

# Accepted Manuscript

Big Data in Finance and the Growth of Large Firms

Juliane Begenau, Maryam Farboodi, Laura Veldkamp

PII: S0304-3932(18)30217-4  
DOI: [10.1016/j.jmoneco.2018.05.013](https://doi.org/10.1016/j.jmoneco.2018.05.013)  
Reference: MONEC 3009

To appear in: *Journal of Monetary Economics*

Received date: 21 April 2018  
Accepted date: 31 May 2018

Please cite this article as: Juliane Begenau, Maryam Farboodi, Laura Veldkamp, Big Data in Finance and the Growth of Large Firms, *Journal of Monetary Economics* (2018), doi: [10.1016/j.jmoneco.2018.05.013](https://doi.org/10.1016/j.jmoneco.2018.05.013)



This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Big Data in Finance and the Growth of Large Firms

Juliane Begenau\*      Maryam Farboodi      Laura Veldkamp  
Stanford GSB & NBER      Princeton      NYU Stern & NBER

June 7, 2018<sup>†</sup>

## Abstract

Two modern economic trends are the increase in firm size and advances in information technology. We explore the hypothesis that big data disproportionately benefits big firms. Because they have more economic activity and a longer firm history, large firms have produced more data. As processor speed rises, abundant data attracts more financial analysis. Data analysis improves investors' forecasts and reduces equity uncertainty, reducing the firm's cost of capital. When investors can process more data, large firm investment costs fall by more, enabling large firms to grow larger.

**JEL Codes:** E2, G1, D8.

**Keywords:** Big data, Fintech, Firm size

---

\*Corresponding author: Juliane Begenau, Stanford University Graduate School of Business, 655 Knight Way, Stanford, CA 94305, [begenau@stanford.edu](mailto:begenau@stanford.edu), 650-724-5661.

<sup>†</sup>This paper was prepared for the Carnegie-Rochester-NYU conference. We thank the conference committee for their support of this work. We also thank Nic Kozeniauskas for his valuable assistance with the data and Adam Lee for his outstanding research assistance.

15 One of the main question in macroeconomics today is why small firms are being replaced  
16 with larger ones. Over the last three decades, the percentage of employment at firms with  
17 less than 100 employees has fallen from 40% to 35% (Figure 1a); the annual rate of new  
18 startups has decreased from 13% to less than 8%, and the share of employment at young  
19 firms (less than 5 years) has decreased from 18% to 8% (Davis and Haltiwanger, 2015).  
20 While small firms have struggled, large firms (more than 1000 employees) have thrived: The  
21 share of the U.S. labor force they employ has risen from one quarter in the 1980s, to about  
22 a third today. At the same time, the revenue share of the top 5% of firms increased from  
23 57% to 67% (Figure 1b).

24 *Figure 1 about here.*

25 One important difference between large and small firms is their cost of capital (Cooley  
26 and Quadrini, 2001). Hennessy and Whited (2007) document that larger firms, with larger  
27 revenues, more stable revenue streams, and more collateralizable equipment, are less risky  
28 to creditors and thus pay lower risk premia. But this explanation for the trend in firm size  
29 is challenged by the fact that while small firms are more volatile, the volatility gap between  
30 small firms and large firms cash flows has not grown.<sup>1</sup> Alternative, the trend in covariance  
31 of firm stock prices with market portfolio, as measured by CAPM  $\beta$ , is not significantly  
32 different across firms of different sizes.

33 If neither volatility nor covariance with market risk has diverged, how could risk premia  
34 and thus the cost of capital diverge? What introduces a wedge between unconditional vari-  
35 ance or covariance and risk is information. Even if the payoff variance is constant, better  
36 information can make payoffs more predictable and therefore less uncertain. Given this new

---

<sup>1</sup>Evidence on the volatility gap between large and small firms is in Appendix A.4. Other hypotheses are that the productivity of large firms has increased or that potential entrepreneurs instead work for large firms. This could be because of globalization, or the skill-biased nature of technological change as in Kozeniauskas (2017). These explanations are not exclusive and may each explain some of the change in the distribution of firm size.

37 data, the conditional payoff variance and covariance fall. More predictable payoffs lower  
38 risk and lower the cost of capital. The strong link between information and the cost of  
39 capital is supported empirically by Fang and Peress (2009), who find that media coverage  
40 lowers the expected return on stocks that are more widely covered. This line of reasoning  
41 points to an information-related trend in financial markets that has affected the abundance  
42 of information about large firms relative to small firms. What is this big trend in financial  
43 information? It is the big data revolution.

44 Our goal is to explore the hypothesis that the use of big data in financial markets has  
45 lowered the cost of capital for large firms relative to small ones, enabling large firms to grow  
46 larger. In modern financial markets, information technology is pervasive and transformative.  
47 Faster and faster processors crunch ever more data: macro announcements, earnings state-  
48 ments, competitors' performance metrics, export market demand, anything and everything  
49 that might possibly forecast future returns. This data informs the expectations of modern  
50 investors and reduces their uncertainty about investment outcomes. More data processing  
51 lowers uncertainty, which reduces risk premia and the cost of capital, making investments  
52 more attractive.

53 To explore and quantify these trends in modern computing and finance, we use a noisy  
54 rational expectations model where investors choose how to allocate digital bits of information  
55 processing power among various firm risks, and then use that processed information to solve  
56 a portfolio problem. The key insight of the model is that the investment-stimulating effect  
57 of big data is not spread evenly across firms. Small firms benefit less. In our model, small  
58 firms are equivalent to young firms, and large firms to old firms. This is consistent with  
59 the data, where age and size are positively correlated. In the model, larger firms are more  
60 valuable targets for data analysis because more economic activity and a longer firm history  
61 generates more data to process. In contrast, all the computing power in the world cannot  
62 inform an investor about a small firm that has a short history with few disclosures. As big

63 data technology improves, large firms attract a more than proportional share of the data  
64 processing. Because data resolves risk, the gap in the risk premia between large and small  
65 firms widens. Such an asset pricing pattern enables large firms to invest cheaply and grow  
66 larger.

67 The data side of the model builds on theory designed to explain human information  
68 processing (Kacperczyk et al., 2016), and embeds it into a standard model of corporate  
69 finance and investment decisions (Gomes, 2001). In this type of model, deviations from  
70 Modigliani-Miller imply that the cost of capital matters for firms' investment decisions. In  
71 our model, the only friction affecting the cost of capital works through the information  
72 channel. The big data allocation model can be reduced to a sequence of required returns for  
73 each firm that depends on the data-processing ability and firm size. These required returns  
74 can then be plugged into a standard firm investment model. To keep things as simple as  
75 possible, we study the big-data effect on firms' investment decision based on a simulated  
76 sample of firms – two, in our case – in the spirit of Hennessy and Whited (2007).

77 The key link between data and real investment is the price of newly-issued equity. Assets  
78 in this economy are priced according to a conditional CAPM, where the conditional variance  
79 and covariance are those of a fictitious investor who has the average precision of all investors'  
80 information. The more data the average investor processes about an asset's payoff, the  
81 lower is the asset's conditional variance and covariance with the market. A researcher who  
82 estimated a traditional, unconditional CAPM would attribute these changes to a relative  
83 decline in the excess returns (alphas) on small firms. Thus, the widening spread in data  
84 analysis implies that the alphas of small firm stocks have fallen relative to larger firms. These  
85 asset pricing moments are new testable model predictions that can be used to evaluate and  
86 refine big data investment theories.

87 This model serves both to expoit a new mechanism and as a framework for measurement.  
88 Obviously, there are other forces that affect firm size. We do not build in many other

89 contributing factors. Instead, we opt to keep our model stylized, which allows a transparent  
90 analysis of the new role that big data plays. Our question is simply how much of the change  
91 in the size distribution is this big data mechanism capable of explaining? We use data in  
92 combination with the model to understand how changes in the amount of data processed over  
93 time affect asset prices of large and small public firms, and how these trends reconcile with  
94 the size trends in the full sample of firms. An additional challenge is measuring the amount  
95 of data. Using information metrics from computer science, we can map the growth of CPU  
96 speeds to signal precisions in our model. By calibrating the model parameters to match the  
97 size of risk premia, price informativeness, initial firm size and volatility, we can determine  
98 whether the effect of big data on firms' cost of capital is trivial or if it is a potentially  
99 substantial contributor to the missing small firm puzzle.

100 **Contribution to the existing literature** Our model combines features from a few dis-  
101 parate literatures. The topic of changes in the firm size distribution is a topic taken up  
102 in many recent papers, including Davis and Haltiwanger (2015), Kozeniauskas (2017), and  
103 Akcigit and Kerr (2017). In addition, a number of papers analyze how size affects the cost  
104 of capital, e.g. Cooley and Quadrini (2001), Hennessy and Whited (2007), and Begenau  
105 and Salomao (2018). We explore a very different force that affects firm size and quantify its  
106 effect.

107 Another strand of literature explores the feedback between information in financial mar-  
108 kets and investment: Maksimovic et al. (1999) models the relationship between a firm's cap-  
109 ital structure and its information acquisition prior to capital budgeting decisions. Bernhardt  
110 et al. (1995) studies the effect of different levels of insider trading on investment. Ozdenoren  
111 and Yuan (2008) studies a setting where asset prices influence fundamentals through coor-  
112 dinated buying and thus self-fulfilling beliefs. Furthermore, there are papers that focus on  
113 long run data or information trends in finance: Asriyan and Vanasco (2014), Biais et al.

114 (2015) and Glode et al. (2012) model growth in fundamental analysis or an increase in its  
115 speed. The idea of long-run growth in information processing is supported by the rise in  
116 price informativeness documented by Bai et al. (2016).

117 Over time, it has gotten easier and easier to process large amounts of data. As in Farboodi  
118 et al. (2017), this growing amount of data reduces the uncertainty of investing in a given  
119 firm. But the new idea that this paper adds to the existing work on data and information  
120 frictions, is this: Intensive data crunching works well to reduce uncertainty about large firms  
121 with long histories and abundant data. For smaller firms, who tend also to be younger firms,  
122 data may be scarce. Big data technology only reduces uncertainty if abundant data exists to  
123 process. Thus as big data technology has improved, the investment uncertainty gap between  
124 large and small firms has widened, their costs of financing have diverged, and big firms have  
125 grown ever bigger.

## 126 **1 Model**

127 We develop a model whose purpose is to understand how the growth in big data technologies  
128 in finance affects firm size and gauge the size of that effect. The model builds on the  
129 information choice model in Kacperczyk et al. (2016) and Kacperczyk et al. (2015).

### 130 **1.1 Setup**

131 This is a repeated, static model. Each period has the following sequence of events. First,  
132 firms choose entry and firm size. Second, investors choose how to allocate their data process-  
133 ing across different assets. Third, all investors choose their portfolios of risky and riskless  
134 assets. At the end of the period, asset payoffs and utility are realized. The next period, new  
135 investors arrive and the same sequence repeats. What changes between periods is that firms  
136 accumulate capital and the ability to process big data grows over time.

