm{IT}_{\mathrm{Harte Hanks}}^{\mathrm{algorithm-only}}|$ .

Given that BHCs can have complicated subsidiaries, especially the larger BHCs, and the fact that Y-9C reporting captures the consolidated balance sheet information and income statement, there should be a nontrivial gap between the IT expense in Harte Hanks based on "Algorithm-only" and the Y-9C. This is because the "Algorithm-only" approach can only capture the IT spending of commercial banking offices of the BHC with the same key string names. However, if the IT expenses provided by Harte Hanks are comprehensive and relatively accurate and our algorithm performs well, as the BHCs sizes become smaller, the difference between Y-9C and "Algorithm-only" should get smaller. Furthermore, if IT expenses estimation provided by Harte Hanks is reliable, then even for BHCs with a nontrivial number of subsidiaries, the "matching iteration adjusting for BHC subsidiaries" approach should allow us to shrink the gap between IT expenses provided by Y-9C and that implied by Harte Hanks to an insignificant level.

As is demonstrated by Table B3, for relatively larger BHCs, there is indeed a non-trivial discrepancy between IT expenses reported in Y-9C and IT spending constructed based on Harte Hanks under the "Algorithm-only" scenario, as is summarized in Column (6). This is especially true for the largest BHCs. Meanwhile, as is shown in Column (3), manual incorporation of subsidiaries' IT spending in Harte Hanks significantly shrink the discrepancy between original IT expenses reported in Y-9C and the "Algorithm-only" scenario of Harte Hanks estimation. Accordingly, the difference between "Manual adjustment + Algorithm-only" and "Algorithm-only", which is summarized in Column (5), tends to have very similar magnitude with Column (6). This means that manual inclusion of subsidiaries' IT spending can effectively match the IT expenses in Harte Hanks with IT expenses reported by BHC's in Y-9C, which indicates that the IT expenses records provided by Harte Hanks can match Y-9C pretty well. Furthermore, as we iterate the algorithm from larger BHCs to smaller BHCs, we find that the discrepancy between "Y9C" and "Algorithm-only" tends to get smaller, which is shown in Column (5). This pattern confirms that as the number of sub-

sidiaries of a BHC gets smaller, the "Algorithm-only" scenario, which is based on matching string name of the major commercial banking office of a BHC alone, can yield a matching that provides a Harte Hanks-version IT spending measure that is decently aligned with the level reported in Y-9C filings. As discussed above, these calculations help validate the reliability of our updated matching algorithm and further suggest that IT spending provided by Harte Hanks can match IT expenses reported by Y9C decently well.

#### **B.2.3** Comparison Outcomes

In this subsection, we conduct the comparison between the IT spending data in Harte Hanks and that in Y-9C at the BHC-year level. As is described above, since larger BHCs tend to have complicated subsidiaries, which necessitates manual incorporation of subsidiaries' IT spending into "Algorithm-only" version of Harte Hanks, we developed an iteration algorithm for the manual incorporation for the larger BHC's. For smaller BHC's, since the subsidiaries are much fewer, we apply a stopping rule of the iteration at the 36-th BHC sorted by BHC total asset sizes and directly compare the Y-9C IT expense report with the "Algorithm-only" version of Harte Hanks. Our comparison covers all BHC's with total asset size (averaged during 2010-2019) over 1 billions in years at which Y-9C has valid measure for us to make the comparison. We exhibit the comparison outcomes in the following two parts:

- 1. Part I: For each of the top 50 BHCs sorted by average total assets between 2010 and 2019, we plot the dollar amount of IT expenses of BHCs reported in Y-9C filings and the IT spending aggregated based on the matched sample adjusted for BHC subsidiaries from Harte Hanks. For each of these BHCs, the IT spending is plotted at each year that allows us to make meaningful comparison (i.e., the reporting in Y-9C is not missing or zero).
- 2. Part II: On the entire matched sample constructed for this comparison, we regress the IT spending amount reported in Y-9C on the IT spending aggregated from Harte Hanks based on the match sample that has adjusted for the BHC subsidiaries as described above at the BHC-year level.

Figure B9 displays the dollar amount comparison between IT expenses reported in Y-9C and the IT spending implied by Harte Hanks based on "matching iteration adjusting for BHC subsidiaries" for the top 50 BHCs. The blue dashed lines show the the dollar amount of the consolidated IT spending of BHCs in the Harte Hanks data and the black dashed lines plot the dollar amount of IT spending reported in Y-9C. The unit of observations are in hundreds of millions of dollars (or 10,000,000 dollars). As is shown by the figures, for majority of these large BHC's, after manually adjusted for their subsidiaries the Harte Hanks data yields very similar magnitude of dollar amount of IT spending and over-time trend of IT spending to those reported in Y-9C.

It is worthy noting that among the larger banks, the Harte Hanks IT spending measure of Citigroup exhibits a relatively significant underestimation when compared to Y-9C. One potential explanation for this disparity is that Citigroup holds over 30% of its assets in overseas markets, i.e., locations outside of the U.S. Consequently, the Harte Hanks data, which provides information solely for IT budgets of U.S.-based sites, fails to capture the spending of Citigroup's overseas offices and subsidiaries. The importance of Citigroup's overseas subsidiaries is briefly discussed in Section B.2.1.

We then conduct this comparison on the entire constructed sample for BHCs of asset size over 1 billions and on years valid comparisons could be made. Figure B10 shows the over-time evolution of 10-th/50-th/90-th percentile comparison of IT expenses in Y-9C and IT expenses in Harte Hanks dataset. In both figures, IT expenses are scaled by total noninterest expenses. The demonstration reveals a strong correlation between the IT expenses reported at the BHC level in Harte Hanks and those in Y-9C. This correlation is evident not only in the bank-year evolution but also in the over-time evolution of the interquartile range. Finally, Figure B11 shows binned scatter plot of IT expenses in Y-9C and IT expenses in the Harte Hanks dataset.

