

Singapore Management University

Institutional Knowledge at Singapore Management University

Research Collection Lee Kong Chian School Of
Business

Lee Kong Chian School of Business

8-2022

Information acquisition and expected returns: Evidence from EDGAR search traffic

Frank Weikai LI

Singapore Management University, wkli@smu.edu.sg

Chengzhu SUN

Hong Kong Polytechnic University

Follow this and additional works at: https://ink.library.smu.edu.sg/lkcsb_research



Part of the [Finance and Financial Management Commons](#), and the [Portfolio and Security Analysis Commons](#)

Citation

LI, Frank Weikai and SUN, Chengzhu. Information acquisition and expected returns: Evidence from EDGAR search traffic. (2022). *Journal of Economic Dynamics and Control*. 141, 1-20.

Available at: https://ink.library.smu.edu.sg/lkcsb_research/5322

This Journal Article is brought to you for free and open access by the Lee Kong Chian School of Business at Institutional Knowledge at Singapore Management University. It has been accepted for inclusion in Research Collection Lee Kong Chian School Of Business by an authorized administrator of Institutional Knowledge at Singapore Management University. For more information, please email cherylds@smu.edu.sg.

Information Acquisition and Expected Returns: Evidence from EDGAR Search Traffic*

Frank Weikai Li[†] Chengzhu Sun[‡]

This Draft: October 2017

Abstract

This paper examines expected return information embedded in investors' information acquisition activity. Using a novel dataset containing investors' access of company filings through SEC's EDGAR system, we reverse engineer their expectations over future payoffs and show that the abnormal number of IPs searching for firms' financial statements strongly predict future returns. The return predictability stems from investors allocating more effort to firms with improving fundamentals and following exogenous shock to underpricing. A long-short portfolio based on our measure of information acquisition activity generate monthly abnormal return of 80 basis points and does not reverse over the long-run. In addition, the return predictability is stronger among firms with larger and lengthy financial filings that are more costly to process. Collectively, these findings support theoretical predictions that costly information acquisition reveals the value of information.

JEL classification: G12, G14

Keywords: Information Acquisition, EDGAR Search, Return Predictability, Market Efficiency

*We are grateful to Utpal Bhattacharya, Claire Hong, Hongqi Liu, Baolian Wang, Jialin Yu and seminar participants at CUHK Shenzhen for helpful comments. All errors are our own.

[†]Singapore Management University, Lee Kong Chian School of Business (Email: wkli@smu.edu.sg)

[‡]Hong Kong University of Science and Technology (Email: csunab@connect.ust.hk)

1 Introduction

Information acquisition and dissemination is key to understanding asset price movements and market efficiency. When information is costly to acquire and price is only partially revealing, economic agents will expend resources and effort to become informed (Grossman and Stiglitz (1980); Verrecchia (1982)), and in doing so, move prices closer to fundamental value. A central prediction from theories of costly information acquisition is that more investors will choose to become informed when they perceive greater benefits from doing so, holding the cost of information acquisition constant. Although theories offer clear and rich predictions, empirical evidence on the relation between information acquisition behavior and value of information is sparse in financial markets, potentially due to the difficulty of directly measuring information acquisition activities of investors.

In this paper, we take advantage of a novel dataset containing investors' access of regulatory filings through SEC's EDGAR (Electronic Data Gathering, Analysis, and Retrieval) system to study the implications of information acquisition activities on firm value. Because EDGAR system is the main sources of firms' regulatory filings and SEC maintains a log file of all activity performed by users on EDGAR, we are able to directly observe investors' information acquisition activity for a broad cross-section of firms over more than 10 year sample period.

Our research objectives are two-fold in this paper. First, we examine the determinants of investors' information acquisition through EDGAR website. Motivated by theories of information acquisition¹, we posit that information acquisition activities should be negatively related to the cost of gathering information and positively related to the value of information. To this end, we use the number of unique IP addresses searching for SEC filings through EDGAR as proxy for investors' information acquisition. We then run cross-sectional regression of our information acquisition proxy on several firm characteristics associated with information cost. Specifically, we hypothesize that firms with higher investor visibility and attention, and better information environment will attract more information acquisition.

¹There is a large theoretical literature on information acquisition, e.g., Grossman and Stiglitz (1980), Diamond and Verrecchia (1981), Verrecchia (1982), Hellwig (1980), Admati (1985), Mele and Sangiorgi (2015).

as these stocks are more accessible in investors' mind and less costly to analyze. We also expect investors to have stronger incentive to acquire information on firms with higher valuation uncertainty. Using firm size to proxy for investor visibility, trading volume to proxy for investor attention (Gervais, Kaniel, and Mingelgrin (2001); Barber and Odean (2007)), analyst coverage to proxy for information environment (Hong, Lim, and Stein (2000)), and idiosyncratic volatility to proxy for valuation uncertainty (Zhang (2006)), we find evidence consistent with theories. These four firm characteristics explain 55% of cross-sectional variation of information acquisition across firms. Further tests show that information acquisition through EDGAR also increases following negative past return performance and among firms with lower institutional ownership, but these additional characteristics does not significantly improve the explanatory power over our baseline model.

After implementing a simple characteristic-based model of expected information acquisition, we proceed to examine our second research question. That is, abnormal level of information acquisition should be positively related to the expected benefits of trading on information. This is based on the simple premise that when attention-constrained investors decide how to allocate their time and effort, they will have a strong preference for firms with the largest price appreciation or depreciation potential. In reality, due to short-sale constraints, investors will more likely engage in costly information acquisition when the expected return of a stock is positive.

To test this, we extract the part of number of IPs unexplained by firm characteristics to reverse engineer investors' expectation over future payoffs. Consistent with the idea that information acquisition embeds the value of information, we show that the abnormal number of IPs (denoted as AIP) requesting EDGAR filings strongly predict subsequent stock returns. An equal-weighted, monthly rebalanced, long-short strategy that buys stocks in highest decile of AIP and sells stocks in the lowest decile of AIP generates 52 to 82 basis points per month after adjusting for the Carhart (1997) four factors and highly significant. Adjusting for the recently proposed factor models—the Fama and French (2016) five-factor model, the Hou, Xue, and Zhang (2015) q-factor model, or the Stambaugh and Yuan (2016) mispricing factor model—does not affect the return spread of the long/short portfolio much. The abnormal

return of AIP strategy is much weaker for value-weighted portfolios. The high-minus-low AIP strategy is only around 30 basis points per monthly and mostly insignificant. This is expected given short-sale constraint is less binding among big stocks, so the direction of the information contained in AIP is more ambiguous for big stocks. Using several proxies of short-sale constraints, we confirm that the positive expected return information embedded in information acquisition is more pronounced among stocks that are difficult to short ex-ante.

The return predictability associated with abnormal number of IPs persists for two quarters, and does not reverse in the subsequent months. The persistence in return predictability alleviates concerns that our findings is driven by temporary price pressure caused by noise traders that reverse over the long-run (Da, Engelberg, and Gao (2011)).

In a Fama-MacBeth regression setting, we confirm that AIP has additional explanatory power for future stock returns when we control for the standard cross-sectional return predictors such as firm size, book-to-market ratio, momentum, short-term reversal, idiosyncratic volatility, turnover and institutional ownership. The return predictability of AIP is also **not** affected by alternative explanations such as post-earnings announcement drift, earnings announcement premium and investor disagreement. Looking into different types of EDGAR filings, we find the return predictability of AIP comes mainly from those searching for firms' annual accounting report 10-K (AIP_10K). As gathering and analyzing 10-K report is more costly than other SEC filings and more reflecting deliberate information acquisition behavior, the stronger predictability of AIP_10K is consistent with theories of costly information acquisition. To further substantiate our argument, we use the file size and word count of 10-Ks as proxy for the complexity of financial disclosure (Loughran and McDonald (2014)), and find return predictability of AIP is indeed stronger among firms with large and lengthy 10-Ks.

Having established the robustness of the return predictability of abnormal number of IPs, we test the sources of return predictability. The underlying assumption in this paper is that investors rationally allocate more effort and resources towards underpriced stocks with high expected return. As mispricing implies the separation of stock prices from firms' fundamental value, we conjecture two non-mutually exclusive channels through which investors can

identify mispricing. The first channel is that investors' information acquisition activity reveals their favorable expectation of firms' fundamental performance that are yet to be priced in market. Consistent with the first channel, we find AIP strongly predicts future change of firms' quarterly Return-on-Assets and analyst consensus forecast, after controlling for past profitability and other determinants of firms' fundamental performance. A second channel is that investors identify mispricing by observing changes in stock prices due to exogenous reasons. Supporting the second channel, we show the abnormal number of IPs searching for EDGAR filings increases significantly for firms experiencing mutual fund outflow-induced selling pressure. Taken together, our evidence suggest that investors expand greater resources and effort towards undervalued stocks and these findings are much more difficult to reconcile with alternative explanations such as omitted risk factor or changes in investor visibility (Merton (1987))².