137 **Firm Decisions** We assume that firms are equity financed. Each firm  $i$  has a profitable  
 138 1-period investment opportunity and wants to issue new equity to raise capital for that  
 139 investment. For every share of capital invested, the firm can produce a stochastic payoff  
 140  $f_{i,t}$ . Thus total firm output depends on the scale of the investment, which is the number of  
 141 shares  $\bar{x}_{i,t}$ , and the output per share  $f_{i,t}$ :

$$y_{i,t} = \bar{x}_{i,t} f_{i,t}. \quad (1)$$

142 The owner of the firm chooses how many shares  $\bar{x}_{i,t}$  to issue. The owner's objective is  
 143 to maximize the revenue raised from the sale of the firm, net of the setup or investment cost

$$\tilde{\phi}(\bar{x}_{i,t}, \bar{x}_{i,t-1}) = \phi_0 \mathbf{1}_{(|\Delta \bar{x}_{i,t}| > 0)} + \phi_1 |\Delta \bar{x}_{i,t}| + \phi_2 (\Delta \bar{x}_{i,t})^2, \quad (2)$$

144 where  $\Delta \bar{x}_{i,t} = \bar{x}_{i,t} - \bar{x}_{i,t-1}$ ,  $\mathbf{1}_{|\Delta \bar{x}_{i,t}| > 0}$  is an indicator function taking the value of one if  $|\Delta \bar{x}_{i,t}|$   
 145 is strictly positive and  $\phi_0, \phi_1, \phi_2 > 0$ . This cost function represents the idea that issuing new  
 146 equity (or buying equity back) has a fixed cost  $\phi_0$  and a marginal cost that is increasing in  
 147 the number of new shares issued. Each share sells at price  $p_{i,t}$ , which is determined by the  
 148 investment market equilibrium. The owner's objective is thus

$$Ev_{i,t} = E[\bar{x}_{i,t} p_{i,t} - \tilde{\phi}(\bar{x}_{i,t}, \bar{x}_{i,t-1}) | \mathcal{I}_{t-1}], \quad (3)$$

149 which is the expected net revenue from the sale of firm  $i$ .

150 The firm makes its choice conditional on the same prior information that all the investors  
 151 have and understanding the equilibrium behavior of investors in the asset market. But the  
 152 firm does not condition on  $p_{i,t}$ . In other words, it does not take prices as given. Rather, the  
 153 firm chooses  $\bar{x}_{i,t}$ , taking into account its impact on the equilibrium price.



154 **Assets** The model features 1 riskless and  $n$  risky assets. The price of the riskless asset is  
 155 normalized to 1 and it pays off  $r_t$  at the end of period  $t$ . Risky assets  $i \in \{1, \dots, n\}$  have  
 156 random payoffs  $f_{i,t} \sim N(\mu, \Sigma)$ , where  $\Sigma$  is a diagonal “prior” variance matrix.<sup>2</sup> We define  
 157 the  $n \times 1$  vector  $f_t = [f_{1,t}, f_{2,t}, \dots, f_{n,t}]'$ .

158 Each asset has a stochastic supply given by  $\bar{x}_{i,t} + x_{i,t}$ , where noise  $x_{i,t}$  is normally  
 159 distributed, with mean zero, variance  $\sigma_x$ , and no correlation with other noises:  $x_t \sim$   
 160  $\mathcal{N}(0, \sigma_x I)$ . As in any (standard) noisy rational expectations equilibrium model, the supply  
 161 is random to prevent the price from fully revealing the information of informed investors.

162 **Portfolio Choice Problem** There is a continuum of measure one of atomless investors.  
 163 Each investor is endowed with beginning-of-period wealth,  $W_t$ .<sup>3</sup> They have mean-variance  
 164 preferences over end-of-period wealth, with a risk-aversion coefficient,  $\rho$ . Let  $\hat{E}_{j,t}$  and  $\hat{V}_{j,t}$   
 165 denote investor  $j$ 's period  $t$  expectations and variances conditioned on all interim information,  
 166 which includes prices and signals. Thus, investor  $j$  chooses how many shares of each asset  
 167 to hold,  $q_{j,t}$  to maximize period  $t$  interim expected utility,  $\hat{U}_{j,t}$ :

$$\hat{U}_{j,t} = \rho \hat{E}_{j,t}[\hat{W}_{j,t}] - \frac{\rho^2}{2} \hat{V}_{j,t}[\hat{W}_{j,t}], \quad (4)$$

168 subject to the budget constraint:

$$\hat{W}_{j,t} = r_t W_t + q'_{j,t}(f_t - p_t r_t), \quad (5)$$

169 where  $q_{j,t}$  and  $p_t$  are  $n \times 1$  vectors of prices and quantities of each asset held by investor  $j$ .

---

<sup>2</sup>We can allow assets to be correlated. To solve a correlated asset problem simply requires constructing portfolios of assets (risk factors) that are independent from each other, choosing how much to invest and learn about these risk factors, and then projecting the solution back on the original asset space. See Kacperczyk et al. (2016) for such a solution.

<sup>3</sup>Since there are no wealth effects in the preferences, the assumption of identical initial wealth is without loss of generality.

170 **Prices** Equilibrium prices are determined by market clearing:

$$\int_0^1 q_{j,t} dj = \bar{x}_t + x_t, \quad (6)$$

171 where the left-hand side of the equation is the vector of aggregate demand and the right-hand  
172 side is the vector of aggregate supply of the assets.

173 **Information sets, updating, and data allocation** At the start of each period, each  
174 investor  $j$  chooses the amount of data that she will receive at the interim stage, before she  
175 invests. A piece of data is a signal about the risky asset payoff. A time- $t$  signal, indexed by  
176  $l$ , about asset  $i$  is  $\eta_{l,i,t} = f_{i,t} + e_{l,i,t}$ , where the data error  $e_{l,i,t}$  is independent across pieces  
177 of data  $l$ , across investors, across assets  $i$  and over time. Signal noise is normally distributed  
178 and unbiased:  $e_{l,i,t} \sim N(0, \sigma_e/\delta)$ . By Bayes' law, choosing to observe  $\mathcal{M}$  signals, each with  
179 signal noise variance  $\sigma_e/\delta$ , is equivalent to observing one signal with signal noise variance  
180  $\sigma_e/(\mathcal{M}\delta)$ , or equivalently, precision  $\mathcal{M}\delta/\sigma_e$ . The discreteness in signals complicates the  
181 analysis, without adding insight. But if we have a constraint that allows an investor to  
182 process  $\bar{\mathcal{M}}/\delta$  pieces of data, each with precision  $\delta/\sigma_e$ , and then we take the limit  $\delta \rightarrow 0$ , we  
183 get a quasi-continuous choice problem. The choice of how many pieces of data to process  
184 about each asset becomes equivalent to choosing  $K_{i,j,t}$ , the precision of investor  $j$ 's signal  
185 about asset  $i$  in period  $t$ . Investor  $j$ 's vector of data-equivalent signals about each asset  
186 is  $\eta_{j,t} = f_t + \varepsilon_{j,t}$ , where the vector of signal noise is distributed as  $\varepsilon_{j,t} \sim \mathcal{N}(0, \Sigma_{\eta,j,t})$ . The  
187 variance matrix  $\Sigma_{\eta,j,t}$  is diagonal with the  $i$ th diagonal element  $K_{i,j,t}^{-1}$ . Investors combine  
188 signal realizations with priors and information extracted from asset prices to update their  
189 beliefs using Bayes' law.

190 Signal precision choices  $\{K_{i,j,t}\}$  maximize start-of-period expected utility,  $U_{j,t}$ , of the  
191 fund's terminal wealth  $\hat{W}_{j,t}$ . Thus the objective is

$$\max_{\{K_{i,j,t}\}_{i=1}^n} E[\hat{U}_{j,t} | \mathcal{I}_{t-1}^+] \quad (7)$$

192

$$\text{where } \mathcal{I}_t = \{\mathcal{I}_{t-1}^+, \eta_{jt}, p_t\} \text{ and } \mathcal{I}_{t-1}^+ = \{\mathcal{I}_{t-1}, x_{t-1}, \tilde{f}_{t-1}\} \quad (8)$$

193 subject to the the budget constraint (5) and three constraints in the information choices.<sup>4</sup>194 The first constraint is the *information capacity constraint*. It states that the sum of the  
195 signal precisions must not exceed the information capacity:

$$\sum_{i=1}^n K_{i,j,t} \leq K_t \quad \text{for each } j, t. \quad (9)$$

196 In Bayesian updating with normal variables, observing one signal with precision  $K_{i,j,t}$  or  
197 two signals, each with precision  $K_{i,j,t}/2$ , is equivalent. Therefore, one interpretation of  
198 the capacity constraint is that it allows the manager to observe  $N$  signal draws, each with  
199 precision  $K_{i,j,t}/N$ , for large  $N$ . The investment manager then chooses how many of those  
200  $N$  signals will be about each shock.<sup>5</sup>201 The second constraint is the *data availability constraint*. It states that the amount of  
202 data processed about the future earnings of firm  $i$  cannot exceed the total data generated  
203 by the firm. Since data is a by-product of economic activity, data availability depends on  
204 the economic activity of the firm in the previous period. In other words, data availability in  
205 time  $t$  is a function of firm size in  $t - 1$ .

$$K_{i,j,t} \leq \hat{K}(x_{i,t-1}) \quad \text{for all } i, j, t. \quad (10)$$

---

<sup>4</sup>See Veldkamp (2011) for a discussion of the use of expected mean-variance utility in information choice problems.

<sup>5</sup>The results are not sensitive to the exact nature of the information capacity constraint. We could instead specify a cost function of data processing  $c(K_{i,j,t})$ . The problem we solve is the dual of this cost function approach. For any cost function, there exists a constraint value  $K_t$  such that the cost function problem and the constrained problem yield identical solutions.

206 This limit on data availability is a new feature of the model. It is also what links firm size  
 207 to the expected cost of capital.<sup>6</sup> We assume that the data availability constraint takes a  
 208 simple, exponential form:  $\hat{K}(x_{i,t-1}) = \alpha \exp(\beta x_{i,t-1})$ .

209 The third constraint is the *no-forgetting constraint*, which ensures that the chosen preci-  
 210 sions are non-negative:

$$K_{i,j,t} \geq 0 \quad \text{for all } i, j, t. \quad (11)$$

211 It prevents the manager from erasing any prior information to make room to gather new  
 212 information about another asset.

## 213 1.2 Equilibrium

214 To solve the model, we begin by working through the mechanics of Bayesian updating. There  
 215 are three types of information that are aggregated in posterior beliefs: prior beliefs, price  
 216 information, and (private) signals. We conjecture and later verify that a transformation  
 217 of prices  $p_t$  generates an unbiased signal about the risky payoffs,  $\eta_{p,t} = f_t + \epsilon_{p,t}$ , where  
 218  $\epsilon_{p,t} \sim N(0, \Sigma_{p,t})$ , for some diagonal variance matrix  $\Sigma_{p,t}$ . Then, by Bayes' law, the posterior  
 219 beliefs about  $f_t$  are normally distributed:  $f_t \sim N(\hat{E}_{j,t}[f_t], \hat{\Sigma}_{j,t})$ , where the posterior mean  
 220 and precision are given by:

$$\hat{E}_{j,t}[f_t] = \hat{\Sigma}_{j,t}(\Sigma^{-1}\mu + \Sigma_{\eta,j,t}^{-1}\eta_{j,t} + \Sigma_{p,t}^{-1}\eta_{p,t}), \quad (12)$$

$$\hat{\Sigma}_{j,t}^{-1} = \Sigma^{-1} + \Sigma_{p,t}^{-1} + \Sigma_{\eta,j,t}^{-1}. \quad (13)$$

221 Next, we solve the model in four steps.

222 *Step 1: Solve for the optimal portfolios, given information sets and issuance.*

---

<sup>6</sup>As our model does not distinguish between size and age, the data availability constraint can also be thought of linking firm age to the expected cost of capital.

223 Substituting the budget constraint (5) into the objective function (4) and taking the  
 224 first-order condition with respect to  $q_{j,t}$  reveals that optimal holdings are increasing in the  
 225 investor's risk tolerance, precision of beliefs, and expected return:

$$q_{j,t}^* = \frac{1}{\rho} \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t). \quad (14)$$

226 *Step 2: Clear the asset market.*

227 Substitute each agent  $j$ 's optimal portfolio (14) into the market-clearing condition (6).  
 228 Collecting terms and simplifying reveals that equilibrium asset prices are linear in payoff risk  
 229 shocks and in supply shocks:

230 **Lemma 1.**  $p_t = \frac{1}{r_t} (A_t + B_t(f_t - \mu) + C_t x_t)$ .

231 A detailed derivation of coefficients  $A_t$ ,  $B_t$ , and  $C_t$ , expected utility, and the proofs of  
 232 this lemma and all further propositions are in the Appendix.