In Table B4, we conduct the regression analysis highlighted in Part II. In particular, we regress the logarithmic of total IT spending in Y-9C on the logarithmic of total IT spending in the Harte Hanks dataset constructed based on "matching iteration adjusting for BHC subsidiaries," with the following regression specification:

$$ln(\text{Total IT}_{i,t}^{\text{Y-9C}}) = \alpha + \beta ln(\text{Total IT}_{i,t}^{\text{Harte Hanks}}) + \epsilon_{i,t}.$$

As is demonstrated in the table, the coefficient of the above estimation is close to 1 (slope estimate 0.935), and the intercept is close to 0 (intercept estimate 0.037). This estimates thus suggests that after properly incorporate the IT expenses of associated subsidiaries, one can use the Harte Hanks data set to construct a panel of bank's IT spending of each year that consistently matches with those reported in Y-9C.

# Figure B5. Snapshot of summary statistics from Kovner et al. (2014)

This figure shows the screenshot of the Panel B of Table 2 of Kovner et al. (2014), where the authors provide summary statistics on the components of other noninterest expenses reported in Y-9C.

Panel B: Components of Other Noninterest Expense, as a Percentage of Total Other Noninterest Expense: 2008-12

		Individual Observations								
Component (Author-Defined)	Industry	p0.5	p5	p25	p50	p75	p95	p99.5	Mean	Standard Deviation
Corporate overhead	18.63	0.00	2.43	10.29	16.26	22.70	34.58	50.95	17.07	10.07
Information technology and data processing	12.63	0.00	0.64	8.21	13.84	19.81	29.91	45.01	14.54	8.69
Consulting and advisory	11.07	0.00	0.00	0.00	2.31	5.78	12.73	29.97	3.74	5.23
Legal	6.68	0.00	0.00	0.00	3.53	6.19	12.43	24.71	4.16	4.71
Retail banking	6.35	0.00	0.00	0.00	6.41	13.48	29.64	55.24	9.24	10.55
FDIC assessments and other government	5.81	0.00	0.00	6.80	11.53	16.95	25.54	37.34	12.26	7.58
Other financial services	3.01	0.00	0.00	0.00	0.00	0.00	4.00	15.85	0.56	2.72
Directors' fees and other compensation	0.25	0.00	0.00	0.00	0.00	3.45	6.99	14.60	1.91	2.85
Miscellaneous	1.76	0.00	0.00	0.00	0.00	0.00	5.75	24.91	0.84	3.98
Total classified	66.20	4.02	35.11	55.83	66.87	75.05	85.72	95.35	64.32	15.73
Unclassified	33.80									

Source: Board of Governors of the Federal Reserve System, Consolidated Financial Statements of Bank Holding Companies (FR Y-9C data).

# Figure B6. Reporting issue with Y-9C

This figure shows a screenshot of JP Morgan's Y-9C report, where the two major items related to IT spending as classified by Kovner et al. (2014), "BHCKC017" ("Data processing expenses") and "BHCKF559" ("Telecommunication expenses"), are missing.

Memo Items 7.a through 7.p are to be completed annually on a calendar year-to-date basis in the December report only by holding companies with less than \$5 billion in total assets. Holding companies with \$5 billion or more in total assets should report these items on a quarterly basis.			
7. Other noninterest expense (from Schedule HI, item 7.d, above) (only report amounts greater			
than \$100,000 that exceed 7 percent of the sum of Schedule HI, item 7.d):			
a. Data processing expenses	C017	0	M.7a.
b. Advertising and marketing expenses		3579000	M.7.b.
c. Directors' fees	4136	0	M.7.c.
d. Printing, stationery, and supplies	C018	0	M.7.d.
e. Postage	8403	0	M.7.e.
f. Legal fees and expenses	4141	0	M.7.f.
g. FDIC deposit insurance assessments	4146		M.7.g.
h. Accounting and auditing expenses.	F556	0	M.7.h.
i. Consulting and advisory expenses	F557	4244000	M.7. i.
j. Automated teller machine (ATM) and interchange expenses		0	M.7. j.
k. Telecommunications expenses	F559	0	M.7.k.
I. Other real estate owned expenses	Y923	0	M.7. I.

Asset-size test is based on the total assets reported as of prior year June 30 report date.

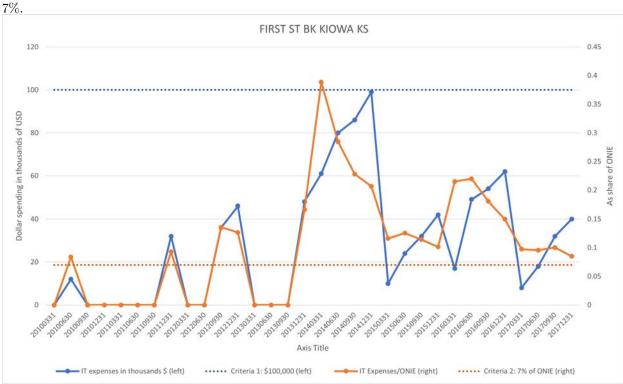
12/2019

Last Update: 20210719.160120 RSSD ID: 1039502

FR Y-9C Page 5 of 72

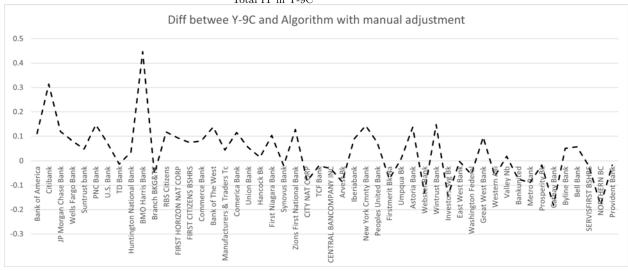
### Figure B7. Reporting issue with Y-9C

This figure demonstrates the reporting issue of missing value frequently seen in Y-9C report. The blue solid line (left axis) shows the reported IT expenses in thousands dollars by "FIRST ST BK KIOWA KS," and the orange solid line (right axis) shows the reported IT spending as a share of total ONIE (total other noninterest expenses). The blue dashed line (left-axis) shows the reporting criteria in dollar terms, which \$100,000; the orange dashed line (right axis) shows the reporting criteria of spending as a share of total ONIE, which is