Finally, we provide some suggestive evidence on the types of investors conducting informed searches on firms' fundamentals through EDGAR. We show that abnormal number of IPs positively predicts net purchases by hedge funds in the following quarter. In contrast, AIP does not have predictability for net purchases by mutual fund managers. These results are consistent with the idea that investors searching for financial filings through EDGAR are more sophisticated than those searching through Google search engine, and hedge funds may potentially be part of these sophisticated investors.³

This paper contributes to several strands of the literature. First, our results offer strong empirical evidence supporting theories of information acquisition that costly information acquisition is positively related to the expected benefits of information. Using the novel EDGAR log file dataset, we construct a direct measure of investors' information acquisition activity, and show its strong predictability for firms' future returns and fundamentals. Du (2015) shows that the number of web visits to SEC filings of insider trades predicts post-filing stock return in the short-run. Although similar in spirit, our paper differs as we

²Alternative explanations based on omitted risk factor or changes in investor base all work through discount-rate channel, while the return predictability of AIP operates (partially) through cashflow channel.

³Drake, Quinn, and Thornock (2017) document that EDGAR users tend to have higher education level and more likely to work in major cities with more accounting and finance jobs.

study a much broader sample of SEC regulatory filings and longer horizon returns. We also test the channels underlying the return predictability results. Using EDGAR search data, Cohen et al.(2017) document that mutual funds tend to track a particular set of firms and insiders, and their tracked trades generate abnormal performance. Lee and So (2017) study the information content of analysts' coverage decisions and show the abnormal number of analyst coverage positively predict future firm performance. By extracting all internet users' information acquisition activities through EDGAR site, our measure captures the expected return information embedded in the collective behavior of a much larger set of market participants—millions of unique users. In addition, analysts' incentives are found to be distorted by generating trading commissions for their brokerage houses or currying favor with firm management (Ke and Yu (2006)), while these distortions are less likely among EDGAR users. Empirically, we show the return predictability of AIP is not affected after controlling for analyst coverage proxies.

Our paper also contributes to the growing literature on the effect of investor attention and information acquisition on asset prices and capital market efficiency. Da, Engelberg, and Gao (2011) show that retail investors' attention as captured by Google search volume causes transitory price pressures on attention-grabbing stocks. Using news searching activity via Bloomberg terminal as proxy for institutional investors' attention, Ben-Rephael, Da, and Israelsen (2017) find that institutional attention facilitates the timely incorporation of fundamental information into asset prices. More related to our study, Drake, Roulstone, and Thornock (2015) show that EDGAR-based information acquisition affects the pricing of earnings-related news. These papers all examine the effect of information acquisition on the pricing of *public* announced news. Our paper is different as we directly infer investors' *private* expectation of firm value through their collective actions.

Finally, our work contributes to the emerging literature on extracting intelligence latent in the collective "wisdom of crowds". Chen, De, Hu, and Hwang (2014) find evidence that investors' social media posts help predict stock return. Lee, Ma, and Wang (2015) show that investors' co-search pattern via EDGAR website could help identify peer firms better than traditional industry benchmark. Huang (2016) finds that consumer opinion expressed on

firms' products in Amazon.com contain value-relevant information about firm fundamentals and stock prices. Similarly, Green, Huang, Wen, and Zhou (2017) document employer reviews on Glassdoor reveals valuable information about their employers' fundamentals. Our paper is complementary to the above mentioned studies as we infer agents' expectation not via what they "say", but through what they actually "do".

Our results that information acquisition activity predicts future returns do not imply market is inefficient. As pointed out by Grossman and Stiglitz (1980), a fully efficient market where prices instantaneously reflect all available information cannot sustain an equilibrium when information is costly to acquire and analyze. Rather, our evidence is mostly consistent with the idea of "efficiently inefficient markets" (Pedersen (2015)), where competition among investors makes market almost efficient, but the market remains so inefficient so that these investors are compensated for their costs of acquiring information.

The remainder of this paper is organized as follows. Section 2 describes the data, presents summary statistics and examines the determinants of information acquisition. Section 3 shows that abnormal level of information acquisition reveals investors' expectation over future return. Section 4 tests the channels underlying the return predictability results. Section 5 concludes.

2 Data and Methodology

2.1 Sample Construction

Our IP search volume data comes from Security and Exchange Commission's (SEC) EDGAR log file database which records all website search traffic for SEC filings since 2003.⁴ Each search record contains information on user's unique Internet Protocol (IP) address (anonymized)⁵, timestamp, searched company identified by Central Index Key (CIK) and searched specific filing identified by the unique SEC accession number.⁶ Following Lee, Ma,

⁴The data is available for download at <https://www.sec.gov/data/edgar-log-file-data-set.html>.

⁵The EDGAR log file dataset provides the first three octets of the IP address with the fourth octet obfuscated with a 3 character string that preserves the uniqueness of the last octet without revealing the full identity of the IP.

⁶The detailed log file record elements are described at https://www.sec.gov/files/EDGAR_variables_FINAL.pdf

and Wang (2015) and Ryans (2017), we firstly filter the raw log data to eliminate the requests made by robots or by automated web crawlers since such massive and indiscriminate requests are uninformative for our research question.⁷ Next, we match CIK in EDGAR log filings to that in COMPUSTAT to identify public companies, and retrieve the filing type and filing date for each requested file by linking the accession number back to the Master Index files maintained by the SEC.⁸ We classify these filings into six groups: 10-K, 10-Q, 8-K, insider, registration, and proxy.⁹ Finally, we calculate the monthly IP search volume for each filing category at firm level by counting the total number of unique IP address that searched one category of SEC filings of a specific company within one month window. We define IP_total as the total number of unique IP addresses searching all six types of EDGAR filings. Drake, Roulstone, and Thornock (2015) documents that periodic accounting reports are the type of SEC filings most frequently requested by investors through EDGAR website. As a result, we also compute two additional measures of information acquisition targeting specifically firms' periodic accounting reports. IP_funtl (IP_10K) is the total number of unique IP address searching 10-K, 10-Q and 8-K files (10-K files). Our sample runs from January 2003 to December 2014.¹⁰

It is important to mention that there are other sources for investors to access financial filings such as firm's investor relations website and Yahoo! Finance. Data vendors such

⁷First, following Lee, Ma, and Wang (2015), we exclude the searching records to each unique user who download more than 50 unique firms' filings in one day. The user is identified by unique IP address. Secondly, following Ryans (2017) and Drake, Roulstone, and Thornock (2015), we remove log records that reference an index (idx=1), as index pages only provide the links to filings rather than filing itself. Third, following Ryans (2017), we keep the request records with successful document delivery (code=200). Next, we further exclude the search records to users with filing requests more than 25 per minute or more than 500 per day, or with more than 3 unique CIK searching per minute. Finally, we only keep one search record for a specific filing (unique accession number) to each user in a given day. This step is to avoid the duplicated records due to users' multiple views for same document especially after the adoption of XBRL filing in 2009. For user who view financial reporting of XBRL adopted firm in interactive data format, every click on a different footnote will generate a new search record although it references the same document.

⁸Further details about the EDGAR index files can be found at <https://www.sec.gov/edgar/searchedgar/accessing-edgar-data.htm>

⁹We define 10-K category as filing type in "10-K", "10-K/A", "10-K405", "10-K405/A", "10-KSB", "10KSB", "10-KSB-A", "10KSB/A", "10-KT", "NT 10-K", "10-KSB40"; 10-Q category as filing type in "10-Q", "10-Q/A", "10QSB", "10-QSB", "10QSB-A", "NT 10-Q", 8-K category as filing type in "8-K", "8-K/A"; Insider category as filing type in "SC 13G", "SC-13D", "SC 13G/A", "SC 13D/A", "3", "4", "5"; registration category as filing type in "S-1", "S-1/A", "S-3", "S-3/A", "S-3ASR", "424B5", "424B4", "424B3", "424B2", "FWP"; proxy category as filing type in "DEF 14A", "DEF 14C", "DEFA14A", "DEFM14A", "DEFR14A", "DEFM14C".

¹⁰There are significant gaps in the data prior to March 2003 and between September 2005 and May 2006, due to lost or corrupt log file. As a result, we exclude these months from our sample in our analysis.

as Bloomberg and FactSet also provide investors access to these reports. As a result, our analysis of the EDGAR server log cannot capture all the views/downloads that the entire universe of investors are conducting on company filings. However, the EDGAR server still possesses several advantages compared to other information sources. First, it is questionable that investors primarily use the company website to retrieve SEC filings. As an example, Monga and Chasan (2015) quote General Electric CFO Jeffrey Bornstein, who noted that GE's 2013 annual report was downloaded from their investor relations website just 800 times.¹¹ For the same annual report, the EDGAR logs record 21,987 (4,325) downloads in the year (two months) following its filing. Secondly, some firms, such as Google (Alphabet, Inc) and ExxonMobil, forward investors directly to the EDGAR website to obtain their SEC filings. For such cases where the investor relations department links the investors to the EDGAR site, these views/downloads will be captured in the SEC server. Third, other sources of the company information often condense income statement and balance sheet information into pre-specified bins. As a result, some critical components of firms' financial information may be misrepresented. Lastly, investors could better assess firm's future prospects by reading the qualitative information contained in 10-K filings, which is not available in these data consolidators (Loughran and McDonald (2011)).