233 In this model, agents learn from prices because prices are informative about the asset  
 234 payoffs  $f_t$ . Next, we deduce what information is implied by Lemma 1. Price information  
 235 is the signal about  $f_t$  contained in prices. The transformation of the price vector  $p_t$  that  
 236 yields an unbiased signal about  $f_t$  is  $\mu + \eta_{p,t} \equiv B_t^{-1}(p_t r_t - A_t)$ . Note that applying the  
 237 formula for  $\eta_{p,t}$  to Lemma 1 reveals that  $\eta_{p,t} = f_t + \varepsilon_{p,t}$ , where the signal noise in prices  
 238 is  $\varepsilon_{p,t} = B_t^{-1} C_t x_t$ . Since we assumed that  $x_t \sim N(0, \sigma_x I)$ , the price noise is distributed  
 239  $\varepsilon_{p,t} \sim N(0, \Sigma_{p,t})$ , where  $\Sigma_{p,t} \equiv \sigma_x B_t^{-1} C_t C_t' B_t^{-1}$ . Substituting in the coefficients  $B_t$  and  $C_t$   
 240 from the proof of Lemma 1 shows that public signal precision  $\Sigma_{p,t}^{-1}$  is a diagonal matrix with  
 241  $i$ th diagonal element  $\sigma_{p,i,t}^{-1} = \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x}$ , where  $\bar{K}_{i,t} \equiv \int K_{i,j,t} dj$  is the average capacity allocated  
 242 to asset  $i$ .

243 This market-clearing asset price reveals the firm's cost of capital. We define the cost of  
 244 capital as follows.

245 **Definition 1.** *The cost of capital for firm  $i$  is the difference between the (unconditional)*  
 246 *expected payout per share the firm will make to investors, minus the (unconditional) expected*  
 247 *price per share that the investor will pay to the firm:  $E_t[f_{i,t}] - E_t[p_{i,t}]$ .*

248 Because  $x_t$  is a mean-zero random variable and the payoff  $f_t$  has mean  $\mu$ , the uncondi-  
 249 tional expected price is  $E_t[p_{i,t}] = A_{i,t}/r$ . Therefore, the expected cost of capital for firm  $i$  is  
 250  $\mu - A_{i,t}/r$ .

251 *Step 3: Compute ex-ante expected utility.*

252 Substitute optimal risky asset holdings from equation (14) into the first-period objective  
 253 function (7) to get:  $U_{j,t} = \rho r_t W_t + \frac{1}{2} E_t[(\hat{E}_{j,t}[f_t] - p_t r_t)' \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t)]$ . Note that the ex-  
 254 pected excess return  $(\hat{E}_{j,t}[f_t] - p_t r_t)$  depends on signals and prices, both of which are unknown  
 255 at the start of the period. Because asset prices are linear functions of normally distributed  
 256 shocks,  $\hat{E}_{j,t}[f_t] - p_t r_t$ , is normally distributed as well. Thus,  $(\hat{E}_{j,t}[f_t] - p_t r_t)' \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t)$   
 257 is a non-central  $\chi^2$ -distributed variable. Computing its mean yields:

$$U_{j,t} = \rho r_t W_t + \frac{1}{2} \text{tr}(\hat{\Sigma}_{j,t}^{-1} V_{j,t} [\hat{E}_{j,t}[f_t] - p_t r_t]) + \frac{1}{2} E_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t]' \hat{\Sigma}_{j,t}^{-1} E_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t]. \quad (15)$$

258 *Step 4: Solve for information choices.*

259 Note that in expected utility (15), the choice variables  $K_{i,j,t}$  enter only through the  
 260 posterior variance  $\hat{\Sigma}_{j,t}$  and through  $V_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t] = V_{j,t}[f - p_t r_t] - \hat{\Sigma}_{j,t}$ . Since there  
 261 is a continuum of investors, and since  $V_{j,t}[f - p_t r_t]$  and  $E_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t]$  depend only on  
 262 parameters and on aggregate information choices, each investor takes them as given.

263 Since  $\hat{\Sigma}_{j,t}^{-1}$  and  $V_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t]$  are both diagonal matrices and  $E_{j,t}[\hat{E}_{j,t}[f_t] - p_t r_t]$  is a  
 264 vector, the last two terms of (15) are weighted sums of the diagonal elements of  $\hat{\Sigma}_{j,t}^{-1}$ . Thus,  
 265 (15) can be rewritten as  $U_{j,t} = r_t W_t + \sum_i \lambda_{i,t} \hat{\Sigma}_{j,t}^{-1}(i, i) - n/2$ , for positive coefficients  $\lambda_{i,t}$ .  
 266 Since  $r_t W_t$  is a constant (in each period  $t$ ) and  $\hat{\Sigma}_{j,t}^{-1}(i, i) = \Sigma^{-1}(i, i) + \Sigma_{p,t}^{-1}(i, i) + K_{i,j,t}$ , the

267 information choice problem is:

$$\max_{K_{1,j,t}, \dots, K_{n,j,t} \geq 0} \sum_{i=1}^n \lambda_{i,t} K_{i,j,t} + \text{constant}, \quad (16)$$

$$\text{s.t.} \quad \sum_{i=1}^n K_{i,j,t} \leq K_t, \quad (17)$$

$$K_{i,j,t} \leq \alpha \exp(\beta x_{i,t-1}) \quad \forall i, \forall j, \quad (18)$$

$$\text{where } \lambda_{i,t} = \bar{\sigma}_{i,t} [1 + (\rho^2 \sigma_x + \bar{K}_{i,t}) \bar{\sigma}_{i,t}] + \rho^2 \bar{x}_{i,t}^2 \bar{\sigma}_{i,t}^2, \quad (19)$$

271 where  $\lambda_{i,t}$  is the marginal value of information  $\bar{\sigma}_{i,t}^{-1} = \int \hat{\Sigma}_{j,t}^{-1}(i, i) dj$  is the average precision  
 272 of posterior beliefs about firm  $i$ . The latter's inverse, average variance  $\bar{\sigma}_{i,t}$ , is decreasing in  
 273  $\bar{K}_{i,t}$ . Equation (19) is derived in the Appendix.

274 This is not a concave objective, so a first-order approach will not find an optimal data  
 275 choice. To maximize a weighted sum (16) subject to an unweighted sum (17), the investor op-  
 276 timally assigns all available data, as per (18), to the asset(s) with the highest weight. If there  
 277 is a unique  $i_t^* = \operatorname{argmax}_i \lambda_{i,t}$ , then the solution is to set  $K_{i_t^*,j,t} = \min(K_t, \alpha \exp(\beta x_{i_t^*,t-1}))$ .

278 In many cases, after all data processing capacity is allocated, there will be multiple assets  
 279 with identical  $\lambda_{i,t}$  weights. That is because  $\lambda_{i,t}$  is decreasing in the average investor's signal  
 280 precision. When there exist asset factor risks  $i, l$  s.t.  $\lambda_{i,t} = \lambda_{l,t}$ , then investors are indifferent  
 281 about which assets' data to process. The next result shows that this indifference is not a  
 282 knife-edge case. It arises whenever the aggregate amount of data processing capacity is  
 283 sufficiently high.

284 **Lemma 2.** *If  $\bar{x}_{i,t}$  is sufficiently large  $\forall i$  and  $\sum_i \sum_j K_{i,j,t} \geq \underline{K}$ , then there exist risks  $l$  and*  
 285  *$l'$  such that  $\lambda_{l,t} = \lambda_{l',t}$ .*

286 This is the big data analog to Grossman and Stiglitz (1980)'s strategic substitutability in  
 287 information acquisition. The more other investors know about an asset, the more informative

288 prices are and the less valuable it is for other investors to process data about the same asset.  
 289 If one asset has the highest marginal utility for signal precision, but capacity is high, then  
 290 many investors will learn about that asset, causing its marginal utility to fall and equalize  
 291 with the next most valuable asset data. With more capacity, the highest two  $\lambda_{i,t}$ 's will be  
 292 driven down until they equate with the next  $\lambda_{i,t}$ , and so forth. This type of equilibrium is  
 293 called a “waterfilling” solution (see, Cover and Thomas (1991)). The equilibrium uniquely  
 294 pins down which assets are being learned about in equilibrium, and how much is learned  
 295 about them, but not which investor learns about which asset.

296 *Step 5: Solve for firm equity issuance.* How a firm chooses  $\bar{x}_t$  depends on how issuance  
 297 affects the asset price. Supply  $\bar{x}_t$  enters the asset price in only one place in the equilibrium  
 298 pricing formula, through  $A_t$ . From Appendix equation (33), we see that

$$A_t = \mu - \rho \bar{\Sigma}_t \bar{x}_t. \quad (20)$$

299  $\bar{x}_t$  has a direct effect on the second term. But also an indirect effect through information  
 300 choices that show up in  $\bar{\Sigma}_t$ .

301 The firm's choice of  $\bar{x}_t$  satisfies its first order condition:

$$E[p_t | \mathcal{I}_{t-1}] - \bar{x}_t \left( \rho \bar{\Sigma}_t - \rho \bar{x}_t \frac{\partial \bar{\Sigma}_t}{\partial \bar{x}_t} \right) - \tilde{\phi}'_1(\bar{x}_t, \bar{x}_{t-1}) = 0. \quad (21)$$

302 The first term is the benefit of more issuance. When a firm issues an additional share,  
 303 it gets expected revenue  $E[p_t | \mathcal{I}_{t-1}]$  for that share. The second term tells us that issuance  
 304 has a positive and negative effect on the share price. The negative effect on the price is that  
 305 more issuance raises the equity premium ( $\rho \bar{\Sigma}_t$  term). The positive price effect is that more  
 306 issuance makes data on the firm more valuable to investors. When investors process more  
 307 data on the firm, it lowers their investment risk, and on average, raises the price they are  
 308 willing to pay ( $\partial \bar{\Sigma}_t / \partial \bar{x}_t$  term). This is the part of the firm investment decision that the rise



309 of big data will affect.

310 The third term, the capital adjustment cost ( $\tilde{\phi}'_1(\bar{x}_t, \bar{x}_{t-1})$ ), reveals why firms grow in size  
 311 over time. Firms have to pay to adjust relative to their initial size. Since firms' starting size  
 312 is small they want to grow, but rapid growth is costly. So, they grow gradually. Each time  
 313 a firm starts larger, choosing a higher  $\bar{x}_t$  becomes less costly because the size of the change,  
 314 and thus the adjustment cost is smaller.

315 Note that in our static model adjustment costs perform a slightly different function  
 316 compared to a dynamic model. In the static model, the main role of adjustment costs is to  
 317 link the initial and final size together, in order to generate cross-sectional differences in the  
 318 marginal value of information. Larger firms can afford to choose a larger final scale. The  
 319 larger the final scale, the higher the marginal value of information.

## 320 2 Parameter Choice

321 In order to quantify the potential effect of big data on firm size, we need to perform a quan-  
 322 titative exercise. What changes exogenously at each date is the total information capacity  
 323  $K_t$ . We normalize  $K_t = 1$  in 1980 and then grow  $K_t$  continuously, at the rate of 36.8% per  
 324 year:  $K_{t+1} = K_t e^{0.368}$ . This rate of growth corresponds to the average rate of growth of CPU  
 325 speed, as illustrated in Figure 2. We simulate the model in this fashion from 1980-2030. As  
 326  $K_t$  increases over time, constraint (9) becomes looser, allowing for a larger overall sum of  
 327 signal precisions.