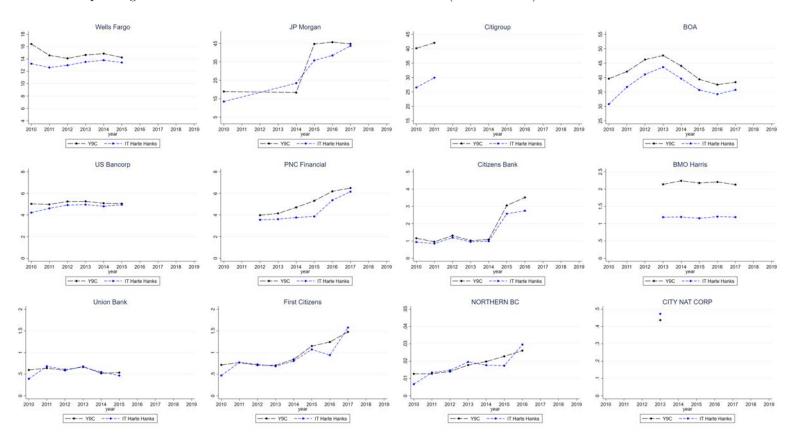
# Figure B8. Difference between algorithm and manual matching shrink as bank sizes gets smaller

The figures below show the comparison of IT spending in the matched Harte Hanks Database and the IT spending in Y9C. The black dashed line shows the difference of total IT spending in Y-9C and total IT spending in Harte Hanks with algorithm combined with manual adjustment of subsidiaries as a share of total IT spending in Y-9C. Mathematically, the black dashed line shows  $\frac{\text{Total IT in Y-9C-Total IT in Harte Hanks}}{\text{Total IT in Y-9C}}$ 



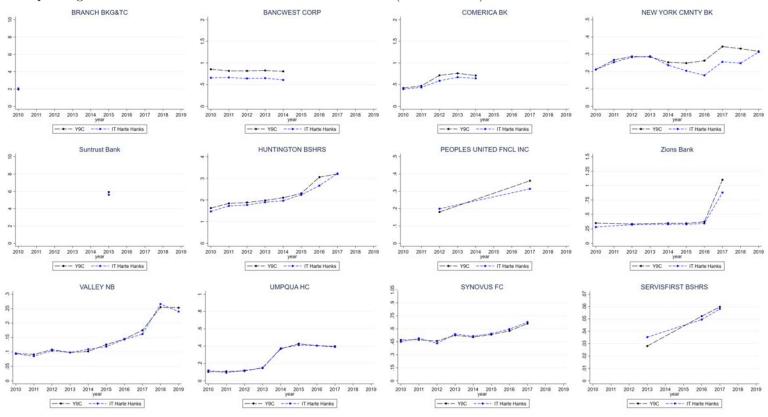
## Figure B9. Comparison with Y9C: Top 50 BHCs

This figure illustrates the evolution of total IT to non-interest expense for the top 50 banks conditional on the availability of data. The sample is based on banks-year observations with non-missing and non-zero IT spending in Y9C. The units are in hundreds of million dollars (or  $10^8$  dollars).



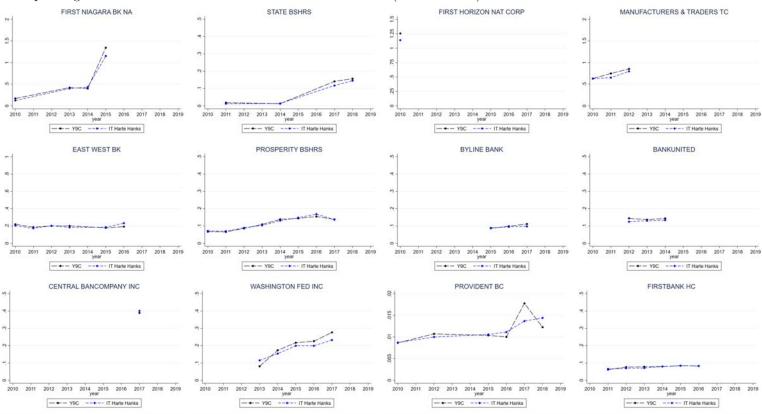
## Figure B9. Comparison with Y9C: Top 50 BHCs (cont'd)

This figure illustrates the evolution of total IT to non-interest expense for the top 50 banks conditional on the availability of data. The sample is based on banks-year observations with non-missing and non-zero IT spending in Y9C. The units are in hundreds of million dollars (or 10<sup>8</sup> dollars).



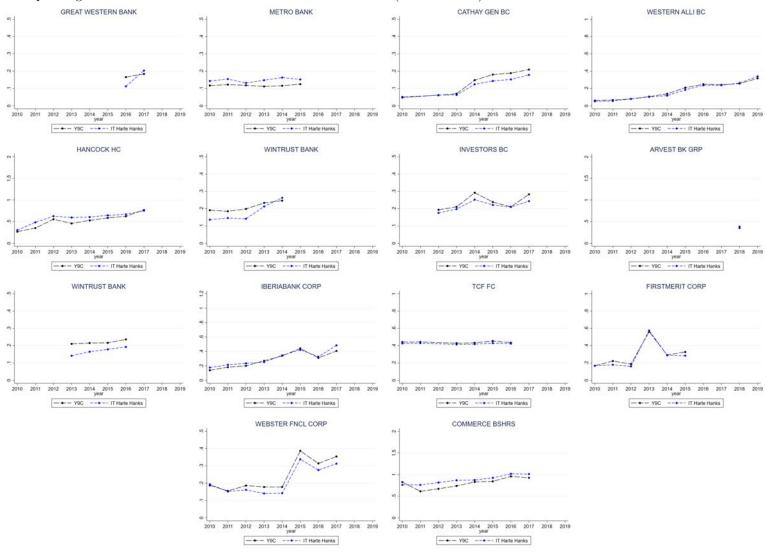
## Figure B9. Comparison with Y9C: Top 50 BHCs (cont'd)

This figure illustrates the evolution of total IT to non-interest expense for the top 50 banks conditional on the availability of data. The sample is based on banks-year observations with non-missing and non-zero IT spending in Y9C. The units are in hundreds of million dollars (or  $10^8$  dollars).



## Figure B9. Comparison with Y9C: Top 50 BHCs (cont'd)

This figure illustrates the evolution of total IT to non-interest expense for the top 50 banks conditional on the availability of data. The sample is based on banks-year observations with non-missing and non-zero IT spending in Y9C. The units are in hundreds of million dollars (or 10<sup>8</sup> dollars).



## Figure B10. Comparison with Y-9C: Distribution of IT spending overtime

This figure plots the evolution of the distribution of BHC-level IT spending as a share of non-interest expense between IT spending at BHC level constructed from Harte Hanks and that from Y9C. The figure plots the 10th percentile, median, and 90th percentile of IT spending over non-interest expense.