We obtain monthly stock returns from the Center for Research in Security Prices (CRSP) and annual accounting data from Compustat. Our sample of stocks starts with all common stocks traded on NYSE, Amex, and NASDAQ. We adjust the stock returns by delisting. If a delisting return is missing and the delisting is performance-related, we set the delisting return to be -30% (Shumway (1997)).

We use standard control variables in our empirical analysis. *Size* (LnME) is defined as the natural logarithm of market capitalization at the end of June in each year. *Book-to-market ratio* (LnBM) equals to the most recent fiscal year-end report of book value divided by the market capitalization at the end of calendar year t-1. Book value equals the value of common stockholders' equity, plus deferred taxes and investment tax credits, minus the book value of preferred stock. *Momentum* (Mom) is defined as the cumulative holding-period return from month t-12 and t-2. We follow the literature by skipping the most recent month's return

¹¹<https://www.wsj.com/articles/the-109-894-word-annual-report-1433203762>.

when constructing the *Momentum* variable. The *short term reversal measure* (REV) is the prior month’s return. *Turnover12* is the monthly trading volume over shares outstanding, averaged within past 12 months. Since the dealer nature of the NASDAQ market makes its turnover difficult to compare with the turnover observed on NYSE and AMEX, we follow Gao and Ritter (2010) by adjusting trading volume for NASDAQ stocks.¹² *Institutional ownership* (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by total shares outstanding. *Idiosyncratic volatility* (IVOL) is the standard deviation of the residuals from the regression of daily stock excess returns on Fama and French (1993) 3-factor returns within a month (Ang, Hodrick, Xing, and Zhang (2006)). Institutional ownership data of stocks are available from Thomson Reuters (formerly CDA/Spectrum) Institutional Holdings database (13F). Coverage is the log one plus number of analyst following a firm from I/B/E/S. We download the file size and number of words of 10-Ks for all publicly-traded firms from WRDS SEC Analytics.

We also construct measures for trading activities by hedge funds. Using the list of hedge funds provided by Jiang (2014), we retrieve their quarterly holdings from Thomson Reuters CDA/Spectrum Institutional (13F) database. We define net purchases of stock i by hedge funds in quarter q as follows:

$$NetBuy_{i,q} = \frac{Shrown_{i,q}}{Shrout_{i,q}} - \frac{Shrown_{i,q-1}}{Shrout_{i,q-1}} \quad (1)$$

where $Shrown_{i,q}$ is the number of shares of firm i held by hedge funds in quarter q , and $Shrout_{i,q}$ is firm i ’s number of shares outstanding in quarter q . To provide a basis for comparison, we construct a similar measure for mutual fund investors.

2.2 Summary Statistics

Panel A of Table 1 reports the time-series average of the cross-sectional means and standard deviations of the variables for the full sample. The average number of unique IPs searching for all six types of EDGAR filings is 155 in a month. The cross-sectional standard

¹²Specifically, we divide NASDAQ volume by 2.0, 1.8, 1.6, and 1.0 for the periods before February 2001, between February 2001 and December 2001, between January 2002 and December 2003, and after January 2004, respectively.

deviation is 317, indicating large cross-sectional variation among firms. Consistent with Drake, Roulstone, and Thornock (2015), annual reports is the mostly frequently searched SEC filings, with an average of 60 IPs requesting in a month. IPs searching for 10-Q and 8-K files are relatively less frequent. The average institutional ownership in our sample is 55%, reflecting the rapid growth of assets managed by institutional investors during our sample period. The remaining summary statistics are well known and do not require additional discussion.

Panel B reports the pairwise rank correlation among our variables. As we can see, the three IP variables are highly correlated. This is expected as periodic accounting reports consist of the largest fraction of EDGAR search requests. The number of IPs is also highly correlated with firm size, analyst coverage and turnover, suggesting firms with high investor visibility and attention have more EDGAR users. The number of IPs is negatively correlated with stock idiosyncratic volatility. However, this is mainly due to the size effect: small firms with high return volatility attract less EDGAR searching. As we will see later, once we control for firm size, the number of IPs becomes positively correlated with idiosyncratic volatility, potentially because the incentives of acquiring information is larger when firm valuation is more uncertain (Grossman and Stiglitz (1980)).

2.3 Cross-sectional Determinants of Number of IPs

Theories of endogenous information acquisition suggest that information acquisition activity is a function of both the cost of acquiring information and the benefits of trading on acquired information. To isolate investors' expected payoffs of information acquisition, we need a model of expected information acquisition activities. To this end, we develop and implement a simple characteristics-based model of expected information acquisition and identify the discrepancies between realized and expected level of information acquisition. Calculating these discrepancies requires proxies for information acquisition and firm characteristics useful in estimating the expected level of information acquisition activities.

Our proxy for information acquisition activity is the number of unique IP addresses searching for EDGAR filings for each firm at given month. To mitigate data mining concerns,

we use three measures capturing information acquisition activities for different types of EDGAR filings. *IP_total* is the total number of unique IPs searching for all types of EDGAR filings, and *IP_funtl* (*IP_10K*) is the total number of unique IPs searching 10-K, 10-Q and 8-K files (10-K files). Our choice of firm characteristics is guided by information acquisition theories. Specifically, we hypothesize that firms with higher visibility and investor attention will attract more information acquisition, as these firms are more accessible in investors' mind. We also conjecture that the strength of firms' information environment will affect information acquisition, although the direction of effect is not clear. On the one hand, firms with abundant public information will be less costly to analyze, so we expect information acquisition will increase with the quality of firm's information environment. On the other hand, better information environment means the stock is less likely to be mispriced ex-ante, so investors' incentive to acquire private information will be reduced. Finally, we expect investors to have stronger incentive to acquire information on firms with higher valuation uncertainty. Following prior literature, we use firm size to proxy for investor visibility, trading volume to proxy for investor attention (Gervais, Kaniel, and Mingelgrin (2001); Barber and Odean (2007)), analyst coverage to proxy for information environment¹³ (Hong, Lim, and Stein (2000)), and idiosyncratic volatility to proxy for valuation uncertainty (Zhang (2006)).

We calculate the abnormal number of IPs by fitting monthly cross-sectional regressions of the raw number of IPs to isolate the components of number of IPs not attributable to firms' size, turnover, analyst coverage and idiosyncratic volatility. To mitigate the effect of outliers, we use the log of one plus number of IPs when estimating firms' abnormal number of IPs. Specifically, we calculate abnormal number of IPs for firm i in month t by estimating the following regression:

$$\text{Log}(1 + IP_{i,t}) = \beta_0 + \beta_1 \text{LnME}_{i,t} + \beta_2 \text{Coverage}_{i,t} + \beta_3 \text{Turnover}_{i,t} + \beta_4 \text{IVOL}_{i,t} + \epsilon_{i,t} \quad (2)$$

where *LnME* is the log of market capitalization in month t , *Coverage* is the log of one plus

¹³Another motivation for including analyst coverage is that Lee and So (2017) shows analyst coverage contains information about future stock return. By including analyst coverage as a regressor, any expected return information embedded in number of IPs will be incremental to that contained in analyst coverage proxies.

analyst coverage, Turnover12 is the monthly turnover averaged over past 12 months, and IVOL is the daily idiosyncratic volatility following Ang, Hodrick, Xing, and Zhang (2006). We define abnormal number of IPs for each firm-month as the regression residuals from equation (1). We use the notation AIP to refer to the abnormal number of IPs, where higher values correspond to firms that have greater number of IPs searching their EDGAR filings given their size, trading volume, analyst coverage and volatility.

Table 2 reports the time-series average coefficients from estimating equation (1). The three panels correspond to three different measures of IPs as dependent variables. To see the improvement of R^2 , We add the explanatory variables one by one from Column (1) to Column (7). Consistent with our hypothesis, information acquisition activities increase with firm size (t-stat=69.44), as larger firms are more visible to investors. Size alone explains 40% of cross-sectional variation of number of IPs. Column (2) and (3) show that information acquisition increases with the strength of firms' information environment and investor attention, proxied by analyst coverage and turnover, respectively. Column (4) further shows that number of IPs increases with return volatility after controlling for firm size. This finding suggests that investors' demand for information is larger for firms with more uncertain value. Column (4) also shows that these four firm characteristics used in equation (1) explains 55% of the cross-sectional variation of number of IPs on average. The results are similar in Panel B and C where the dependent variable are IP_fundl and IP_10K, respectively.

The four firm characteristics used in equation (1) were selected based on theories and parsimony but may omit other firm characteristics that drive variation in the expected level of information acquisition activity. For example, investors may be attracted to firms with extreme past performance and glamour characteristics (Barber and Odean (2007)). To examine the explanatory power of other firm characteristics, we add stock's past 12-month return, book-to-market ratio and institutional ownership iteratively from Column (5) to Column (7). The results suggest more investors searching for EDGAR filings when the firm performs poorly in the past year and behave like value stocks. However, adding these additional three characteristics improve the average R^2 of equation (1) by only 0.5%, suggesting the limited incremental explanatory power of past return performance, book-to-

market and institutional ownership. In the robustness test below, we show that the inclusion of other firm characteristics does not significantly affect the return predictability of AIP.