328 *Figure 2 about here.*

329 We also need to choose values for the model parameters. For  $\mu, \sigma, \sigma_x, r$ , we use the same  
 330 values as in the numerical example in the supplementary appendix to Kacperczyk et al.  
 331 (2016). The next parameter to pick is risk aversion. Risk aversion clearly matters for the

332 level of the risky asset price. But it is tough to identify. The reason for the difficulty is that if  
 333 we change risk aversion and then re-calibrate the mean, persistence and variance parameters  
 334 to match price coefficients and variance at the new risk aversion level, the predictions of  
 335 the model are remarkably stable. Roughly, doubling variance and halving risk aversion  
 336 mostly just redefines units of risk. We set  $\rho = 0.1$ . For firm investment costs, we use  
 337 parameter estimates from Hennessy and Whited (2007). Using annual data from 1988-2001,  
 338 they estimate the cost of external investment funding as  $\Lambda(x) = \phi_0 + \phi_1\tilde{x} + \phi_2\tilde{x}^2$ , where  $\tilde{x}$   
 339 are the proceeds raised from equity flotation. This amount raised corresponds to the change  
 340 in issuance  $\Delta\bar{x}_t$  in our model. Their parameter estimates for  $\phi_0$ ,  $\phi_1$ , and  $\phi_2$  are reported  
 341 in Table 1. The data availability parameters  $\alpha$  and  $\beta$  are chosen to give our mechanism  
 342 a shot at meaningful results. We choose parameters so that the constraint (18) binds for  
 343 small firms only, for about the first decade. This pins both parameters to a narrow range.  
 344 If this constraint did not bind, there would be little difference between small and large firm  
 345 outcomes. If the availability constraint was binding for all firms, then there would be no  
 346 effect of big data growth because there would be insufficient data to process with the growing  
 347 processing power.

348 *Table 1 about here.*

### 349 **3 Quantitative Results**

350 Our main results use the simulated model to understand how the growth of big data affects  
 351 the evolution of large and small firms and how large that effect might be. We start by  
 352 exploring how the rise in big data availability changes how data is allocated. Then, we  
 353 explore how changes in data the investors observe affect the firm's cost of capital. Finally,  
 354 we turn to the question of how much the change in data and the cost of capital affect the  
 355 evolution of firms that start out small and firms that start out large.

356 In presenting our results, we try to balance realism with simplicity, which illuminates the  
 357 mechanism. If we put in a large number of firms, it is, of course, more realistic. But this  
 358 would also make it harder to see what the trade-offs are. Instead, we characterize the firm  
 359 distribution with one representative large firm and one representative small firm. The two  
 360 firms are identical, except that the large firm starts off with a larger size  $\bar{x}_0 = 10,000$ . The  
 361 small firm starts off with  $\bar{x}_0 = 2,000$ . Starting in 1980, we simulate our economy with the  
 362 parameters listed in Table 1 with one period per year until 2030.

### 363 3.1 Data Allocation Choices

364 The reason that data choice is related to firm size in the model is that small firms are  
 365 equivalent to young firms. Young firms do not have a long history of data that can be  
 366 processed.<sup>7</sup> They cannot offer investors the data they need to accurately assess risk and  
 367 return. Data comes from having an observable body of economic transactions. A long  
 368 history with a large amount of economic activity generates this data. In the simulation,  
 369 small firms are those that have more recently entered.

370 *Figure 3 about here.*

371 But the question is, how does the rise of investors' ability to process big data interact  
 372 with this size effect? Since investors are constrained in how much data they can process  
 373 about young, small firms, the increase in data processing ability results in more data being  
 374 processed about the large firm. We can see this in Figure 3 where the share of data processed  
 375 on the large firm rises and the share devoted to the small firm falls (left panel). In the right  
 376 panel, we see that investors are not processing fewer bits of data about the small firm. In  
 377 fact, as the firm grows, little by little, more data is available. As more data is available,  
 378 more small firm data is processed and data precision rises.

---

<sup>7</sup>In the data, small firms are typically younger firms.

379 Eventually, the small firm gets large enough and produces a long enough data history  
 380 that it outgrows its data availability constraint. The availability constraint was pushing data  
 381 choices for the two firms apart, creating the visual bump in Figure 3. As the constraint re-  
 382 laxes, the bump gives way to a slow, steady convergence. But, even once the data availability  
 383 constraint stops binding, investors still process more data on the larger firm. A secondary  
 384 effect of firm size is that data has more value when it is applied to a larger fraction of an  
 385 investor's portfolio. An investor can use a data set to guide his investment of one percent of  
 386 the value of his portfolio. But he gains a lot more when he uses that data to guide investment  
 387 of fifty percent of his portfolio. Big assets constitute more of the value share of the average  
 388 investor's portfolio. Therefore, information about big assets is more valuable.

389 Mathematically, we can see firm size  $\bar{x}_{i,t}$  enter in the marginal value of information  $\lambda_{i,t}$   
 390 in (19). Of course, this firm size is the firm's final size that period. But the final size is linked  
 391 to the firm's initial size through the adjustment cost (2). Firms that are initially larger will  
 392 have a larger final size because size adjustment is costly. This larger final size is what makes  
 393  $\lambda_{i,t}$ , the marginal value of data, higher.

394 In the limit, the small firm keeps growing faster than the large firm and eventually catches  
 395 up. When the two firms approach the same size, the data processing on both converges to  
 396 an equal, but growing amount of data processing.

## 397 **3.2 Capital costs**

398 The main effect of data is to systematically reduce a firm's average cost of capital. Recall that  
 399 the capital cost is the expected payoff minus the expected price of the asset (Definition 1).  
 400 Data does not change the firm's payoff, but it does change how a share of the firm is priced.  
 401 The systematic difference between expected price and payoff is the investor's compensation  
 402 for risk. Investors are compensated for the fact that firm payoffs are unknown, and therefore  
 403 buying a share requires bearing risk. The role of data is to help the investor predict that firm

404 payoff. In doing so, data reduces the compensation for risk. Just like a larger data set lowers  
405 the variance of an econometric estimate, more data in the model reduces the conditional  
406 variance of estimated firm payoffs. An investor who has a more accurate estimate is less  
407 uncertain and bears less risk from holding the asset. The representative investor is willing  
408 to pay more, on average, for a firm that they have good data on. Of course, the data might  
409 reveal problems at the firm that lower the investor's valuation of it. But on average, more  
410 data is neither to reveal positive nor negative news. What data does on average improve is  
411 the precision and resolution of risk. Resolving the investors' risk reduces the compensation  
412 the firm needs to pay the investor for bearing that risk, which reduces the firm's cost of  
413 capital.

414 *Figure 4 about here.*

415 Figure 4 shows how the large firm, with its more abundant data, has a lower cost of  
416 capital. With definition 1 in mind, we can think of the cost of capital as the value per share  
417 delivered to investors. The value per share mechanically depends on the expected payout  
418 per share, which may vary across firm size. In order to compare the cost of capital across  
419 firms of different sizes, we normalize firms' cost of capital with their expected payout per  
420 share.

421 More abundant data does not reduce the cost of capital evenly and proportionately.  
422 There is a second force at play here. The second force is that firm size matters. Because a  
423 firm is large, it represents a larger share of the investor's portfolio. In CAPM-speak, large  
424 firms have a higher beta, and therefore need to offer investors a higher risk compensation.  
425 To induce investors to hold lots of a risk, the compensation per unit of risk must rise. To  
426 induce investors to hold a small amount of a risk is cheap, because small risks wash out in a  
427 large portfolio. Thus, because large firms have more equity outstanding and are more highly  
428 correlated with market risk, a large firm with the same volatility and conditional variance

429 as a small firm, would face a higher cost of capital.

430 As firm size and data evolve together, initially, data dominates. The cost of capital for  
431 the large firm falls, from around 50% of earnings per share to close to 1%, because more  
432 processing power is reducing the risk of investing in that firm. The small firm cannot initially  
433 benefit much from higher processing power because it is a young firm and has little data  
434 available to process. As the small firm grows older, the data availability constraint loosens,  
435 investors can learn from the firm's track record, risk falls and the cost of capital comes back  
436 down. Where the two lines merge is where the small firm finally out-grows its data availability  
437 constraint. From this point on, the only constraint on processing data on either firm is the  
438 total data processing power  $K$ . Large and small firms evolve similarly. The only difference  
439 between the two firms, after the inflection point where the data availability constraint ceases  
440 to bind, is that the small firm continues to have a slightly smaller accumulated stock of  
441 capital. Because the small firm continues to be slightly smaller, it has slightly less equity  
442 outstanding, and a slightly lower cost of capital due to the second force described above.  
443 Once data is abundant, small and large firms converge gradually over time.

### 444 **3.3 The Evolution of Firms' Size**

445 In order to understand how big data has changed the size of firms, it is useful to look at how  
446 a large firm and a small firm evolve in this economy. Then, we turn off various mechanisms  
447 in the model to understand what role is played by each of our key assumptions. Once the  
448 various mechanisms are clear, we contrast firm evolution in the 1980's to the evolution of  
449 firms in the post-2000 period.

450 Recall that firms have to pay to adjust, relative to their previous size. Since firms'  
451 starting size is small, but rapid growth is costly, firms grow gradually. Figure 5 shows that  
452 both the large and small firms grow. However, the rates at which they grow differ. One  
453 reason growth rates differ is that small firms are further from their optimal size. If this were

454 the only force at work, small firms would grow by more each period and that growth rate  
 455 would gradually decline for both firms, as they approach their optimal size.

456 *Figure 5 about here.*

457 Instead, Figure 5 reveals that small firms sometimes grow faster and sometimes slower  
 458 than their large firm counterparts. For much of the start of their life, the small firms grow  
 459 more slowly than the large firms do. These variations in growth rates are due to investors'  
 460 data processing decisions. This is the force that can contribute to the change in the size of  
 461 firms.

462 The level of the size can be interpreted as market capitalization, divided by the expected  
 463 price. Since the average price ranges from 7 to 15 in this model, these are firms with zero  
 464 to 12 million dollars of market value outstanding. In other words, these are not very large  
 465 firms.

466 **The Role of Growing Big Data** Plotting firm outcomes over time as in Figure 5 conflates  
 467 three forces, all changing over time. The first thing changing over time is that firms are  
 468 accumulating capital and growing bigger. The second change is that firms are accumulating  
 469 longer data histories, which makes more data for processing available. The third change is  
 470 that technology enables investors to process more and more of that data over time. We want  
 471 to understand how each of these contributes to our main results. Therefore, we turn off  
 472 features of the model one-by-one, and compare the new results to the main results, in order  
 473 to understand what role each of these ingredients plays.

474 *Figure 6 about here.*

475 To understand the role that improvements in data processing play, we turn off the growth  
 476 of big data and compare results. We fix  $K_t = 5 \forall t$ . Data processing capacity is frozen at

477 its 1985 level. Firms still have limited data histories and still accumulate capital. Figure 6  
478 shows that this small change in the model has substantial consequences for firm dynamics.  
479 Comparing Figures 5 and 6, we can see the role big data plays. In the world with fixed  
480 data processing, instead of starting with rapid growth and growing faster as data processing  
481 improves, the large firm growth rate starts at the same level as before, but then steadily  
482 declines as the firm approaches its stationary optimal size. We learn that improvements in  
483 data processing are the sources of firm growth in the model and are central to the continued  
484 rapid growth of large firms.

485 **The Role of Limited Data History** One might wonder, if large firms attract more data  
486 processing, is that alone producing larger big firms? Is the assumption that small firms have  
487 a limited data history really important for the results? To answer this question, we now turn  
488 off the assumption of limited data history. We maintain the growing data capacity and firm  
489 capital accumulation from the original model.

490 *Figure 7 about here.*

491 Figure 7 reports results for the model with unlimited firm data histories, but limited  
492 processing power, to the full model. Comparing Figures 5 and 7, we can see the difference  
493 that data availability makes. In the world where firms have unlimited data histories, small  
494 firms quickly catch up to large firms. There is no persistent difference in size. Small firms are  
495 far below their optimal size. So they invest rapidly. Investment makes them larger, which  
496 increases data processing immediately. Quickly, the initially small and large firms become  
497 indistinguishable. Adjustment costs are a friction preventing immediate convergence. But it  
498 is really the presence of the data availability constraint that creates the persistent difference  
499 between firms with different initial size.