## Figure B11. Comparison with Y-9C: Overall sample correlation

This figure shows the binned-scatter plot of log total IT spending of our measure (y-axis) vs that from Y9C (x-axis) at the bank-year level. The natural log of IT spending reported in Y-9C as shown in x-axis is divided into 20 bins.

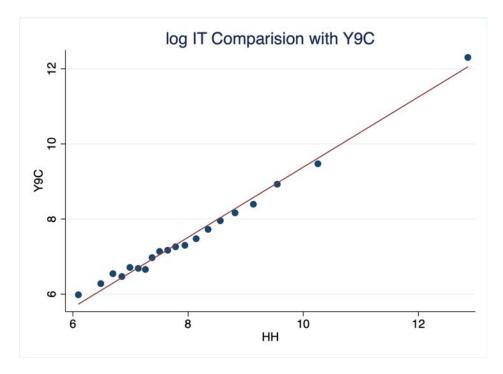


Table B3. Comparison with Y-9C: Manually Adjusted for Subsidiaries v.s. Algorithm-only Matching

The table below shows the dollar amount of IT spending in Harte Hanks under "algorithm-only," "algorithm plus manual adjustment for subsidiaries names" and the dollar amount of IT spending in Y9C. Column (2)-(4) are show the average absolute dollar amount of IT spending in billion of US dollars during 2010-2019. Column (5) shows the absolute value of log difference between IT spending under "algorithm plus manual adjustment for subsidiaries" and "algorithm-only" in Harte Hanks, that is, " $|\log(\mathrm{IT}_{\mathrm{Harte\ Hanks}}^{\mathrm{algorithm+manual\ adjustment}} - \log(\mathrm{IT}_{\mathrm{Harte\ Hanks}}^{\mathrm{algorithm-only}}|$ ." Column (6) shows the absolute value of log difference between IT

spending under "algorithm-only" in Harte Hanks and IT spending in Y9C, that is, " $|\log(\text{IT}_{Y9C}) - \log(\text{IT}_{\text{Harte Hanks}}^{\text{algorithm-only}}|$ .

Banks	Y9C	Manual Adjustment + Algorithm	Algorithm -only	Diff. Manual Adj +Algo vs Algo-only	Diff. Algo-only vs Y9C
(1)	(2)	(3)	(4)	(5)	(6)
Bank of America	4.19	3.72	1.33	1.03	1.14
Citibank	4.11	2.82	1.34	0.74	1.12
JP Morgan Chase Bank	3.45	3.04	1.23	0.90	1.03
Wells Fargo Bank	1.48	1.36	0.66	0.72	0.81
Suntrust bank	0.59	0.56	0.22	0.94	0.99
PNC Bank	0.51	0.44	0.24	0.60	0.76
U.S. Bank	0.51	0.47	0.23	0.70	0.78
TD Bank	0.35	0.35	0.16	0.79	0.77
Huntington National Bank	0.22	0.21	0.16	0.30	0.33
BMO Harris Bank	0.22	0.12	0.02	1.79	2.38
Branch BKG&TC	0.19	0.20	0.13	0.46	0.41
RBS Citizens	0.13 $0.17$	0.15	0.15	1.14	1.27
FIRST HORIZON NAT CORP	0.17	0.13	0.00	0.15	0.24
FIRST CITIZENS BSHRS	0.10	0.09	0.10	0.67	0.74
Commerce Bank	0.10	0.08	0.05	0.46	0.54
Bank of The West	0.09 $0.08$	0.08	0.05	0.39	0.54 $0.54$
Manufacturers & Traders Tc	0.08 $0.07$		0.05	0.42	
		0.07			0.47
Comerica Bank	0.06	0.05	0.05	0.18	0.30
Union Bank	0.06	0.06	0.03	0.50	0.55
Hancock Bk	0.06	0.06	0.05	0.07	0.08
First Niagara Bank	0.06	0.05	0.03	0.48	0.59
Synovus Bank	0.05	0.05	0.05	0.03	0.01
Zions First National Bank	0.05	0.04	0.04	0.13	0.27
CITY NAT CORP	0.04	0.05	0.04	0.11	0.03
TCF Bank	0.04	0.04	0.04	0.21	0.18
CENTRAL BANCOMPANY INC	0.04	0.04	0.04	0.12	0.09
Arvest Bk	0.04	0.04	0.03	0.16	0.07
Iberiabank	0.03	0.03	0.03	0.01	0.10
New York Cmnty Bank	0.03	0.02	0.02	0.21	0.36
Peoples United Bank	0.03	0.03	0.03	0.00	0.08
Firstmerit Bk Na	0.03	0.03	0.03	0.00	0.07
Umpqua Bk	0.03	0.03	0.03	0.02	0.03
Astoria Bank	0.02	0.02	0.02	0.21	0.36
Webster Bank	0.02	0.02	0.02	0.30	0.17
Wintrust Bank	0.02	0.02	0.02	0.10	0.26
Investors Svg Bk	0.02	0.02	0.02	0.13	0.01
East West Bank	0.02	0.02	0.02	0.07	0.06
Washington Federal	0.02	0.02	0.01	0.38	0.33
Great West Bank	0.02	0.02	0.01	0.09	0.20
Western Alli	0.02	0.02	0.01	0.13	0.07
Valley Nb	0.01	0.01	0.01	0.00	0.02
Bankunited	0.01	0.01	0.01	0.38	0.30
Metro Bank	0.01	0.01	0.01	0.01	0.07
Prosperity Bk	0.01	0.01	0.01	0.14	0.12
Cathay Bank	0.01	0.01	0.01	0.07	0.10
Byline Bank	0.01	0.01	0.01	0.02	0.07
Bell Bank	0.01	0.01	0.01	0.39	0.45
SERVISFIRST BSHRS	0.00	0.00	0.00	0.05	0.03
NORTHERN BC	0.00	0.00	0.00	0.13	0.07
Provident Bank	0.00	0.00	0.00	0.04	0.02

# Table B4. Comparison result with Y-9C report

This table presents the results of a simple regression of the logarithm of total IT spending in Y-9C for banks satisfying the reporting asset threshold on the logarithm of total IT spending from Harte Hanks.

	ln(Total IT <sup>HH</sup> )	SEs
$ln(Total\ IT^{Y9C})$	0.935***	0.010
Constant	0.037	0.083
$R^2$	0.661	
N	4947	

### B.3 Additional Analysis on Measurement Related Issues

### B.3.1 Comparison with BEA industry-level data

The cross-validation check with the regulatory Y-9C data helps establish the reliability of the Harte Hanks data in measuring banks' total IT spending. A key feature of our work is the investigation on how banks adopt skills to deal with different types of information when facing credit demand shocks associated with different information nature. Therefore, measuring the detailed composition of IT spending profiles, in particular, the differential responses of software IT spending and communication IT spending is crucial to our analysis. It is thus also important for us to also conduct cross-checking of the Harte Hanks bank IT spending data in this aspect with other available sources.