As there might be nonlinear relationship between abnormal number of IPs and firm characteristics, we further look at average stock characteristics across decile portfolios sorted on abnormal number of IPs searching for 10-K files (AIP_10K). Higher (lower) deciles correspond to firms with abnormally high (low) number of IPs. Panel C of Table 1 reports the time-series average of cross-sectional mean values of each variable for each decile. First, the observation counts show that there are about 330 firms in each decile, suggesting that our measure of information acquisition is available for a broad cross-sectional sample of 3300 firms per month. Second, the table shows that AIP is positively correlated with the raw number of IPs searching for EDGAR filings. Third, AIP is, by construction, uncorrelated with firm size, analyst coverage, turnover, and volatility, although middle portfolios have slightly larger size and turnover. Last, the panel shows that firms in the extreme deciles have lower institutional ownership and more likely to be value stocks.

3 Information Acquisition and Future Stock Returns

When investors expend effort and time to acquire firms' fundamental information, they must perceive some benefits of utilizing such information. Hence a key hypothesis in this paper is that costly information acquisition activities reveal investors' perception of expected payoff. Although in theory, the direction of the information content could be either positive or negative, in reality we expect firms with larger number of IPs searching their EDGAR filings to have better future performance due to short-sale constraints. In addition, the positive predictive power of AIP should be stronger among small firms with stonger short-selling constraints. In this section, we test the relation between abnormal number of IPs and future return using both portfolio sorts and Fama-MacBeth regression.

3.1 Portfolio Sorts

In this section, we show that stocks sorted based on their abnormal number of IPs generate significant return spreads. We conduct the decile portfolio sorts as follows. At the end of each month, we sort stocks into deciles by their AIP. We then compute the average return of each decile portfolio over the next month. This gives us a time series of monthly returns for each decile. We use these time series to compute the average excess return of each decile over the entire sample. As we are most interested in the return spread between the two extreme portfolios, we also report the return to a long–short portfolio (i.e., a zero-investment portfolio that goes long the stocks in the highest AIP decile and shorts the stocks in the lowest decile). Our sample is from January 2003 to December 2014.

Table 3 reports the average monthly excess return of each decile portfolio. Panel A reports the equal-weighted portfolio return, and Panel B reports the value-weighted return. The three columns in each panel correspond to sorting based on the abnormal number of IPs searching for three different types of EDGAR filings. Panel A shows a strong positive relation between AIP and future returns, regardless of which IP variables we use. For sorts based on AIP_total, firms in the highest decile of AIP outperforms the firms in the lowest decile by 71 basis points per month on an equal-weighted basis (t-stat=3.18). The results are stronger when we do the portfolio sorts based on AIP_funtl and AIP_10K. Specifically, the high-minus-low monthly return spread is 100 basis points (t-stat=4.70) based on AIP_10K, which corresponds to an annualized return of 12%. The evidence show that aggregate information acquisition activities across EDGAR users reveal an economically large source of predictable return across firms. The economic magnitude is quite impressive given the fact that many other well-documented asset pricing anomalies are no longer profitable in our sample period (Chordia, Subrahmanyam, and Tong (2014); Green, Hand, and Zhang (2017)).

The larger return spread based on IPs searching for 10K compared to IPs searching for other types of EDGAR filings is consistent with information acquisition theories. Firms' annual report is among the most lengthy and difficult-to-read SEC filings. Annual reports contain detailed annual operating and financial performance and metrics, suggesting that digesting these report require a large amount of effort and attention from investors' part.

Compared to 10-K files, 10-Q and 8-K files are usually shorter and easier to digest, and investors driven to these type of filings more likely respond to contemporaneous news events, rather than reflecting deliberate information acquisition choice. Given the substantial higher cost of acquiring and analyzing 10-K files, the expected benefits perceived by investors should also be larger, which is consistent with our results.

The return spread of high-AIP minus low-AIP portfolio is considerably smaller and less significant when returns are value weighted. The high-minus-low return is only around 30 basis points per monthly, and not significant. This is consistent with our prior that when short-sale constraints are less binding for big firms, the information content embedded through EDGAR searching could be either positive or negative. Investors could take unconstrained short position to benefit from the negative information they obtained through EDGAR filings. This implies that, ex-ante, we do not have a clear directional prediction between abnormal number of IPs and future return. When we take a closer look at the value-weighted portfolio return from decile 1 to decile 10, we find the relation between AIP and average return is an inverted U shape. This could be due to the fact that firms in the top decile of AIP are a mixture of firms with high and low expected return, and in aggregate they cancel out.

Table 4 examines the relation between abnormal number of IPs and firms' future return after controlling for portfolios' exposure to standard asset pricing factors. The table reports the monthly Carhart (1997) four factor alpha for decile portfolios sorted on AIP, as well as the long/short hedge portfolio. The four factor alpha is the intercept from a regression of the portfolio's excess return on the contemporaneous excess market return (MKTRF), the size factor (SMB), the value factor (HML), and the momentum factor (UMD). Panel A shows that AIP continuously to predict strong positive return spread cross-sectionally for equal-weighted portfolios. The four-factor alphas of the long/short portfolio range from 52 to 82 basis points per month and are highly significant. Moreover, in the case of AIP_10K, the alphas are fairly symmetric across deciles. The lowest AIP decile portfolio generates four factor alpha of about -34 basis points (t-stat=-2.84), and the highest AIP decile generates positive alpha of 48 basis points (t-stat=3.30). The evidence suggests that when short-

sale constraints is binding, investors would rationally allocate less effort towards firms with negative expected return. Panel B of Table 4 shows the portfolio alpha for value-weighted returns. Again, we find the results are generally weaker, both economically and statistically. The four-factor alpha of long/short portfolio ranges from 14 to 42 basis points, and are either insignificant or only marginally significant.

To emphasize the importance of measuring the abnormal level of information acquisition activity when uncovering expected return information, we conduct a parallel portfolio tests when ranking firms into deciles based on the raw number of IPs searching for EDGAR filings in Table 5. Panel A reports the equal-weighted excess return and Panel B reports the equal-weighted four-factor alpha. The results show that the raw number of IPs is not significantly correlated with firms' future returns, regardless of which IP variable we use. The monthly four-factor alpha of the long-short portfolio based on raw number of IP ranges from -20 to 9 basis points, and are never significant. The lack of significant predictive power of raw IP suggests that it is important to control for the expected level of information acquisition activities when uncovering investors' expected payoff.

To get a better sense of how the AIP strategy performs if an investor could get access to the EDGAR log file data at monthly frequency, we plot the cumulative returns to the low and high AIP_10K decile portfolio, as well as the long-short hedge portfolio in Figure 1. The blue line shows that one dollar invested in the lowest AIP_10K decile portfolio at the beginning of 2003 will grow to two at the end of 2014. One dollar invested in the highest AIP_10K decile portfolio will grow to 7.5. The grey line shows that one dollar will grow to almost four dollar if investing in the long-short hedged portfolio, with a smooth return path.

3.2 Robustness

In Table A1, we examine the robustness of our portfolio sorts. For brevity, we focus on the sorts based on AIP_10K. The first row shows the return spread when returns are weighted by past month gross return, as suggested by Asparouhova, Bessembinder, and Kalcheva (2013). The gross-return-weighted return spread is 1.1% ($t=5.16$). Row (2) and (3) show that our results barely change when we subtract the characteristic-matched portfolio

(Daniel, Grinblatt, Titman, and Wermers (1997)) or industry-level return from stock return. This suggests the nature of information contained in costly information acquisition behavior is firm-specific. In the fourth row, we augment the Carhart (1997) four-factors with the Pástor and Stambaugh (2003) liquidity factor. The Pástor and Stambaugh (2003) five-factor adjusted alpha is 0.80% ($t=4.23$) for the equal-weighted portfolio and 0.35% ($t=1.78$) for the value-weighted portfolio. The fifth row shows that our results hold when we use the Fama and French (2016) five factors to calculate alphas, with a monthly return spread of 0.69% ($t=3.36$) for the equal-weighted portfolio. This suggests our long-short portfolio is not merely loading on the profitability and investment factor as proposed by Fama and French (2016). The sixth row shows that our results still hold when we use the Stambaugh and Yuan (2016) mispricing factor model to compute alpha. The portfolio generates equal-weighted alpha of 0.89% ($t=4.42$) and value-weighted alpha of 0.27% ($t=1.35$). Using Hou, Xue, and Zhang (2015) Q-factor model also do not change our results, as shown in the seventh row. The eighth row of Table A1 shows that our results survive when we exclude stocks whose market capitalization are in the bottom quintile of NYSE size distribution. Again, the strategy based on AIP generates a monthly four-factor alpha of 0.52% ($t=2.58$) and 0.28% ($t=1.35$) when returns are equal-weighted and value-weighted, respectively. The ninth row reports the long-short alphas if we skip six months between when we sort stocks and when we measure strategy returns. The purpose of this test is to mimick the profits an investor would generate in reality since SEC delay the release of EDGAR log file data by 6 months. The equal-weighted alpha reduces quite a bit in this case, but nonetheless still significant with an equal-weighted four-factor alpha of 0.53% ($t=2.23$). The tenth and eleventh rows show that the long-short portfolio generates significant alpha in two subperiods: one from 2003 to 2008 and another from 2009 to 2014.