500 **Small and Large Firms in the New Millennium** So far, the experiment has been  
501 to drop a small firm and a large firm in the economy in 1980 and watch how they evolve.  
502 While this is useful to explaining the model's main mechanism, it does not really answer the  
503 question of why small firms today struggle more than in the past and why large firms today  
504 are larger than the large firms of the past. To answer these questions, we really want to  
505 compare small and large firms that enter the economy today to small and large firms that  
506 entered in 1980.

507 To do this small vs. large, today vs. 1980 experiment, we use the same parameters as in  
508 Table 1 and use the same starting size for firms. The only difference is that we start with  
509 more available processing power. Instead of starting  $K_t$  at 1, we start it at the 2000 value,  
510 which is about 527.

511 *Figure 8 about here.*

512 Each panel of figure 8 shows the growth rate of a large firm, minus the growth rate of a  
513 small firm. In the left panel, both firms start in 1980, when data processing capacity was  
514 quite limited. In the right panel, both firms enter in the year 2000, when data is abundant.  
515 In both cases, the difference is positive for most of the first decade, meaning that large firms  
516 grow faster than small ones. But in 2000, the difference is much more positive. Relative  
517 to small firms, large firms grow much more quickly. The difference in 2000 growth rates  
518 is nearly twice as large. In both cases, a surviving small firm eventually outgrows its data  
519 availability problem, grows quickly, and then converges to the growth rate of the large firm  
520 (differences converge to 0).

521 In a model with random shocks and exit, many small firms would not survive. Of course,  
522 for some firms, the possibility of future growth would induce them to hang on, preventing  
523 exit. In a world where large firms gain market share much more rapidly, firms would either  
524 exit, unable to compete, or strive to quickly grow large. This illustrates how data pro-

525 censing advances may contribute to the puzzle of missing small firms, by disproportionately  
526 benefiting large firms.

527 For comparison, we examine the growth rates of large and small firms in the U.S. pre-  
528 1980 and in the period 1980-2007. We end in 2007 so as to avoid measuring real effects of  
529 the financial crisis. For each industry sector and year, we select the top 25% largest firms in  
530 Compustat and call those large firms and select the bottom half of the firm size distribution  
531 to be our small firms. Within these two sets of firms, we compute the growth rates of various  
532 measures of firm size and average them, with an equal weight given to each firm. Then, just  
533 as in the model, we subtract the growth rate of large firms from that of small firms. For  
534 most measures, small firms grow more slowly, and that difference grows later in the sample.

535 *Table 2 about here.*

536 At times, the magnitudes of the model's growth rates are quite large, compared to the  
537 data. Of course, the data is averaged over many years and many firms at different points in  
538 their life cycle. This smooths out some of the extremes in the data. If we average the firm  
539 growth in our model from 1980-1985, for firms that enter in 1980, we get 39.5% for large  
540 firms and 7.3% for small firms, a difference of 32.2%. If we average firm growth in our model  
541 from 2000-2005, for firms that enter in 2000, we get 59.9% for large firms and 7.3% for small  
542 firms, a difference of 52.6%.

543 While it is not unheard of for a small firm to double in size, some of this magnitude  
544 undoubtedly reflects some imprecision of our current numerical example vis-à-vis the data.  
545 A larger adjustment cost, or a labor hiring delay, would help to moderate the extremes of firm  
546 size growth. The results also miss many aspects of the firm environment that have changed  
547 in the last four decades. The type of firms entering in the last few years are quite different  
548 than firms of prior years. They have different sources of revenue and assets that might be  
549 harder to value. Firm financing has changed, with a shift toward internal financing. Venture

550 capital funding has become more prevalent and displaced equity funding for many firms,  
551 early in their life cycle. All of these forces would moderate the large effect we document  
552 here.

553 Our results only show that big data is a force with some potential. There is a logical way  
554 in which the growth of big data and the growth of large firms is connected. This channel  
555 has the potential to be quantitatively powerful. The role of big data in firms is thus a topic  
556 ripe for further exploration.

## 557 4 Discussion

558 In this model, there is a one-to-one correspondence between projects and firms. Investors  
559 gather information about the firm and smaller firms have smaller amounts of information.  
560 This then feeds back into real investments of each firm in its single project, and determines  
561 the firm/project size distribution. As such, information processing and big data helps a  
562 small firm less, even if it is investing in a well known technology (for instance, the  $n^{th}$  firm  
563 to drill an oil well). We think this is a reasonable assumptions for publicly listed firms, since  
564 information about the firm is both about its track record as well as the quality of its project.  
565 The investors have a harder time accessing the survival probability of a firm with no track  
566 record relative to a well established firm in a highly competitive industry, which is why we  
567 find our information assumption relevant even in such settings.

568 We should note that our model is not best suited to speak to firm entry. For instance, a  
569 new class of online firms have emerged who use big data to facilitate capital markets' access  
570 to an under-served segment of population, such as personal loans to people with very low  
571 credit scores. Such firms are often small, but they have only emerged as a by product of big  
572 data availability. This trend is fascinating in its own right, yet is outside the scope of our  
573 paper.

574 In the context of the model, firms are equity financed. This implies that their real  
575 investment and thus their cash flow is determined by firms' cost of accessing external capital  
576 markets. Financing is costly for firms since investors require an equity premium to hold  
577 firms' risky shares. However, since more data is available about large firms, big data reduces  
578 the asymmetric information friction relatively more for big firms compared to small firms.  
579 Cheaper access to external capital markets reduces large firms' cost of capital and accelerates  
580 their growth. On the other hand, small firms growth is initially stagnant. However, once  
581 they become sufficiently large, their access to capital markets improve as well, and their  
582 growth rate picks up. This is consistent with information asymmetries being a short-horizon  
583 notion.

## 584 **5 Conclusion**

585 Big data is transforming the modern economy. While many economists have used big data,  
586 fewer think about how the use of data by others affects market outcomes. This paper starts to  
587 explore the ways in which big data might be incorporated in modern economic and financial  
588 theory. One way that big data is used is to help financial market participants make more  
589 informed choices about the firms in which they invest. These investment choices affect the  
590 prices, cost of capital, and investment decisions of these firms. We set up a very simple  
591 model to show how such big data choices might be incorporated and one way in which the  
592 growth of big data might affect the real economy. But this is only a modest first step.

593 One might also consider how firms themselves use data, to refine their products, to  
594 broaden their customer market, or to increase the efficiency of their operations. Such data,  
595 produced as a by-product of economic activity, might also favor the large firms whose abun-  
596 dant economic activity produces abundant data.

597 Another step in a big-data agenda would be to consider the sale of data. In many

598 information models, we think of signals that are observed and then embedded in one's  
599 knowledge, not easily or credibly transferable. However data is an asset that can be bought,  
600 sold and priced on a market. How do markets for data change firms choices, investments,  
601 evolution and their valuations as firms? It is true that data intermediaries like Foursquare  
602 or Amazon help small businesses benefit from each others' data. At the same time, these  
603 intermediaries retain control of the data and extract rents from firms that use it. A firm  
604 that has its own customer data clearly has a real advantage. Whether an intermediary can  
605 find a way for small firms to collectively leverage their data, in a way that mimics a large  
606 firm advantage, remains to be seen.

607 Finally, if data is a storable, sellable, priced asset, then investment in data should be  
608 valued just as if it were investment in a physical asset. Understanding how to price data as  
609 an asset might help us to better understand the valuations of new-economy firms and better  
610 measure aggregate economic activity.

## 611 References

- 612 Akcigit, U., Kerr, W., 2017. Growth through heterogeneous innovations. *Journal of Political Econ-*  
613 *omy* forthcoming.
- 614 Asriyan, V., Vanasco, V., 2014. Informed intermediation over the cycle, Stanford Working Paper.
- 615 Bai, J., Philippon, T., Savov, A., 2016. Have financial markets become more informative? *Journal*  
616 *of Financial Economics* 122 (35), 625–654.
- 617 Begenau, J., Salomao, J., 2018. Firm financing over the business cycle. *Review of Financial Studies*  
618 forthcoming.
- 619 Bernhardt, D., Hollifield, B., Hughson, E., 1995. Investment and insider trading. *The Review of*  
620 *Financial Studies* 8 (2), 501–543.
- 621 Biais, B., Foucault, T., Moinas, S., 2015. Equilibrium fast trading. *Journal of Financial Economics*  
622 116, 292–313.
- 623 Cooley, T. F., Quadrini, V., December 2001. Financial markets and firm dynamics. *American*  
624 *Economic Review* 91 (5), 1286–1310.  
625 URL <http://ideas.repec.org/a/aea/aecrev/v91y2001i5p1286-1310.html>
- 626 Cover, T., Thomas, J., 1991. *Elements of information theory*, 1st Edition. John Wiley and Sons,  
627 New York, New York.
- 628 Davis, S. J., Haltiwanger, J., 2015. Dynamism diminished: The role of credit conditions. in progress.
- 629 Fang, L., Peress, J., 2009. Media coverage and the cross-section of stock returns. *The Journal of*  
630 *Finance* 64 (5), 2023–2052.  
631 URL <http://www.jstor.org/stable/27735164>
- 632 Farboodi, M., Matray, A., Veldkamp, L., 2017. Where has all the big data gone?, Working Paper,  
633 Princeton University.

- 634 Glode, V., Green, R., Lowery, R., 2012. Financial expertise as an arms race. *Journal of Finance* 67,  
635 1723–1759.
- 636 Gomes, J. F., 2001. Financing investment. *The American Economic Review* 91 (5), pp. 1263–1285.  
637 URL <http://www.jstor.org/stable/2677925>
- 638 Grossman, S., Stiglitz, J., 1980. On the impossibility of informationally efficient markets. *American*  
639 *Economic Review* 70(3), 393–408.
- 640 Hennessy, C. A., Whited, T. M., 2007. How costly is external financing? evidence from a structural  
641 estimation. *The Journal of Finance* LXII (4).
- 642 Hennessy, J., Patterson, D., 2011. *Computer Architecture*. Elsevier.
- 643 Kacperczyk, M., Nosal, J., Stevens, L., 2015. Investor sophistication and capital income inequality,  
644 Imperial College Working Paper.
- 645 Kacperczyk, M., Van Nieuwerburgh, S., Veldkamp, L., 2016. A rational theory of mutual funds’  
646 attention allocation. *Econometrica* 84(2), 571–626.
- 647 Kozeniauskas, N., 2017. Technical change and declining entrepreneurship, Working Paper, New  
648 York University.
- 649 Maksimovic, V., Stomper, A., Zechner, J., 1999. Capital structure, information acquisition and  
650 investment decisions in an industry framework. *Review of Finance* 2 (3), 251–271.
- 651 Ozdenoren, E., Yuan, K., 2008. Feedback effects and asset prices. *The journal of finance* 63 (4),  
652 1939–1975.
- 653 Van Nieuwerburgh, S., Veldkamp, L., 2010. Information acquisition and under-diversification. *Re-*  
654 *view of Economic Studies* 77 (2), 779–805.
- 655 Veldkamp, L., 2011. *Information choice in macroeconomics and finance*. Princeton University Press.