In this section, we compare the banking sector's IT spending profile in software spending and communication spending implied by the Harte Hanks data with the profile of different types of information technology input documented by Bureau of Economic Analysis (BEA) at the industry level. Several recent works have utilized BEA data to study industries' input characteristics. For instance, Barkai (2020) utilizes BEA fixed asset tables to capture industry-level capital and labor share; McGrattan (2017) uses the BEA input-output (I-O) table to capture industries' investment in intangible input; Kwon et al. (2022) use BEA fixed asset tables to accumulation of intangible assets at industry-level, etc. The strength and helpfulness of I-O tables lie in that detailed categories of industries' input uses are captured. This allows us to compare with our data to see whether relative values of sub-categories of IT products purchased by banking industry in the BEA data matches with the sub-categories of IT spending in our data.

With the focus being on IT expenditures of the aggregate banking industry, which is equal to the total value of input purchases made by the banking sector, we follow the existing literature and utilize the BEA input-output use table to construct the industry-level input value of IT products. More specifically, we construct total investment in software ("software input") as the sum of the items prepackaged software (ENS1), custom software (ENS2), own account software (ENS3), software publishers (RD40), and computer systems design and related services (RD60). The total investments in communication ("communication input") as the sum of the items communications equipment (EP20) and communication structures (SU20). In calculating the total amount of IT product purchase in each category, we focus on the "Finance and Insurance Industries," which have NAICS codes starting with "52." Since the Harte Hanks data provides industry categorization using SIC codes, we cross-walk the NAICS codes with SIC codes to help us locate the finance industry in the Harte Hanks dataset corresponding to "Finance and Insurance Industries (NAICS52)."

Table B5 provides a summary of our comparison statistics. In column (1), we calculate the total software and communication spending at industry-level in the Harte Hanks dataset during 2010-2019 for banking industry. In column (2), we calculate the total software spending and communication spending of the banking industry using the BEA input-output use table during

2010 and 2019. Importantly, we compare the input profile as captured by the relative sizes between the communication input and software input by comparing  $\frac{\text{\$ Communication spending Harte Hanks}}{\text{\$ Software spending}}$  with  $\frac{\text{\$ Communication input}}{\text{\$ EA}}$ 

\$ Software input

As is reported in the third row of the table, the "communication/software" ratio during 2010-2019 captured by the Harte Hanks dataset is 0.1465, which is very close to the "communication/software" ratio (0.1305) as implied by the BEA.

#### Brief introduction on the BEA input-output use tables and construction methods:

BEA I-O use tables measure how the supply of goods and services is used by different sectors, including domestic purchases by industries, individuals, and government, and exports to foreign purchasers. The construction methods of these I-O use tables are described as follows:

- Companies' income statements (main), take important items for I-O tables such as companies' revenues and expenses. Income statements come from census surveys, Annual Industry Survey, Service Annual Survey, Transportation Survey, etc. These surveys are typically conducted every five years.
- 2. The basic methodology starts with the value of an industry' cost of goods sold, say industry A, for the various intermediate inputs for producing A's output (say one of the inputs come from industry B), minus the net change in inventories of the input throughout the year, this gives an estimate for the use of goods/services from industry B.

I-O Data Sources and Basic Accounting Methods: The input-output table provides the total value of transactions of input a industry expend to make their own production—"A transaction is an economic flow between establishments or from an establishment to a final user." Transactions and values of transactions in BEA I-O tables are estimated based on certain data sources and estimation methodologies depending on industry categories. For industries such as manufacturing and mining, data such as materials consumed, input category controls directly related to a specific commodity, and commodity flow, are available. For industries where these information is not directly available, various sources of data, are used to make estimations. The important information here is the estimation of "input category controls," namely the component of input of goods/commodities/services from other industries. The underlying data sources are the surveys run by industry-specific agencies. For example, Farm Cost and Return Survey (Department of Agriculture), for wholesale, retail, and services, it is the Business Expenses Survey (BES).

More details about the methodologies and accounting method of the BEA data can be found at the BEA official website and the BEA Input-output Accounts Manual.

### B.3.2 Comparison with Modi et al. (2022)

In this section, we provide detailed descriptions for the comparison of bank-year level IT spending in our updated sample with the bank-year level of IT adoption measure constructed by Modi

et al. (2022). Figure ?? shows the binned scatter plot between of bank-year level IT spending in Modi et al. (2022) and the IT spending in the Harte Hanks data under the our updated algorithm described in Section B.1.1. We find a decent correlation of the IT expenses between the two data sets with an R square of 0.6.

We further compare the time trend of IT spending by size deciles between our measure and that from Modi et al. (2022). The illustration is in Figure B12. In Panel A, we follow the same grouping rule as in Modi et al. (2022) to split our sample by bank average asset across the sampling period into 10 groups, and then plot the ratio total of group total IT spending over group total non-interest expenses. For comparison, Panel B of Figure 2 displays the evolution of IT to non-interest expense ratios by bank size deciles in Modi et al. (2022). As shown by the plots, despite having different levels of IT spending intensity (the reason of which is discussed below), both plots show a consistent diverging trend between the largest decile and other asset groups after around 2015.