Our results are insensitive to the specific model of calculating abnormal number of IPs, as we show in Table A2. The first row shows that the long-short portfolio based on AIP_10K calculated using model (7) of equation 2 generates four-factor alpha of 0.66% ($t=3.95$). In the second row, we show that positive relation between AIP and returns holds for change-based specifications, which mitigates concerns that the return predictability of AIP is driven by an

omitted firm-fixed effect not controlled for in our model of AIP or multivariate regressions. The long-short portfolio based on change of AIP_10K relative to its 12-month moving average generates equal-weighted four-factor alpha of 0.88% ($t=4.82$). In the third row, we include the square terms of the four firm characteristics when calculating AIP to account for the nonlinear relation between number of IPs and firm characteristics. The 4-factor alpha is 0.689% and 0.552% for equal- and value-weighted portfolio, respectively. In the fourth row, we control for the lagged number of IPs when calculating AIP, and the alpha is still significant.

3.3 The Role of Firm Size and Arbitrage Frictions

Our previous results show that the long/short portfolio alpha is only significant for equal-weighted return, but not value-weighted return. This raises the concern that the return predictability of AIP strategy only exists among small capitalization stocks. To take a closer look at the role of firm size, we reports the results by size quintiles in Table 6. For each month, we group all stocks into size quintiles based on the NYSE size breakpoints. Within each size quintile, we further sort stocks into quintiles based on AIP_10K. The table reports the Carhart (1997) four-factor alpha for the 25 portfolios: equal-weighted returns in Panel A and value-weighted returns in Panel B. We also report the alpha for each size quintile of the high-AIP minus low-AIP portfolios. The result shows that the return predictability of AIP is strongest among microcap stocks, but interestingly, the result also shows that it is not limited to only the smallest size quintile. The high-minus-low AIP portfolio generates an equal-weighted four-factor alpha of 0.45% ($t=2.29$) and value-weighted alpha of 0.31% ($t=1.89$) in the largest size quintile. The alpha is also significant in the middle size quintile group, but are insignificant in size quintile 2 and 4.

The findings in Table 6 show that the return predictability of AIP is more pronounced for small firms than for large firms, which could be driven by two non-mutually exclusive channels. The first reason is that the latent information embedded in the number of IPs searching EDGAR files could be either positive or negative when short-sale constraints are not binding. Given large firms have less short-sale impediments, the direction of return

predictability among large firms is more ambiguous. An independent channel that could reinforce the weak return predictability among these stocks is that whatever information contained in the EDGAR searches, they are arbitrated away quickly due to less arbitrage frictions (e.g., short-sale costs, liquidity, non-fundamental volatility) among large firms. We now explore the return predictability of AIP with other measures of limits to arbitrage.

Following the literature, we investigate the role of three limits-to-arbitrage measures: idiosyncratic volatility (Stambaugh, Yu, and Yuan (2015); Pontiff (2006)), residual institutional ownership (Nagel (2005)) and residual analyst coverage (Hong, Lim, and Stein (2000)). At the end of each month, we sort all stocks into terciles based on each limits-to-arbitrage variable X . We then independently sort stocks into quintiles based on the abnormal number of IPs searching for 10K files. In Table 7, we report the equal-weighted four-factor alpha of the lowest and highest AIP portfolio among the lowest and highest X group. Consistent with the limits to arbitrage predictions, the alpha of high-minus-low AIP_10K portfolio is more pronounced among stocks with higher idiosyncratic volatility, lower institutional ownership and less analyst coverage. For example, the high-minus-low AIP_10K portfolio generates 1.24% ($t=4.44$) monthly alpha among high volatility stocks, while only 0.23% ($t=1.76$) among low volatility stocks.

3.4 Variation in Complexity of Financial Filings

Our maintained hypothesis is that investors' costly information acquisition activity should be positively related to the expected payoff from using the information. If this is true, we would expect the payoff to be larger when information acquisition/processing cost is higher. To test this prediction, we use the complexity of firm's financial filings to proxy for cost of information acquisition/processing. The idea is intuitive, as more complex filings require more effort and time for investors to process and digest. Following the most recent literature (Loughran and McDonald (2014)), we use the size of firm's 10-K filing and the number of words contained in the filing to proxy for filing complexity.¹⁴

¹⁴Loughran and McDonald (2014) document that 10-K file size is positively associated with high return volatility in a one-month period following 10-K filings, supporting the use of file size as proxy for the linguistic complexity of 10-K disclosure.

To this end, we first get the filing size and number of words contained in firms' most recent 10-K report. As big firms have more business lines and diverse set of operations, they naturally have lengthy and larger 10-K filing.¹⁵ To remove the confounding effect of firm size, we regress the log of filing size and number of words on the log of firm's market capitalization, and use the regression residual as our proxy of filing complexity. At the end of each month, we sort all stocks into terciles based on either the residual file size or the residual word count. We then independently sort stocks into quintiles based on AIP_10K. In Table 8, we report the equal-weighted four-factor alpha of the lowest and highest AIP_10K portfolio among the highest and lowest group of filing complexity. Consistent with theories of endogenous information acquisition, the alpha of high-minus-low portfolio is indeed larger and more significant among firms with more complex financial filing. For example, the high-minus-low AIP_10K portfolio generates 0.92% (t=4.46) monthly alpha among firms with largest 10-K size, and 0.65% (t=3.51) among firms with small 10-K size. The result is similar when we use the word count in 10-K as proxy for disclosure complexity. Overall, the evidence provides strong support to our hypothesis that the more costly information acquisition/processing is, the larger the expected payoff revealed by information acquisition activity.

3.5 Fama-MacBeth Regression

We now test our main hypothesis using the Fama and MacBeth (1973) regression methodology. One advantage of this methodology is that it allows us to examine the predictive power of AIP while controlling for other known predictors of cross-sectional stock returns. This is important because, as shown in Table 1, AIP is correlated with some of these predictors. We conduct the Fama-MacBeth regressions in the usual way. Each month, starting in February 2003 and ending in December 2014, we run the following cross-sectional regression:

$$Ret_{i,t+1} = \beta_0 + \beta_1 AIP_{i,t} + \gamma X_{i,t} + \epsilon_{i,t} \quad (3)$$

¹⁵The rank correlation is 0.34 between 10-K file size and firm size, and 0.40 between word count and firm size.

where $Ret_{i,t+1}$ is return of stock i in month $t + 1$, $AIP_{i,t}$ is the abnormal number of IPs searching for firm i 's EDGAR filings in month t , and X is a set of control variables known to predict returns, including the natural logarithm of the book-to-market ratio (LnBM), the natural logarithm of the market value of equity (LnME), returns from the prior month (Rev), returns from the prior 12-month period excluding month $t-1$ (Mom), institutional ownership (IO), and idiosyncratic volatility (IVOL) and past 12-month turnover (Turnover12).

Table 9 reports the time-series averages of the coefficients on the independent variables, and the t -statistics are Newey-West adjusted with four lags to control for heteroskedasticity and autocorrelation. We report the results for AIP_total from Column (1) to (3), AIP_fundl from Column (4) to (6) and AIP_10K from Column (7) to (9). Column (1), (4) and (7) shows the coefficient on AIP without any other return predictors. The coefficients on all three AIP variables are positive and significant at 1% level. This is consistent with our portfolio sorting results in which stocks with abnormal large number of IPs searching its EDGAR filings have higher expected return. In Column (2), (5) and (8), we add the usual controls including size, book-to-market ratio, past 1-month returns, and past 12-month returns. The coefficients on AIP barely change and retains its strong predictive power. In Column (3), (6) and (9), we further add institutional ownership, turnover and idiosyncratic volatility in the regression, and AIP still positively predicts future returns. The economic magnitude is also quite large. The difference of AIP_10K between the lowest decile portfolio and highest decile portfolio is 2.39, which implies a monthly return spread of 105 basis points between these two extreme deciles. The magnitude estimated from Fama-MacBeth regression is in line with our portfolio sorting results. For the control variables, the sign of coefficients is consistent with previous literature, except for momentum, which attracts a negative coefficient.¹⁶ Due to the short and recent sample period, however, the coefficients on most control variables are not significantly different from zero.

EDGAR searching activity is positively related to scheduled firm events such as earnings announcement (Drake, Roulstone, and Thornock (2015)). Since earnings surprise leads to post-earnings announcement drift (Bernard and Thomas (1989)) and announcement months

¹⁶This is due to the 2009 momentum crash, see Daniel and Moskowitz (2016). The coefficient on Momentum becomes positive once we exclude year 2009 from our sample.

are generally associated with positive stock return (Lamont and Frazzini (2007)), the return predictability of AIP may be driven by these earnings-related return predictability effects. As a robustness check, we add standardized unexpected earnings (SUE) and an earnings announcement month dummy (EAM) in Fama-MacBeth regression. Column (1)-(3) of Table A3 shows that the coefficients on AIP become stronger after controlling for earnings-related variables, suggesting that the information contained in AIP is not driven by earnings-related return predictability effects.

Chen, Hong, and Stein (2002) show that reduction of breadth of institutional ownership is a proxy for overvaluation when short-sale constraint is binding for some investors. To the extent that breadth of ownership is positively correlated with the number of IPs searching EDGAR filings, our result may be rediscovery of their finding. Column (4)-(6) of Table A3 shows this is not the case. The coefficients on AIP barely change after controlling for change of breadth of ownership (dBreadth). The coefficient on change of breadth of ownership is positive but insignificant, probably due to the short sample period.