## 656 A Proofs

### 657 A.1 Useful notation, matrices and derivatives

658 All the following matrices are diagonal with  $ii$  entry given by:

- 659 1. Average signal precision:  $(\bar{\Sigma}_{\eta,t}^{-1})_{ii} = \bar{K}_{i,t}$ , where  $\bar{K}_{i,t} \equiv \int K_{i,j,t} dj$ .
- 660 2. Precision of the information prices convey about shock  $i$ :  $(\Sigma_{p,t}^{-1})_{ii} = \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} = \sigma_{i,p,t}^{-1}$
- 661 3. Precision of posterior belief about shock  $i$  for an investor  $j$  is  $\hat{\sigma}_{i,j,t}^{-1}$ , which is equivalent to

$$(\hat{\Sigma}_{j,t}^{-1})_{ii} = (\Sigma^{-1} + \Sigma_{\eta,j,t}^{-1} + \Sigma_{p,t}^{-1})_{ii} = \sigma_i^{-1} + K_{i,j,t} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} = \hat{\sigma}_{i,j,t}^{-1} \quad (22)$$

- 662 4. Average posterior precision of shock  $i$ :  $\bar{\sigma}_{i,t}^{-1} \equiv \sigma_i^{-1} + \bar{K}_{i,t} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x}$ . The average variance is therefore  
663  $(\bar{\Sigma}_t)_{ii} = \left[ \sigma_i^{-1} + \bar{K}_{i,t} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} \right]^{-1} = \bar{\sigma}_{i,t}$ .
- 664 5. Ex-ante mean and variance of returns: Using Lemma 1 and the coefficients given by (33), (34) and  
665 (35), we can write the return as:

$$\begin{aligned} f_t - p_t r_t &= (I - B_t)(f_t - \mu) - C_t x_t + \rho \bar{\Sigma}_t \bar{x}_t \\ &= \bar{\Sigma}_t \left[ \Sigma^{-1}(f_t - \mu) + \rho \left( I + \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \right) x_t \right] + \rho \bar{\Sigma}_t \bar{x}_t. \end{aligned}$$

666 This expression is a constant plus a linear combination of two normal variables, which is also a normal  
667 variable. Therefore, we can write

$$f_t - p_t r_t = V_t^{1/2} u_t + w_t, \quad (23)$$

668 where  $u_t$  is a standard normally distributed random variable  $u_t \sim N(0, I)$ , and  $w_t$  is a non-random  
669 vector measuring the ex-ante mean of excess returns

$$w_t \equiv \rho \bar{\Sigma}_t \bar{x}_t. \quad (24)$$

670 and  $V_t$  is the ex-ante variance matrix of excess returns:

$$\begin{aligned} V_t &\equiv \bar{\Sigma}_t \left[ \Sigma^{-1} + \rho^2 \sigma_x \left( I + \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \right) \left( I + \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \right)' \right] \bar{\Sigma}_t \\ &= \bar{\Sigma}_t \left[ \Sigma^{-1} + \rho^2 \sigma_x \left( I + \frac{1}{\rho^2 \sigma_x} (\bar{\Sigma}_{\eta,t}^{-1'} + \bar{\Sigma}_{\eta,t}^{-1}) + \frac{1}{\rho^4 \sigma_x^2} \bar{\Sigma}_{\eta,t}^{-1'} \bar{\Sigma}_{\eta,t}^{-1} \right) \right] \bar{\Sigma}_t \\ &= \bar{\Sigma}_t \left[ \Sigma^{-1} + \rho^2 \sigma_x I + (\bar{\Sigma}_{\eta,t}^{-1'} + \bar{\Sigma}_{\eta,t}^{-1}) + \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \bar{\Sigma}_{\eta,t}^{-1} \right] \bar{\Sigma}_t \\ &= \bar{\Sigma}_t \left[ \rho^2 \sigma_x I + \bar{\Sigma}_{\eta,t}^{-1'} + \Sigma^{-1} + \bar{\Sigma}_{\eta,t}^{-1} + \Sigma_{p,t}^{-1} \right] \bar{\Sigma}_t \\ &= \bar{\Sigma}_t \left[ \rho^2 \sigma_x I + \bar{\Sigma}_{\eta,t}^{-1'} + \bar{\Sigma}_t^{-1} \right] \bar{\Sigma}_t. \end{aligned}$$

671 The first line uses  $E[x_t x_t'] = \sigma_x I$  and  $E[(f_t - \mu)(f_t - \mu)'] = \Sigma$ , the fourth line uses (36) and the fifth  
672 line uses  $\bar{\Sigma}_t^{-1} = \Sigma^{-1} + \Sigma_{p,t}^{-1} + \bar{\Sigma}_{\eta,t}^{-1}$ .



673 This variance matrix  $V_t$  is a diagonal matrix. Its diagonal elements are:

$$\begin{aligned} (V_t)_{ii} &= (\bar{\Sigma}_t [\rho^2 \sigma_x I + \bar{\Sigma}_{\eta,t}^{-1} + \bar{\Sigma}_t^{-1}] \bar{\Sigma}_t)_{ii} \\ &= \bar{\sigma}_{i,t} [1 + (\rho^2 \sigma_x + \bar{K}_{i,t}) \bar{\sigma}_{i,t}]. \end{aligned} \quad (25)$$

## 674 A.2 Solving the Model

675 **Step 1: Portfolio Choices** From the FOC, the optimal portfolio is chosen by investor  $j$  is

$$q_{j,t}^* = \frac{1}{\rho} \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t). \quad (26)$$

676 where  $\hat{E}_{j,t}[f_t]$  and  $\hat{\Sigma}_{j,t}$  depend on the skill of the investor.

677 Next, we compute the portfolio of the average investor.

$$\begin{aligned} \bar{q} \equiv \int q_{j,t}^* dj &= \frac{1}{\rho} \int \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t) dj \\ &= \frac{1}{\rho} \left( \int \Sigma_{\eta,j,t}^{-1} \eta_{j,t} dj + \Sigma_{p,t}^{-1} \eta_{p,t} + \bar{\Sigma}_t^{-1} (\mu - p_t r_t) \right) \\ &= \frac{1}{\rho} (\bar{\Sigma}_{\eta,t}^{-1} f_t + \Sigma_{p,t}^{-1} \eta_{p,t} + \bar{\Sigma}_t^{-1} (\mu - p_t r_t)), \end{aligned} \quad (27)$$

678 where the fourth equality uses the fact that average noise of private signals is zero.

679 **Step 2: Clearing the asset market and computing expected excess return** Lemma  
680 1 describes the solution to the market-clearing problem and derives the coefficients  $A_t$ ,  $B_t$ , and  $C_t$  in the  
681 pricing equation. The equilibrium price, along with the random signal realizations determines the interim  
682 expected return  $(\hat{E}_{j,t}[f_t] - p_t r_t)$ . But at the start of the period, the equilibrium price and one's realized  
683 signals are not known. To compute beginning-of-period utility, we need to know the ex-ante expectation and  
684 variance of this interim expected return.

685 The interim expected excess return can be written as:  $\hat{E}_{j,t}[f_t] - p_t r_t = \hat{E}_{j,t}[f_t] - f_t + f_t - p_t r_t$  and  
686 therefore its variance is:

$$V_t[\hat{E}_{j,t}[f_t] - p_t r_t] = V_t[\hat{E}_{j,t}[f_t] - f_t] + V_t[f_t - p_t r_t] + 2 \text{Cov}_t[\hat{E}_{j,t}[f_t] - f_t, f_t - p_t r_t]. \quad (28)$$

687 Combining (12) with the definitions  $\eta_{j,t} = f_t + \varepsilon_{j,t}$  and  $\eta_{p,t} = f_t + \varepsilon_{p,t}$ , we can compute expectation errors:

$$\begin{aligned} \hat{E}_{j,t}[f_t] - f_t &= \hat{\Sigma}_{j,t} [(\Sigma^{-1} \mu + (\Sigma_{\eta,j,t}^{-1} + \Sigma_{p,t}^{-1}) f_t + \Sigma_{\eta,j,t}^{-1} \varepsilon_{j,t} + \Sigma_{p,t}^{-1} \varepsilon_{p,t}) - f_t] \\ &= \hat{\Sigma}_{j,t} [-\Sigma^{-1} (f_t - \mu) + \Sigma_{\eta,j,t}^{-1} \varepsilon_{j,t} + \Sigma_{p,t}^{-1} \varepsilon_{p,t}] \end{aligned}$$

688 Computing the expectation, we obtain  $E_t[\hat{E}_{j,t}[f_t] - f_t] = \hat{\Sigma}_{j,t} \hat{\Sigma}_{j,t}^{-1} \mu - \mu = 0$  and its variance is  $V_t[\hat{E}_{j,t}[f_t] -$   
689  $f_t] = \hat{\Sigma}_{j,t} [\Sigma^{-1} + \Sigma_{\eta,j,t}^{-1} + \Sigma_{p,t}^{-1}] \hat{\Sigma}_{j,t} = \hat{\Sigma}_{j,t}$ .

690 From (23) we know that  $V_t[f_t - p_t r_t] = V_t$ . To compute the covariance term, we can rearrange the  
691 definition of  $\eta_{p,t}$  to get  $p_t r_t = B_t \eta_{p,t} + A_t - B_t \mu$  and  $\eta_{p,t} = f_t + \varepsilon_{p,t}$  to write

$$f_t - p_t r_t = (I - B_t) f_t - A_t - B_t \varepsilon_{p,t} + B_t \mu \quad (29)$$

$$= \rho \bar{\Sigma}_t \bar{x}_t + \bar{\Sigma}_t \Sigma^{-1} (f_t - \mu) - (I - \bar{\Sigma}_t \Sigma^{-1}) \varepsilon_{p,t} \quad (30)$$

692 where the second line comes from substituting the coefficients  $A_t$  and  $B_t$  from Lemma 1. Since the constant  
693  $\rho \bar{\Sigma}_t \bar{x}_t$  does not affect the covariance, we can write

$$\begin{aligned} \text{Cov}_t[\hat{E}_{j,t}[f_t] - f_t, f_t - p_t r_t] &= \text{Cov}[-\hat{\Sigma}_{j,t} \Sigma^{-1}(f_t - \mu) + \hat{\Sigma}_{j,t} \Sigma_{p,t}^{-1} \varepsilon_{p,t}, \bar{\Sigma}_t \Sigma^{-1}(f_t - \mu) - (I - \bar{\Sigma}_t \Sigma^{-1}) \varepsilon_{p,t}] \\ &= -\hat{\Sigma}_{j,t} \Sigma^{-1} \Sigma \Sigma^{-1} \bar{\Sigma}_t - \hat{\Sigma}_{j,t} \Sigma_{p,t}^{-1} \Sigma_{p,t} (I - \Sigma^{-1} \bar{\Sigma}_t) \\ &= -\hat{\Sigma}_{j,t} \Sigma^{-1} \bar{\Sigma}_t - \hat{\Sigma}_{j,t} (I - \Sigma^{-1} \bar{\Sigma}_t) = -\hat{\Sigma}_{j,t} \end{aligned}$$

Substituting the three variance and covariance terms into (28), we find that the variance of excess return is  $V_t[\hat{E}_{j,t}[f_t] - p_t r_t] = \hat{\Sigma}_{j,t} + V_t - 2\hat{\Sigma}_{j,t} = V_t - \hat{\Sigma}_{j,t}$ . Note that this is a diagonal matrix. Substituting the expressions (25) and (22) for the diagonal elements of  $V_t$  and  $\hat{\Sigma}_{j,t}$  we have

$$(V_t[\hat{E}_{j,t}[f_t] - p_t r_t])_{ii} = (V_t - \hat{\Sigma}_{j,t})_{ii} = (\bar{\sigma}_{i,t} - \hat{\sigma}_{i,j,t}) + (\rho^2 \sigma_x + \bar{K}_{i,t}) \bar{\sigma}_{i,t}^2$$

694 In summary, the excess return is normally distributed as  $\hat{E}_{j,t}[f_t] - p_t r_t \sim \mathcal{N}(w_t, V_t - \hat{\Sigma}_{j,t})$ .

695 **Step 3: Compute ex-ante expected utility** Ex-ante expected utility for investor  $j$  is  $U_{j,t} =$   
696  $E_t[\rho \hat{E}_{j,t}[\hat{W}_{j,t}] - \frac{\rho^2}{2} \hat{V}_{j,t}[\hat{W}_{j,t}]]$ . In period 2, the investor has chosen his portfolio and the price is in his infor-  
697 mation set, therefore the only payoff-relevant, random variable is  $f_t$ . We substitute the budget constraint in  
698 the optimal portfolio choice from (26) and take expectation and variance conditioning on  $\hat{E}_{j,t}[f_t]$  and  $\hat{\Sigma}_{j,t}$   
699 to obtain  $U_{j,t} = \rho r_t W_t + \frac{1}{2} E_t[(\hat{E}_{j,t}[f_t] - p_t r_t)' \hat{\Sigma}_{j,t}^{-1} (\hat{E}_{j,t}[f_t] - p_t r_t)]$ .