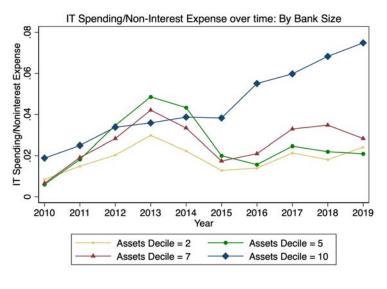
It is worthy noting that IT expenses in our updated sample are systemically higher than those in Modi et al. (2022). Modi et al. (2022) extract IT expenses among the top-3 reported unclassified non-interest expenses in Call Report using keywords such as "web," "IT," etc. As discussed by the authors, the IT expenses in Modi et al. (2022) could be underestimating Banks' IT expenses for two reasons. "First, some expenditures may be reported without reference to their IT content....Second, an IT-related expenditure may not be reported because it does not meet the criteria of being among the top three non-interest expenses and exceeding the minimum reporting threshold." These factors may help explain the discrepancy between the two datasets, for example, the measure constructed in Modi et al. (2022) does not include "equipment maintenance" and IT-related "consulting fees," but these expenses are included in our dataset.

### B.3.3 Comparison with Feyer et al. (2021)

In Figure B13, we provide a comparison of the average data storage costs in our data set with data storage costs constructed in Feyen et al. (2021). We calculate the average storage cost in our data set as the ratio between total storage amount (measured in PB) and the total dollar amount spent on storage devices. We calculate the average storage cost for each banks and then take the average among all banks in each year. In the top panel of the Figure B13, we plot the average storage costs measured in million dollars per PB of banks in our sample, which is equivalent to dollar per GB as is shown in the Figure 2 of Feyen et al. (2021). As is shown in the two figures, the average cost per dollar both declined from around \$0.7 per BG in 2010 to around \$0.15 per GB in 2020.

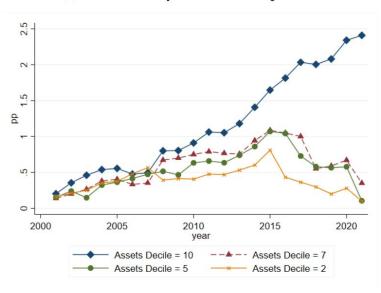
# Figure B12. Comparison with Modi et al. (2022) – Time trend and cross-size comparison

The figures compares the time trend of banks' IT spending between Harte Hanks and that from Modi et al. (2022). For panel A, the vertical axis is banks' total IT spending scaled by non-interest expenses. The asset size groups are categorized by 10 asset deciles. For Panel B, the y-axis is total IT spending as a share of non-interest expenses in percentage points (pp), and Panel B is the panel B of figure 2 in Modi et al. (2022).



(a) Figure A

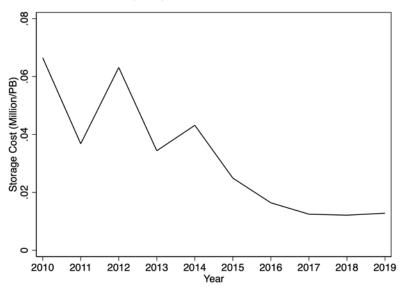
### (b) Normalized by non-interest expenses

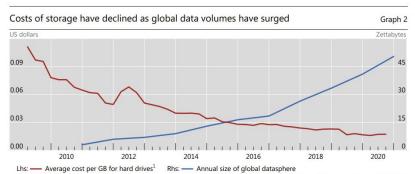


(b) Figure B

### Figure B13. Comparing storage cost in the data with other sources

This figures below show the comparison of per-unit storage costs in our sample and the per-unit storage cost documented by Feyen et al. (2021). The per-unit storage cost in Harte Hanks is calculated as million dollars per PB. The storage costs in Feyen et al. (2021) is calculated as dollars per GB. 1 PB is equal to  $10^6$  GB.





From Sep 2017, data extrapolated using the growth rate in price per MB from http://www.jcmit.net/diskprice.htm. The increase in 2012 is explained by flooding in Thailand, where one-third of hard drives were produced globally. One zettabyte is one trillion gigabytes.

Sources: BACKBLAZE: jcmit.net/diskprice; SEAGATE (2018), The digitization of the world from edge to core, November.

## Table B5. Comparison with BEA

This table compares the 10-year average of spending from 2011 to 2019 in communication-related technology and software-related technology between Harte Hanks and BEA. For data from BEA, software technology is defined as the sum of investments in prepackaged software, custom software, own account software, software publishers, and computer systems design and related services; communication technology is defined as the sum of investments in communications equipment and communication structures.

Million USD (2010-2019 Average)	Harte Hanks (1)	BEA (2)
Software	67687.21	59042.33
Communication	9918.00	7706.11
Ratio	14.65%	13.05%

### B.4 Details of IT Data Collection by Harte Hanks

In this section, we provide supplemental materials and additional analysis regarding the IT data collection by Harte Hanks, which are summarized as below:

- A polynomial regression analysis to further investigate the possibility that site-level IT spending is imputed by modeling, based on a regression of IT spending on a 3-order polynomial of the number of employees and revenue at the site level, controlling for firm fixed effects. See Table B6.
- A calculation comparing IT spending at headquarters and non-headquarters, which helps address issues related to how large bank-level IT purchases are being accounted in the Harte Hanks dataset. See Table B7.
- A calculation of the within-bank variation in the IT spending across sites. See Table B8.

Table B6. Third-order Polynomial Fit between IT Spending and Revenue and Number of Employees

This table assesses the branch-level relationship between IT spending and branch characteristics using third-order polynomials. Total, Hardware, Software, and Communication are respectively the branch-level spending on Total IT, Hardware IT, Software IT, and Communication IT. Emp is the number of employees, and Rev is the total amount of revenues.

	Total	Hardware	Software	Communication
	(1)	(2)	(3)	(4)
Emp	-1.307	-0.603	-0.236	-0.327
	(1.676)	(0.574)	(0.155)	(0.226)
$\mathrm{Emp}^2$	0.000**	0.000**	0.000**	0.000**
	(0.000)	(0.000)	(0.000)	(0.000)
${ m Emp^3}$	-0.000*	-0.000**	-0.000**	-0.000**
	(0.000)	(0.000)	(0.000)	(0.000)
Rev	41.056***	10.738***	2.296***	2.905***
	(2.802)	(0.834)	(0.164)	(0.199)
$\mathrm{Rev}^2$	-0.000**	-0.000**	-0.000***	-0.000***
	(0.000)	(0.000)	(0.000)	(0.000)
$\mathrm{Rev}^3$	0.000	0.000	0.000***	0.000**
	(0.000)	(0.000)	(0.000)	(0.000)
Zipcode FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
N	11960349	11960349	11960349	11960349
$\mathbb{R}^2$	0.575	0.545	0.502	0.558
$AdjR^2$	0.470	0.430	0.380	0.454

## Table B7. IT spending at Headquarters and Non-headquarters

This table assesses whether headquarters spend asymmetrically more on IT. Headerquarter=1 is a dummy variable that labels if a site in the Harte Hanks data is labeled as a headquarter. In each column, we include bank×year fixed effects.