Having established that abnormal number of IPs predicts one-month-ahead returns, our next analysis examine the persistence of this predictive relation. This test could help rule out an alternative explanation, that the short-run predictability is due to temporary price pressure driven by investors' demand for attention-grabbing stocks. For example, Da, Engelberg, and Gao (2011) show that an increase in Google Search Volume for a stock predicts higher stock prices in the short-run that eventually reverse back within a year. As we hypothesize that AIP contains expected return information driven by firms' fundamental changes, the return predictability of AIP should not reverse over the long-run. To test this, we run Fama and MacBeth (1973) regression of cumulative returns from month $t + j$ to $t + k$ on the abnormal number of IPs searching 10-K filings in EDGAR database (AIP_10K) at month t . The result is report in Table 10. We separately show the return predictability of AIP_10K for the next quarter return skipping the immediate month in Column (1), the second quarter return in Column (2), the second half-year return in Column (3) and the second year return in Column (4). The table shows that lagged value of AIP also significantly predicts returns for up to 2 quarters, and eventually levels off for longer horizon. The coefficient on AIP

is always positive and never reverse, mitigating concerns that the predictive power of AIP comes from transitory price pressure that reverses subsequently. Investors searching firm fundamentals through EDGAR system appear to be more sophisticated than those searching through Google Search Engine, and their aggregate information acquisition activities contains value-relevant information about firms that slowly diffuse into stock prices.

3.6 Which Types of EDGAR Filings Matter Most?

Given the high correlation between the three types of IP measures as shown in Table 1, we next examine whether the expected return information embedded in the three AIP variables are incremental to each other. To test this, we run a horse race by including all three AIP variables in the Fama-MacBeth regression. The result is reported in Table 11. Column (1) reports the result without other controls, and Column (2) includes all the usual return predictors. The results show clearly that the return predictability of AIP comes mainly from those searching for firms' annual report 10K. While AIP_10K retains its strong predictive power, the coefficient on AIP_total and AIP_fundl becomes insignificant. Acquiring and analyzing 10K report is more costly than other SEC filings and more reflect deliberate information acquisition behavior. The result is thus consistent with our hypothesis that costly information acquisition activity contains expected benefits from utilizing such information.

3.7 IPs or Searches?

Our measure of information acquisition activity essentially equal weights each investor searching through EDGAR regardless of the number of searches they requested through EDGAR system during one month window. An alternative measure of information acquisition activity is the total number of searches on a firm requested by investors through EDGAR system. This measure is problematic because, as documented by Drake, Roulstone, and Thornock (2015), the number of requests through EDGAR is dominated by a small fraction of investors who access EDGAR very frequently, and their activities are over-

represented in this alternative measure.¹⁷ Under the assumption that information is more likely dispersed among a large group of market participants (Hayek (1945)), we think that our measure of abnormal number of IPs should be more powerful in terms of backing out the latent information embedded in "the wisdom of crowd". Nevertheless, to test which measure of information acquisition activity has stronger return predictability, we conduct a horse race between the abnormal number of searches (Asearch) and abnormal number of IPs (AIP) using Fama-MacBeth regression approach. Using the same decomposition method, we extract the abnormal number of searches for each firm as the residual from a monthly regression of log one plus raw number of EDGAR requests for SEC filings on the same set of firm characteristics used in equation (1).

The result is reported in Table 12. We look at searches/IPs for all types of EDGAR files in Column (1) and (2), 10K, 10Q and 8K in Column (3) and (4), and searches for annual reports only in Column (5) and (6). Column (1), (3) and (5) show that the return predictability of Asearch is generally positive but weaker compared to that of AIP. Column (2), (4) and (6) show that once we control for AIP, the coefficient on Asearch is no longer significant and even changes sign. Importantly, the coefficients on AIP are still positive and highly significant. The result supports our use of number of IPs as a cleaner measure of aggregate information acquisition activity, and indirectly supports the underlying assumption that private information is dispersed among market participants.

4 Channels

The key hypothesis in this paper is that information acquisition activity embeds expected return information because investors rationally expend greater effort to analyzing firms that are underpriced with large price appreciation potential. As mispricing implies the separation of stock prices from firms' fundamental value, there are two non-mutually exclusive channels through which investors can identify mispricing. The first channel is that investors' costly information acquisition contains their favorable expectation of firms' fundamental performance

¹⁷Drake, Roulstone, and Thornock (2015) document that 86% of the users accessing EDGAR do so infrequently and only around 2% of the users access EDGAR actively during a given quarter.

that are not fully priced in by market. A second channel is that investors identify mispricing by observing changes in stock prices that are unwarranted by firms' fundamental changes. In this section, we test both channels.

4.1 Predicting Fundamental Performance

We first test whether information acquisition via EDGAR reveals novel information about firms' fundamental performance change. We use two measures of firms' fundamental performance. The first one is the change of quarterly Return-on-Asset (dROA) from four quarters ago, which takes into account of the seasonality of firms' operating performance. The second measure is the monthly forecast revision of analysts' consensus Earnings-per-Share (EPS) forecast (FREV) scaled by stock prices 12 months ago, as a higher frequency measure of firms' fundamental performance. We run panel regression of dROA and FREV on lagged AIP, controlling for other firm characteristics that are correlated with firms' fundamental performance, including size, book-to-market, past 12-month return, analyst coverage, turnover, institutional ownership and idiosyncratic volatility and lagged quarterly ROA. Since quarterly Return-on-Assets is measured at quarterly frequency, we calculate the AIP at quarterly frequency as the monthly AIP averaged within a quarter. We also control for time fixed effect and standard errors are double clustered by firm and time following Petersen (2009). If the return predictability of AIP is partially driven by its predictive power for firm fundamentals, the coefficient on AIP should be significantly positive.

Table 13 reports the results of predicting fundamental performance based on AIP. The dependent variable is the change of quarterly ROA from Column (1) to Column (3), and analyst forecast revision from Column (4) to Column (6). We show the predictability result for all three AIP measures. The coefficients on AIP are significantly positive for both measures of fundamental performance, regardless of which AIP measures we use. The economic magnitude is non-trivial. For example, Column (3) shows that an interquartile increase in AIP_10K is associated with an increase of 0.22 percentage points in dROA, which is about 17% of the interquartile range of quarterly change of ROA. This finding suggests that information acquisition via EDGAR contains investors' expectation about firms' future operating

performance and even leads analysts’ revision of their forecast of firms’ fundamental. It is worth pointing out that the predictability of AIP is obtained after controlling for other determinants of firms’ fundamental performance. For example, past 12-month return strongly positively predict change of ROA and analyst forecast revision, while turnover and idiosyncratic volatility negatively predict fundamental performance. Overall, the test supports our hypothesis that the source of return predictability comes from investors allocating greater effort towards firms with improving fundamentals.

4.2 Underpricing Driven by Outflow-induced Fire Sale

A second channel through which mispricing could occur is exogenous shock to stock prices that are unwarranted by fundamentals. One such example is index addition event, as Shleifer (1986), Wurgler and Zhuravskaya (2002), and Chang, Hong, and Liskovich (2014) show that forced buying from index-tracking institutional investors around such events could lead to large price pressure on affected stocks. However, index addition events are rare, which limits its applicability in our setting. In this paper, we use mutual fund outflow-induced fire sell as an exogenous shock to stock price. Coval and Stafford (2007), Khan, Kogan, and Serafeim (2012) and Edmans, Goldstein, and Jiang (2012) find that mutual funds sell firm’s shares roughly in proportion to its portfolio weights when facing severe outflows. The forced selling behavior results in significant downward price pressure that persists for more than a year. This is a relatively exogenous and clean measure of underpricing as it is associated with who is selling—funds facing large investor redemptions—rather than what is being sold, so it is unlikely driven by (unobserved) changes in firms’ fundamental performance.

To test whether investors expend more effort and time towards firms experiencing fire sell-induced underpricing, we examine change of abnormal number of IPs following flow-induced fire sale. Specifically, we run the following Fama and MacBeth (1973) regression:

$$dAIP_{i,q+1} = \beta_0 + \beta_1 Outflow_{i,q} + \beta_2 X_{i,q} + \epsilon_{i,q+1} \quad (4)$$

where $Outflow_{i,q}$ is the flow-induced fire sale measure calculated following Edmans, Goldstein, and Jiang (2012), which reflects fund outflow expressed as a percentage of firms’ shares

outstanding. Our dependent variable $dAIP_{i,q+1}$ is the within-firm change of AIP in quarter $q+1$ following mutual fund outflows. X is a set of firm characteristics that may affect change of AIP.

Table 14 reports the result. Again we show the result for all three AIP measures. Column (1), (3) and (5) show the coefficients on "Outflows" are significantly negative without other controls, for all three IP measures. The negative coefficient means that more investors are searching for the EDGAR filing of firms that are underpriced due to exogenous reasons. Column (2), (4) and (6) show that the negative relation between outflows-induced selling pressure and change in AIP is robust after controlling for firms' size, book-to-market ratio, analyst coverage, idiosyncratic volatility, turnover, institutional ownership and past returns, suggesting our findings are likely driven by variation in underpricing.