700 Define  $m_t \equiv \hat{\Sigma}_{j,t}^{-1/2} (\hat{E}_{j,t}[f_t] - p_t r_t)$  and note that  $m_t \sim \mathcal{N}(\hat{\Sigma}_{j,t}^{-1/2} w_t, \hat{\Sigma}_{j,t}^{-1/2} V_t \hat{\Sigma}_{j,t}^{-1/2} - I)$ . The second  
701 term in the  $U_{i,t}$  is equal to  $E[m_t' m_t]$ , which is the mean of a non-central Chi-square. Using the formula, if  
702  $m_t \sim N(E[m_t], \text{Var}[m_t])$ , then  $E[m_t' m_t] = \text{tr}(\text{Var}[m_t]) + E[m_t]' E[m_t]$ , we get

$$U_{j,t} = \rho r_t W_t + \frac{1}{2} \text{tr}(\hat{\Sigma}_{j,t}^{-1/2} V \hat{\Sigma}_{j,t}^{-1/2} - I) + \frac{1}{2} w_t' \hat{\Sigma}_{j,t}^{-1} w_t = \rho r_t W_t + \frac{1}{2} \text{tr}(\hat{\Sigma}_{j,t}^{-1} V) - \text{tr}(I) + \frac{1}{2} w_t' \hat{\Sigma}_{j,t}^{-1} w_t.$$

703 Finally, we substitute the expressions for  $\hat{\Sigma}_{j,t}^{-1}$  and  $w_t$  from (22) and (24):

$$\begin{aligned} U_{j,t} &= \rho r_t W_t - \frac{N}{2} + \frac{1}{2} \sum_{i=1}^N \left( \sigma_i^{-1} + K_{i,j,t} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} \right) (V_t)_{ii} + \frac{\rho^2}{2} \sum_{i=1}^N \bar{x}_{i,t}^2 \bar{\sigma}_{i,t}^2 \left( \sigma_i^{-1} + K_{i,j,t} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} \right) \\ &= \frac{1}{2} \sum_{i=1}^N K_{i,j,t} [(V_t)_{ii} + \rho^2 \bar{x}_{i,t}^2 \bar{\sigma}_{i,t}^2] + \rho r_t W_t - \frac{N}{2} + \frac{1}{2} \sum_{i=1}^N \left( \sigma_i^{-1} + \frac{\bar{K}_{i,t}^2}{\rho^2 \sigma_x} \right) [(V_t)_{ii} + \rho^2 \bar{x}_{i,t}^2 \bar{\sigma}_{i,t}^2] \\ &= \frac{1}{2} \sum_{i=1}^N K_{i,j,t} \lambda_{i,t} + \text{constant} \end{aligned} \quad (31)$$

$$704 \lambda_{i,t} = \bar{\sigma}_{i,t} [1 + (\rho^2 \sigma_x + \bar{K}_{i,t}) \bar{\sigma}_{i,t}] + \rho^2 \bar{x}_{i,t}^2 \bar{\sigma}_{i,t}^2 \quad (32)$$

705 where the weights  $\lambda_{i,t}$  are given by the variance of expected excess return  $(V_t)_{ii}$  from (25) plus a term that  
706 depends on the supply of the risk.

707 **Step 4: Information choices** The attention allocation problem maximizes ex-ante utility in (31)  
708 subject to the information capacity, data availability and no-forgetting constraints (17), (18) and (11).  
709 Observe that  $\lambda_{i,t}$  depends only on parameters and on aggregate average precisions. Since each investor  
710 has zero mass within a continuum of investors, he takes  $\lambda_{i,t}$  as given. Since the constant is irrelevant,  
711 the optimal choice maximizes a weighted sum of attention allocations, where the weights are given by  $\lambda_{i,t}$

(equation (19)), subject to a constraint on an un-weighted sum. This is not a concave objective, so a first-order approach will not deliver a solution. A simple variational argument reveals that allocating all capacity to the risk(s) with the highest  $\lambda_{i,t}$  achieves the maximum utility. For a formal proof of this result, see Van Nieuwerburgh and Veldkamp (2010). Thus, the solution is given by:  $K_{i,j,t} = K_t$  if  $\lambda_{i,t} = \max_k \lambda_{k,t}$ , and  $K_{i,j,t} = 0$ , otherwise. There may be multiple risks  $i$  that achieve the same maximum value of  $\lambda_{i,t}$ . In that case, the manager is indifferent about how to allocate attention between those risks. We focus on symmetric equilibria.

### 719 A.3 Proofs

#### 720 Proof of Lemma 1

721 *Proof.* Following Admati (1985), we know that the equilibrium price takes the following form  $p_t r_t =$   
722  $A_t + B_t(f_t - \mu) + C_t x_t$  where

$$A_t = \mu - \rho \bar{\Sigma}_t \bar{x}_t \quad (33)$$

$$B_t = I - \bar{\Sigma}_t \Sigma^{-1} \quad (34)$$

$$C_t = -\rho \bar{\Sigma}_t \left( I + \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \right) \quad (35)$$

723 and therefore the price is given by  $p_t r_t = \mu + \bar{\Sigma}_t \left[ (\bar{\Sigma}_t^{-1} - \Sigma^{-1})(f_t - \mu) - \rho(\bar{x}_t + x_t) - \frac{1}{\rho \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} x_t \right]$ . Further-  
724 more, the precision of the public signal is

$$\Sigma_{p,t}^{-1} \equiv \left( \sigma_x B_t^{-1} C_t C_t' B_t^{-1'} \right)^{-1} = \frac{1}{\rho^2 \sigma_x} \bar{\Sigma}_{\eta,t}^{-1'} \bar{\Sigma}_{\eta,t}^{-1} \quad (36)$$

725 □

726 **Proof of Lemma 2** See Kacperczyk et. al (2016).

### 727 A.4 Firm Volatility Data

728 The introduction of our paper claims that differential trends in the volatility of large and small firms' earnings  
729 is not a plausible explanation for the different trends in the cost of capital. To support this claim, we explore  
730 whether the volatility of large and small firms has diverged. We find some fluctuations, but no consistent  
731 trend in the difference.

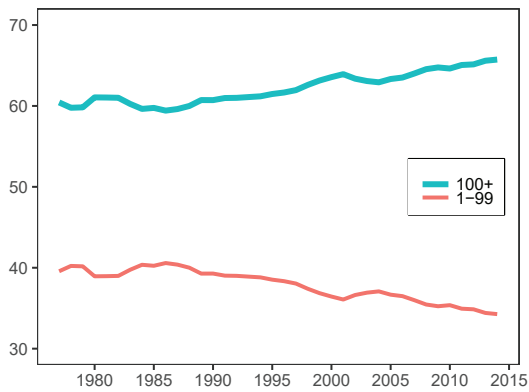
732 Our volatility measure is based on the annual growth rate in earnings calculated at the firm level from  
733 quarterly CRSP/Computat data from 1960 - 2016. Earnings are constructed by multiplying basic earnings  
734 per share, excluding extraordinary items (EPS) by the number of shares used to calculate EPS. We measure  
735 the volatility of earnings growth as the rolling standard deviation over the past 20 quarters. The firms are  
736 split by size in a number of ways: firstly, the firms in the sample are split by whether or not they are (at that  
737 time) a member of the S&P500 index. Secondly, we split firms by whether or not they were in the top half  
738 of the earnings distribution in each quarter. Lastly, we consider only firms in the bottom and top quartiles  
739 of the earnings distribution. In the plots below, the dashed lines are the median volatility, whilst the solid  
740 line is the trend extracted from this series using a HP filter, with  $\lambda = 1600$ .

741 *Figure 9 about here.*

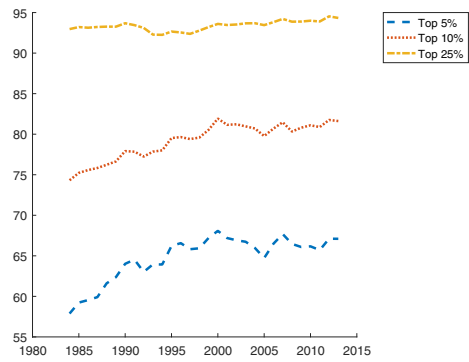
742 Figure 9 plots volatility for large and small firms, over time. Whether the firms are cut at the median,  
743 the top and bottom quartiles, or by membership in the S&P 500, in every case, there are fluctuations in

744 volatility, and there are long-run increases in volatility. But there is no consistent long-run trend in the gap  
745 between the different sized firms.

**Figure 1: Large Firms Growing Relatively Larger.** The left panel uses the Business Dynamics Statistics data published by the Census Bureau (from Kozeniauskas, 2017). It contains all firms with employees in the private non-farm sector in the United States. The right panel uses Compustat/CRSP data. Top  $x\%$  means the share of all firm revenue earned by the  $x\%$  highest-revenue firms.

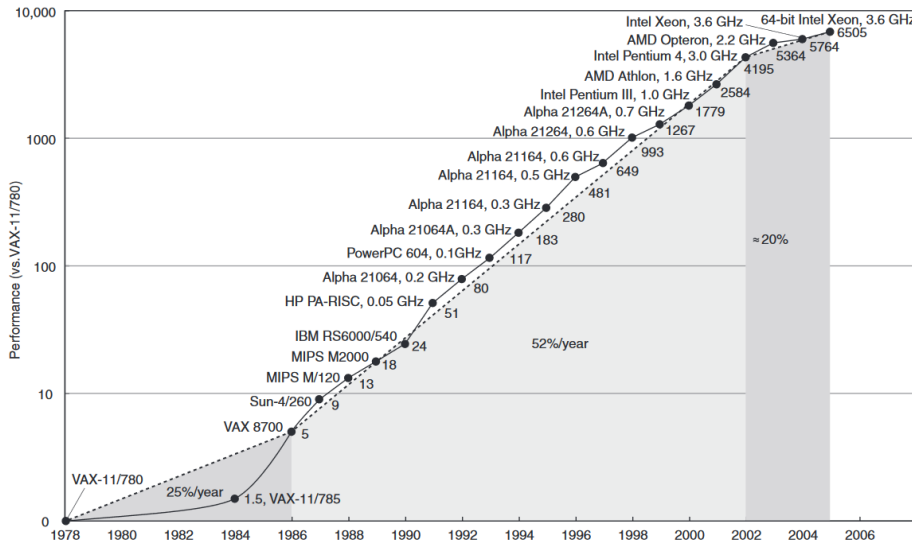


(a) Employment % of large/small firms



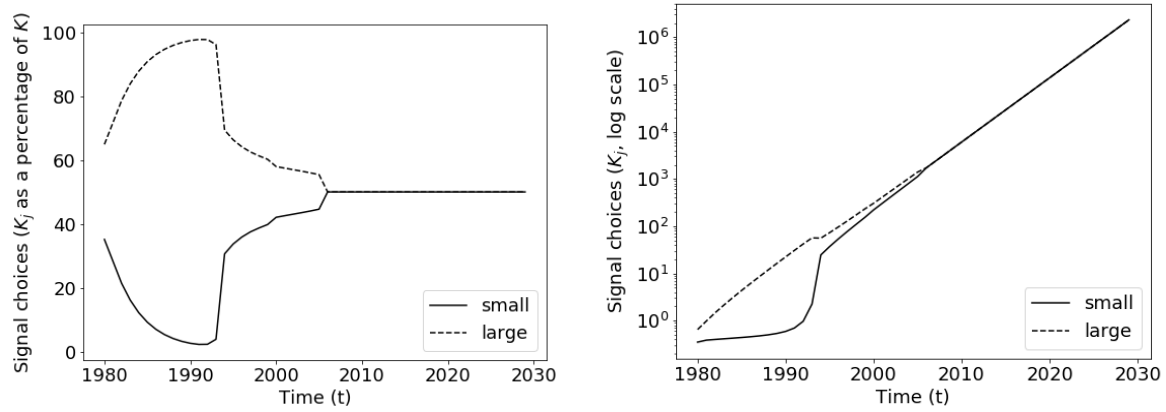
(b) Revenue % at firms in top sales percentile