	Total/Revenue	Communication/Revenue	Software/Revenue	Hardware/Revenue	Other/Revenue
	(1)	(2)	(3)	(4)	(5)
Headquarter=1	0.788	0.00591	0.217	0.0164	-0.00803
	(0.519)	(0.0326)	(0.151)	(0.0556)	(0.0110)
Fixed effects			Bank × Year		
AdR-squared	0.332	0.304	0.342	0.358	0.276
N	804000	804000	804000	804000	808000

Standard errors in parentheses

<sup>\*</sup> p < 0.10, \*\* p < 0.05, \*\*\* p < 0.01

# Table B8. Variation in within-bank IT Spending across branches

This table presents the within-bank variation in IT spending across counties. Panel A (B) gives the average/median of IQR (standard deviation) over average total IT spending for each bank across different counties.

Panel A: 1	IQR/Mean	Panel B: SD/Mean		
Mean	Median	Mean	Median	
67.81%	47.97%	166.02%	118.22%	

### C Other Data Construction

This Online Appendix consists of three sections. In Section 1, we describe how we extract and clean IT spending data in Harte Hanks. In Section 2, we provide data-construction details on how to map banks in the IT spending data set with "Summary of Deposits" and "Call Report." In Section 3, we provide the construction details of the other supporting data sets utilized in the paper.

### C.1 Construction of IT Data

In this part, we provide details on how to extract the relevant information in the original Harte Hanks IT data sets and create the panel data of banks' characteristics and IT spending at bank-county-year level. Harte Hanks collects the establishment-level (hereafter "site-level") information on IT spending and the characteristic annually. For each given year, the site-level IT spending and site characteristics are saved in two different files, "IT Spend" and "Site Description," respectively. We extract site level variables from the two files and combine them together to get the panel data of site-level IT spending and characteristics. First, for each year, from the site characteristics file, we restrict the data to one-digit SIC code equal to 6. Then we keep "site ID," company name, location (zip code), homepage url, revenue, and number of employees as our site level characteristics variables. "Site ID" is the unique identifier of the site across years in the Harte Hanks data. Second, from the IT Spend file, we get the site-level IT budget data including total budget, communication budget, software budget, services budget, hardware budget, etc, as well as the site ID. Then we merge the site characteristics and site-level IT Spending using "Site ID." Repeating the process for each year gives us a panel data set of site-level IT spending information and site level characteristics.

Next, we aggregate the number of sites, IT spending variables, revenues and employees at the zip code-year-bank level. Most sites include a url variable that labels the homepage website address of the bank. When aggregating site-level variables into county level, we first aggregate the variables by url. For those sites without a url, we aggregate by the cleaned company names. Cleaned company names are defined as the lowercase of company names after removing "national association", "n.a.", "fsb", "s.b." etc. This gives us IT spending profile and revenue and employee profile of a bank at zip code and year level.

Finally, we crosswalk zip codes to fips code (the commonly used county identifier) and aggregate all the variables at the county level, this gives us banks' IT budget and characteristics at bank-county-year level. When mapping zip codes into fips code, we noticed that some zip codes are mapped into multiple different counties. This is because some zip code areas are at the border of multiple different counties and some of the businesses or residents reside in one county while the rest of the zip code's businesses or residents are located in the other counties. For instance, zip code 49963 is mapped into both "Houghton, MI" and "Ontonagon, MI." To correctly account for the IT spending of banks located in zip codes like this into the two counties, we multiply the IT

spending of a bank in this zip code with the ratios of addresses in this zip code that belongs to the two counties, before aggregating to county level IT spending. In the above example, 23% of the addresses in zip code 49663 belongs to "Houghton, MI," while 77% of the addresses in zip code 49663 belongs to "Ontonagon, MI." We multiply a bank's IT spending in zip code 49663 by 0.23 and aggregate this adjusted number to "Houghton, MI"; we multiple a bank's IT spending in zip code 49663 by 0.77 before and aggregate the adjusted number to "Ontonagon, MI." We obtain the the information on the ratio of a zip code's addresses that belong to each county for each zip code from the Office of Policy Development and Research. We use "TOTAL RATIO" provided by the Office of Policy Development and Research, which is the ratio of all types of addresses in the zip code that belongs to a county, to adjust for the spending before aggregation.<sup>5</sup>

#### C.1.1 Matching Bank Names in Two Datasets

We now explain the procedure to match bank names in the IT Dataset and those in Summary of Deposits.

Matching at Bank-Year Level This subsection describes how do we match bank names in Summary of Deposits (hereafter SOD) and bank names in the IT data at the bank level and construct the panel data containing bank IT and bank characteristics at bank-year level.

To start with, we take the bank names from SOD data set and obtain the banks' homepage from Google. The first step is to extract a smaller set of site names in the site-level bank IT data that are similar to the names of the banks in SOD. We drop the suffixes ", national association", "national association", ", fsb", "fsb", ", n.a.", "n.a.", "f.s.b.", " f. s. b.", ", f. s. b.", ", s.b.", ", s/b", ", s.b.", ", ssb", ", s.s.b.", " (west), fsb", ", fsb", ", fsb", ", a fsb", and ", a federal savings and loan association", "bank", and "national bank", etc in the SOD data. We split the names into at most two key words by spaces. For example, Wells Fargo Bank is labeled as "Wells" and "Fargo." This is because many site names in the IT data set, which is going to be merged later, are written without spaces. In the Wells Fargo Bank case, the site names could be written as "wellsfargo bank" or "thewellsfargobank." Given that most sites in the IT data set also include a url variable that label the website address of the bank's homepage, we conduct the matching using url first, and if matching with url doesn't work, we match using keywords in names constructed above. For those sites with url, we first outer merge the names from the "Site Description" files with the url of the banks' website address (after dropping "www." and ".com"), and retain the sites whose url contains the url of the banks' website address. Then we can match sites names with SOD bank names with the url's. For those sites without a url, we outer merge the names from the "Site Description" files with the key words constructed in SOD, and only keep the sites of which the names contain all the keywords from SOD. These above procedures give us the extracted site names whose names as close to names of banks in SOD. In the last step, we assign these extracted sites names with a

<sup>&</sup>lt;sup>5</sup>See the link to Office of Policy Development and Research's webpage for the relevant files.

bank name from the names of banks in SOD that has the largest Levenshtein score. We aggregate the site level IT Spending using the matched bank names in the Bank IT data and merge with SOD through the matched bank names. This gives us the panel data of banks IT spending, total assets and deposits in SOD, matched bank names, and the bank identifier RSSDID in SOD, at bank-year level. The bank identifier RSSDID is also utilized to merge with other data sets such as "Call Reports" and HMDA, etc.