In sum, by using mutual fund outflow-induced selling pressure to identify stock-level underpricing, our test also supports the second channel that part of the return predictability we document is attributable to investors allocating more attention and resources towards firms experiencing exogenous sources of undervaluation that deviates from firm fundamentals.

4.3 Information Acquisition and Institutional Trading

As investors' information acquisition through EDGAR contain value-relevant information about stocks, a natural question that emerges is who are these sophisticated investors? Although we don't have the identify of those searching through EDGAR system, hedge fund managers appear to fit the profile of informed investors in the equity market. A growing literature shows that hedge funds possess stock picking skills and are able to identify stock-level mispricing (Brunnermeier and Nagel (2004); Jiao, Massa, and Zhang (2016); Agarwal, Jiang, Tang, and Yang (2013)). If hedge fund managers rationally allocate more resources and efforts towards firms with improving fundamentals, they should also trade in the direction of the latent information indicated by EDGAR search traffic.

To examine whether abnormal number of IPs predict hedge fund trades, we run Fama-MacBeth regression of net purchases by hedge funds and mutual funds in quarter q on lagged AIP, and control for other stock characteristics that might influence fund trading decisions.

Net purchase is measured as quarterly change of hedge fund holding on a stock, where holding is expressed as a fraction of firm’s shares outstanding. Since hedge fund trades are inferred from quarterly holding reports, we calculate the AIP at quarterly frequency as the monthly AIP averaged within a quarter. Specifically, in each quarter, we run the following cross-sectional regression:

$$NetBuy_{i,q} = \beta_0 + \beta_1 AIP_{i,q-1} + \beta_2 Hold_{i,q-1} + \gamma X_{i,q-1} + \epsilon_{i,q} \quad (5)$$

where $NetBuy_{i,q}$ is either the net purchases by hedge funds or those by mutual funds in quarter q , $AIP_{i,q-1}$ is abnormal number of IPs searching for firm i ’s EDGAR filings in quarter $q-1$, $Hold_{i,q-1}$ is either the hedge fund ownership or mutual fund ownership at end of quarter $q-1$, and $X_{i,q-1}$ is a vector of firm characteristics at end of quarter $q-1$, including firm size, book-to-market, analyst coverage, volatility, turnover, institutional ownership and momentum. If hedge funds contribute to informed searches through EDGAR, the coefficient on AIP should be significant and positive.

Table 15 reports the time-series averages of the cross-sectional regression coefficients. The dependent variable from Column (1) to Column (3) is net buying by hedge funds. The coefficient on AIP in the regression of hedge funds’ net purchases is positive and significant for all three measures of AIP. In terms of economic magnitude, one interquartile increase in AIP_10K is associated with an increase of 0.26 percentage points in hedge funds’ net purchases, which is about 25% of the interquartile range of net buying by hedge funds. The economic magnitude is reasonable given that not all hedge funds are fundamental investors and they also have other information sources to aid their investment decisions. In contrast, Column (4) to (6) show that abnormal number of IPs searching EDGAR filing does not significantly predict mutual funds’ net purchase.

Overall, the evidence suggests that either hedge funds are part of these sophisticated investors making informed searches through EDGAR system, or their own information source is consistent with the latent information embedded in ”the wisdom of crowds”. Either interpretation would provide support to our premise in this paper that investors’ information acquisition activity reveals their expected benefits of trading on such information.

5 Conclusion

In this paper, we examine the expected return information contained in investors' costly information acquisition activities. Specifically, we use a novel dataset of investors' requests to company filings through EDGAR system to back out their expectations over future payoffs. To this end, we develop and implement a simple characteristic-based model to decompose the total number of IPs searching for EDGAR filings into an abnormal and expected components and show that the abnormal number of IPs searching for firms' financial reports positively predict subsequent stock returns. A long-short portfolio that buys stocks with abnormal number of IPs in the top decile and sells stocks in the bottom decile generates an equal-weighted monthly four-factor alpha of up to 82 basis points that does not reverse in the long run. We also find that abnormal number of IPs predicts firms' ascending fundamental performance and also increases following exogenous underpricing, suggesting investors rationally allocate greater resources and effort towards firms with large price appreciation potential. Lastly, information acquisition via EDGAR also predicts subsequent purchases by hedge fund managers, suggesting that sophisticated investors are making informed searches on firms with largest potential payoffs.

Taken together, our findings provide empirical support to theoretical models of endogenous information acquisition that costly information acquisition activity is positively associated with the value of information (Grossman and Stiglitz (1980)). Our research also highlights the promise of using the collective wisdom of investors—extracted from their EDGAR search behavior—to study expected returns and other important economic outcomes.

References

- Admati, A. R., 1985, “A noisy rational expectations equilibrium for multi-asset securities markets,” *Econometrica: Journal of the Econometric Society*, pp. 629–657.
- Agarwal, V., W. Jiang, Y. Tang, and B. Yang, 2013, “Uncovering hedge fund skill from the portfolio holdings they hide,” *The Journal of Finance*, 68(2), 739–783.
- Amihud, Y., 2002, “Illiquidity and stock returns: cross-section and time-series effects,” *Journal of Financial Markets*, 5(1), 31–56.
- Ang, A., R. J. Hodrick, Y. Xing, and X. Zhang, 2006, “The cross-section of volatility and expected returns,” *The Journal of Finance*, 61(1), 259–299.
- Asparouhova, E., H. Bessembinder, and I. Kalcheva, 2013, “Noisy prices and inference regarding returns,” *The Journal of Finance*, 68(2), 665–714.
- Barber, B. M., and T. Odean, 2007, “All that glitters: The effect of attention and news on the buying behavior of individual and institutional investors,” *The Review of Financial Studies*, 21(2), 785–818.
- Ben-Rephael, A., Z. Da, and R. D. Israelsen, 2017, “It Depends on Where You Search: Institutional Investor Attention and Underreaction to News,” *The Review of Financial Studies*, p. hhx031.
- Bernard, V. L., and J. K. Thomas, 1989, “Post-earnings-announcement drift: delayed price response or risk premium?,” *Journal of Accounting research*, pp. 1–36.
- Brunnermeier, M. K., and S. Nagel, 2004, “Hedge Funds and the Technology Bubble,” *Journal of Finance*, pp. 2013–2040.
- Carhart, M. M., 1997, “On persistence in mutual fund performance,” *The Journal of finance*, 52(1), 57–82.
- Chang, Y.-C., H. Hong, and I. Liskovich, 2014, “Regression discontinuity and the price effects of stock market indexing,” *The Review of Financial Studies*, 28(1), 212–246.
- Chen, H., P. De, Y. Hu, and B.-H. Hwang, 2014, “Wisdom of crowds: The value of stock opinions transmitted through social media,” *The Review of Financial Studies*, 27(5), 1367–1403.

- Chen, J., H. Hong, and J. C. Stein, 2002, “Breadth of ownership and stock returns,” *Journal of Financial Economics*, 66(2), 171–205.
- Chordia, T., A. Subrahmanyam, and Q. Tong, 2014, “Have capital market anomalies attenuated in the recent era of high liquidity and trading activity?,” *Journal of Accounting and Economics*, 58(1), 41–58.
- Coval, J., and E. Stafford, 2007, “Asset fire sales (and purchases) in equity markets,” *Journal of Financial Economics*, 86(2), 479–512.
- Da, Z., J. Engelberg, and P. Gao, 2011, “In search of attention,” *The Journal of Finance*, 66(5), 1461–1499.
- Daniel, K., M. Grinblatt, S. Titman, and R. Wermers, 1997, “Measuring mutual fund performance with characteristic-based benchmarks,” *The Journal of Finance*, 52(3), 1035–1058.
- Daniel, K., and T. J. Moskowitz, 2016, “Momentum crashes,” *Journal of Financial Economics*, 122(2), 221–247.
- Diamond, D. W., and R. E. Verrecchia, 1981, “Information aggregation in a noisy rational expectations economy,” *Journal of Financial Economics*, 9(3), 221–235.
- Drake, M. S., P. J. Quinn, and J. R. Thornock, 2017, “Who Uses Financial Statements? A Demographic Analysis of Financial Statement Downloads from EDGAR,” *Accounting Horizons*.
- Drake, M. S., D. T. Roulstone, and J. R. Thornock, 2012, “Investor information demand: Evidence from Google searches around earnings announcements,” *Journal of Accounting Research*, 50(4), 1001–1040.
- , 2015, “The determinants and consequences of information acquisition via EDGAR,” *Contemporary Accounting Research*, 32(3), 1128–1161.
- , 2016, “The usefulness of historical accounting reports,” *Journal of Accounting and Economics*, 61(2), 448–464.
- Du, Z., 2015, “Endogenous Information Acquisition: Evidence from Web Visits to SEC Filings of Insider Trades,” working paper, Working paper, Kellogg School of Management.
- Edmans, A., I. Goldstein, and W. Jiang, 2012, “The real effects of financial markets: The