**Figure 2:** The evolution of processing performance over the period 1978–2007  
Hennessy and Patterson (2011)

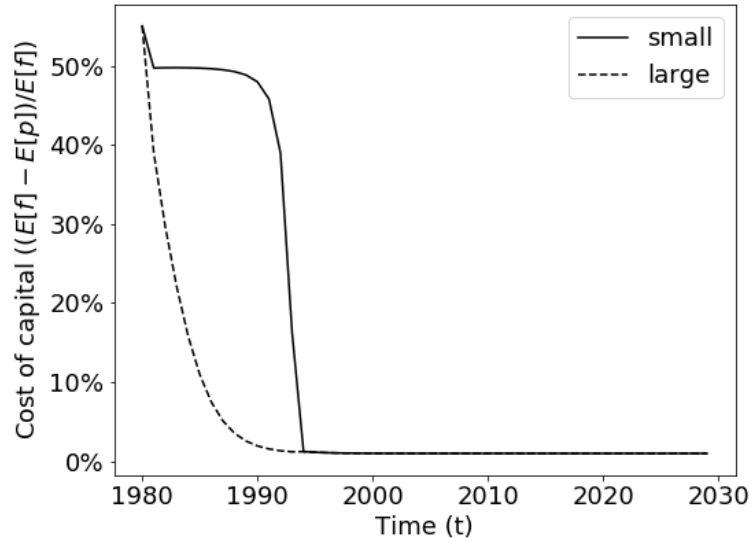


**FIGURE 1.16 Growth in processor performance since the mid-1980s.** This chart plots performance relative to the VAX 11/780 as measured by the SPECint benchmarks (see Section 1.8). Prior to the mid-1980s, processor performance growth was largely technology-driven and averaged about 25% per year. The increase in growth to about 52% since then is attributable to more advanced architectural and organizational ideas. By 2002, this growth led to a difference in performance of about a factor of seven. Performance for floating-point-oriented calculations has increased even faster. Since 2002, the limits of power, available instruction-level parallelism, and long memory latency have slowed uniprocessor performance recently, to about 20% per year. Copyright © 2009 Elsevier, Inc. All rights reserved.

**Figure 3: Investors' Data Choices** The left panel shows the share of the total data processed for each firm. The right panel shows the number of bits processed about each firm.

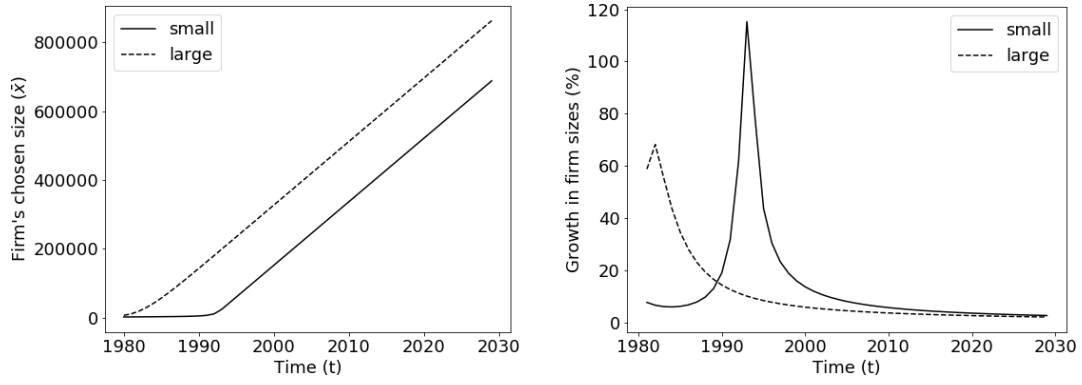


**Figure 4: Cost of Capital for a New Firm** The solid line represents the cost of capital per share,  $E_t[f_{i,t}] - E_t[p_{i,t}]$ , normalized by average earnings per share,  $E_t[f_{i,t}]$ , of the small firm ( $\bar{x}_0 = 2000$ ). The dashed line is the (normalized) cost of capital of the large firm ( $\bar{x}_0 = 10000$ ). Simulations use parameters listed in Table 1.

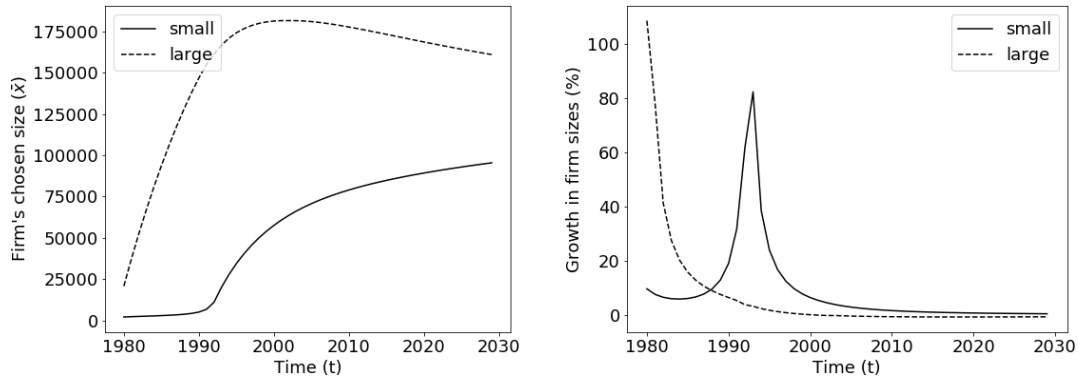




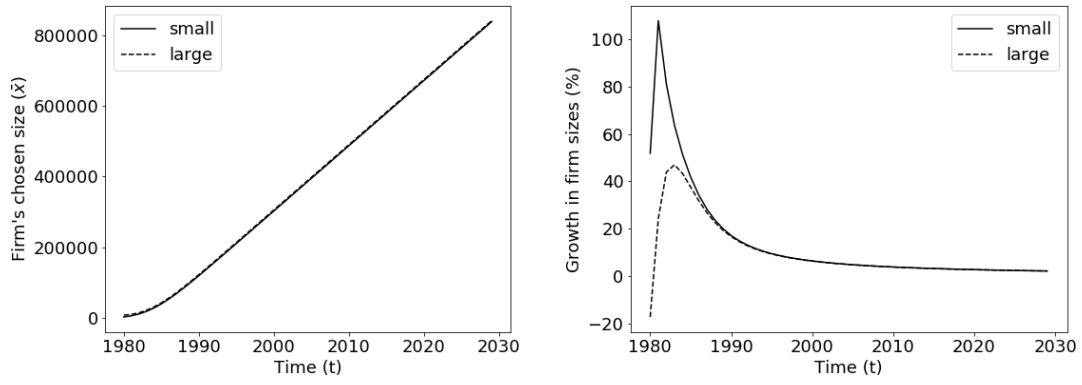
**Figure 5: The Evolution of Small and Large Firms (level and growth rate)** These figures plot firm size  $\bar{x}_t$  (left) and growth in firm size,  $(\bar{x}_t/\bar{x}_{t-1} - 1) \times 100$  (right), for a small firm, with starting size 2000 and a large firm with starting size 10000. Simulation parameters are those listed in Table 1.



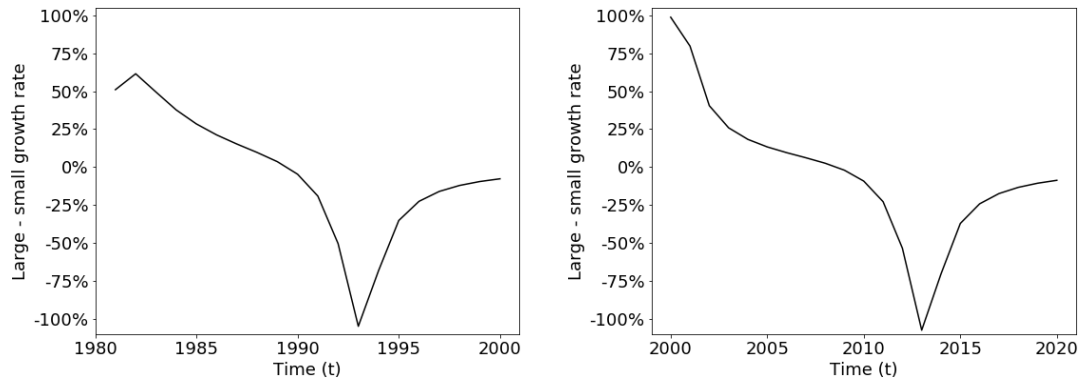
**Figure 6: Without Improvements in Data Processing, Firm Size Converges** These results use the same simulation routine and parameters to plot the same quantities as in Figure 5. The only difference is that these results hold data processing capacity fixed at  $K_t = 5 \forall t$ .

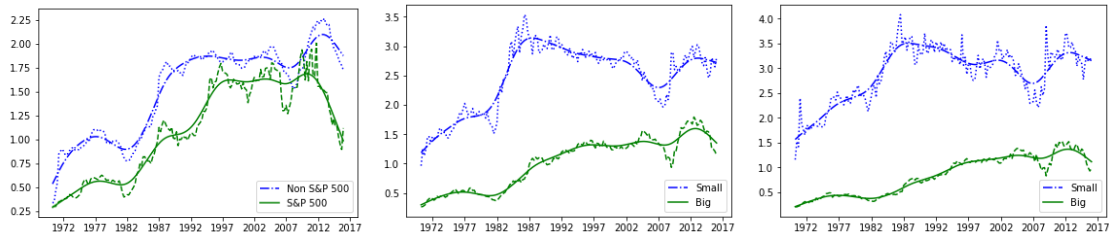


**Figure 7: With Unlimited Data Histories, Small and Large Firms Converge Quickly.** These results use the same simulation routine and parameters to plot the same quantities as in Figure 5. The only difference is that these results set the data availability parameters ( $\alpha, \beta$ ) to be large enough such that the data availability constraint never binds.



**Figure 8: Large Firms Grow Faster in 2000 than in 1980.** Both panels plot a difference in the growth rate of size ( $(\bar{x}_t/\bar{x}_{t-1} - 1) \times 100$ ). The difference is the growth rate of a large firm ( $\bar{x}_t$  starting at 10,000) minus the growth rate of a small firm ( $\bar{x}_t$  starting at 2000). Both are the result of simulations using parameters in Table 1. The left panel shows the difference in firm growth for firms that start in 1980, with  $K_{1980} = 1$ . The right panel shows the difference in firm growth for firms that start in 2000, with  $K_{2000} \approx 527$ .





**Figure 9: 20 month rolling standard deviation of growth in earnings: Large and small firms.** Left panel: Median volatility by whether a firm is a member of the S&P500 index. Middle panel: Median volatility by whether a firm has earnings above or below the median each quarter. Right panel: Median volatility by whether a firm's earnings are in the top or bottom quartile each quarter

**Table 1:** Parameters used in the numerical example

$\mu$	$\sigma$	$\sigma_x$	$r_t$	$\phi_0$	$\phi_1$	$\phi_2$	$\rho$	$\alpha$	$\beta$
15	0.55	0.5	1.01	0.598	0.091	0.0004	0.1	0.249	0.0002

**Table 2: Large Firm Growth Minus Small Firm Growth from Compustat**

For each industry sector and year, large firms are the top 25% largest firms in Compustat; small firms are the bottom half of the firm size distribution. Growth rate is the annual log-difference. Reported figures are equal-weighted averages of growth rates over firms and years.

	prior to 1980	1980 - 2007
Assets	2.1%	8.2%
Investment	14.2%	16.0%
Assets with Intangibles	0.3%	1.1%
Capital Stock	-0.9%	3.7%
Sales	1.4%	2.4%
Market Capitalization	1.1%	8.9%