Matching at Bank-County-Year Level In this subsection, we describe how we match the banks in the IT data constructed in Section 1 with banks in Summary of Deposits (SOD) and merge the IT data with bank characteristics in SOD at county level. The output of this matching procedure generates the panel data on banks' IT spending (from the IT data) and bank assets, deposits, and bank identifier (from the SOD) at bank-county-year level.

We get total assets, year, total deposits and bank names from the SOD data set. We convert the bank names in SOD data set to its lower case. We drop the suffixes ", national association", "national association", "fsb", "fsb", "n.a.", "n.a.", "f.s.b.", "f. s. b.", ", f. s. b.", ", s.b.", ", s.b.", ", s.b.", "fsb", "fsb", ", fsb", ", fsb", ", a fsb", and ", a federal savings and loan association". We then collapse the above SOD information by zip code, bank name and year. This gives us a panel of banks in SOD with bank names, location information (zip codes), total assets, total deposits, RSSDID, and year.

To match the SOD panel data with the IT spending panel data, we first merge these two data sets using zip code and year, this gives us all the possible pairs of bank names in SOD and IT spending for each combination of zip code and year. Then for each bank name showing up in IT spending data at the zip code-year level, we calculate the Levenshtein distance of the string names between this bank name string and all the string names showing up in SOD within the same zip code and year. For the merged observations, we keep the RSSDID (the unique bank identifier) from the SOD data set with the highest Levenshtein score and we keep the observations with the calculated highest Levenshtein score larger than 2/3. This gives us a panel data of banks at zip code level that match with banks in SOD at zip code level, other variables include IT spending, bank identifier (RSSDID), total assets and total deposits.

Finally, we employ the same method described in Part I to aggregate the matched bank IT spending panel data at zip code-year level to county-year level by adjusting the ratio of addresses of a zip code that could potentially show up in multiple counties. This gives us a panel data of banks' IT spending and banks' deposits and assets, and RSSDID (bank identifier) at county-year level.

### C.2 Construction of Other Data Sets

Call Report In this subsection, we describe how do we construct loan portfolio information and bank-level control variables using "Call Reports." We get the banks' balance sheet information from

"Call Reports" quarterly data for the year of 2010-2019. We collapse the key variables by the last quarter of a bank within a year. The linkage between Call Reports and our IT data set is through RSSDID. We define C&I loan share as the "ciloans" scaled by "qavgloans," we define personal loan share as (personal loans) "persloans" scaled by "qaveloans," and we define agriculture loan share as "agloans" scaled by "qavgloans." Banks' Profitability ("prof") is defined as net income (netinc) scaled by "qavgassets," Equity/assets is defined as "equity" scaled by "assets," Deposits/assets is defined as "dep" scaled by "assets," Salary/assets is defined as "sal" scaled by "assets," and number of employees per thousand dollar assets is defined as number of employees (nume multiplied by 1000) scaled by assets (this is because number of employees is in the unit of 1000), we define revenue per employee as income scaled by number of employees.

Home Mortgage Disclosure Act Data This subsection describes the construction of refinancing and origination amount for each bank in each year in a county.

We use the panel of "HMDA nationwide records" files to construct origination and refinance volumes. We define loan as origination if "loan purpose" is equal to 1 and define loan as refinance if loan "loan purpose" is equal to 3. We aggregate the total loan amount of each bank (identified by respondent id) at state code-county code and year level, for origination and refinance, respectively. We then construct the fips code (county identifier) by combining the state code with county code. Finally, we crosswalk respondent id to RSSDID provided in the "HMDA institution" files.

Freddie Mac Single-Family Loan-Level Data Set This subsection describe how do we construct the potential mortgage repayment savings using the Freddie Mac Single Family Loan Data Set at the county-year level.

We first use the Historical data to get the average interest rate between 2010 and 2015 at the zip code-maturity-FICO group level. Specifically for each year, we assign loans into 12 FICO bins: <620, ..., 780-800, 800-820, and >820. We then calculate the average interest rate by year, zip code, FICO group and maturity. We then use loans originated between 1999 and 2009 from the Historical Time Data to get the payment savings. Specifically, for each loan originated between 1999 and 2009, we first keep those that are not pre-paid or defaulted between 2010 and 2015, and the calculate the remain balance separately based on the loans' original interest rate and the hypothetical interest rate as the zip code-maturity-fico group average from 2010 to 2015. We then take the average payment saving by each zip code, and aggregate the data to the county level.

Mergent Intellect This subsection describes how we construct the bank hierarchical structure data using Mergent Intellect.

We download all the information of banks' family trees in Mergent Intellect with two-digit SIC code "60" and "61." We replace "Domestic Parent Name" with "Company Name" if an entity's "Domestic Parent Name" is missing. We then sum up the number of "Headquarters," number of "Single Location" and number of "Branch" offices within the "Domestic Parent Name." Banks with

only one type of locations is defined as having 1 layer in their hierarchy; banks with two different types of locations is defined as having 2 layers in their hierarchy and banks with three different types of locations is defined as having 3 layers in their hierarchy.

To match the bank names in the cleaned version of Mergent Intellect as described above, we link bank names in Mergent Intellect with the institution names provided by FDIC and then link the matched results with banks in our sample. Specifically, we first remove the words "Bank", "INC", "National Association", "LLC", "CORPORATION", "COMPANY", "THE", "CORP", "SERVICES" from the headquarters' names in Mergent Intellect, and unify names of entities within the Mergent Intellect, then we append cities' names where the banks are located to bank names. Next, we repeat the same process with our sample data. Finally, we merge bank names and cities in the Mergent Intellect with bank names and cities in our sample data using Jarodistance algorithm and keep the matched pairs with the highest Jarodistance score.