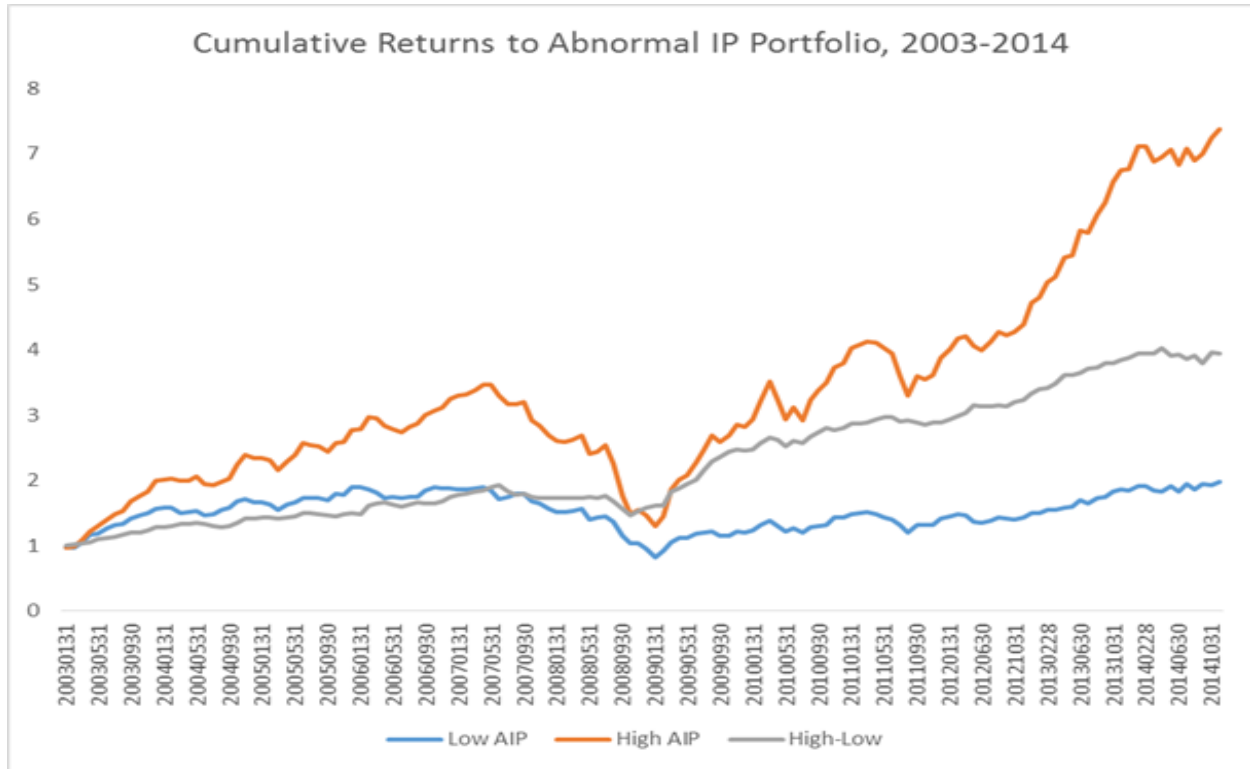
- impact of prices on takeovers,” *The Journal of Finance*, 67(3), 933–971.
- Fama, E. F., and K. R. French, 1993, “Common risk factors in the returns on stocks and bonds,” *Journal of Financial Economics*, 33(1), 3–56.
- , 2006, “Profitability, investment and average returns,” *Journal of Financial Economics*, 82(3), 491–518.
- , 2016, “Dissecting anomalies with a five-factor model,” *Review of Financial Studies*, 29(1), 69–103.
- Fama, E. F., and J. D. MacBeth, 1973, “Risk, return, and equilibrium: Empirical tests,” *The Journal of Political Economy*, pp. 607–636.
- Frazzini, A., R. Israel, and T. J. Moskowitz, 2012, “Trading costs of asset pricing anomalies,” *Fama-Miller Working Paper*, pp. 14–05.
- Gao, X., and J. R. Ritter, 2010, “The marketing of seasoned equity offerings,” *Journal of Financial Economics*, 97(1), 33–52.
- Gervais, S., R. Kaniel, and D. H. Mingelgrin, 2001, “The high-volume return premium,” *The Journal of Finance*, 56(3), 877–919.
- Green, J., J. R. Hand, and X. F. Zhang, 2017, “The characteristics that provide independent information about average us monthly stock returns,” *The Review of Financial Studies*, p. hhx019.
- Green, T. C., R. Huang, Q. Wen, and D. Zhou, 2017, “Wisdom of the Employee Crowd: Employer Reviews and Stock Returns,” .
- Grossman, S. J., and J. E. Stiglitz, 1980, “On the impossibility of informationally efficient markets,” *The American economic review*, 70(3), 393–408.
- Hayek, F. A., 1945, “The use of knowledge in society,” *The American economic review*, pp. 519–530.
- Hellwig, M. F., 1980, “On the aggregation of information in competitive markets,” *Journal of economic theory*, 22(3), 477–498.
- Hong, H., T. Lim, and J. C. Stein, 2000, “Bad news travels slowly: Size, analyst coverage, and the profitability of momentum strategies,” *The Journal of Finance*, 55(1), 265–295.

- Hou, K., C. Xue, and L. Zhang, 2015, “Digesting Anomalies: An Investment Approach,” *Review of Financial Studies*, 28(3), 650–705.
- Huang, J., 2016, “The customer knows best: The investment value of consumer opinions,” *Browser Download This Paper*.
- Jegadeesh, N., and S. Titman, 1993, “Returns to buying winners and selling losers: Implications for stock market efficiency,” *The Journal of Finance*, 48(1), 65–91.
- Jiang, W., 2014, “Leveraged speculators and asset prices,” .
- Jiao, Y., M. Massa, and H. Zhang, 2016, “Short selling meets hedge fund 13F: An anatomy of informed demand,” *Journal of Financial Economics*, 122(3), 544–567.
- Ke, B., and Y. Yu, 2006, “The effect of issuing biased earnings forecasts on analysts’ access to management and survival,” *Journal of Accounting Research*, 44(5), 965–999.
- Khan, M., L. Kogan, and G. Serafeim, 2012, “Mutual fund trading pressure: Firm-level stock price impact and timing of SEOs,” *The Journal of Finance*, 67(4), 1371–1395.
- Lamont, O., and A. Frazzini, 2007, “The earnings announcement premium and trading volume,” working paper, National Bureau of Economic Research.
- Lee, C. M., P. Ma, and C. C. Wang, 2015, “Search-based peer firms: Aggregating investor perceptions through internet co-searches,” *Journal of Financial Economics*, 116(2), 410–431.
- Lee, C. M., and E. C. So, 2017, “Uncovering expected returns: Information in analyst coverage proxies,” *Journal of Financial Economics*, 124(2), 331–348.
- Loughran, T., and B. McDonald, 2011, “When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks,” *The Journal of Finance*, 66(1), 35–65.
- , 2014, “Measuring readability in financial disclosures,” *The Journal of Finance*, 69(4), 1643–1671.
- , 2017, “The use of EDGAR filings by investors,” *Journal of Behavioral Finance*, 18(2), 231–248.
- Mele, A., and F. Sangiorgi, 2015, “Uncertainty, information acquisition, and price swings in asset markets,” *The Review of Economic Studies*, 82(4), 1533–1567.

- Merton, R. C., 1987, “A simple model of capital market equilibrium with incomplete information,” *The journal of finance*, 42(3), 483–510.
- Monga, V., and E. Chasan, 2015, “The 109,894-Word Annual Report: As Regulators Require More Disclosures, 10-Ks Reach Epic Lengths; How Much Is Too Much?,” *Wall Street Journal*.
- Nagel, S., 2005, “Short sales, institutional investors and the cross-section of stock returns,” *Journal of Financial Economics*, 78(2), 277–309.
- Newey, W. K., and K. D. West, 1987, “A Simple, Positive Semi-Definite, Heteroskedasticity and Autocorrelation Consistent Covariance Matrix,” *Econometrica*, pp. 703–708.
- Pástor, L., and R. F. Stambaugh, 2003, “Liquidity Risk and Expected Stock Returns,” *Journal of Political Economy*, 111(3), 642–685.
- Pedersen, L. H., 2015, *Efficiently inefficient: how smart money invests and market prices are determined*. Princeton University Press.
- Petersen, M. A., 2009, “Estimating standard errors in finance panel data sets: Comparing approaches,” *Review of Financial Studies*, 22(1), 435–480.
- Pontiff, J., 2006, “Costly arbitrage and the myth of idiosyncratic risk,” *Journal of Accounting and Economics*, 42(1), 35–52.
- Ryans, J. P., 2017, “Using the EDGAR Log File Data Set,” .
- Shleifer, A., 1986, “Do demand curves for stocks slope down?,” *The Journal of Finance*, 41(3), 579–590.
- Shumway, T., 1997, “The delisting bias in CRSP data,” *The Journal of Finance*, 52(1), 327–340.
- Stambaugh, R. F., J. Yu, and Y. Yuan, 2015, “Arbitrage asymmetry and the idiosyncratic volatility puzzle,” *The Journal of Finance*.
- Stambaugh, R. F., and Y. Yuan, 2016, “Mispricing factors,” *The Review of Financial Studies*, 30(4), 1270–1315.
- Verrecchia, R. E., 1982, “Information acquisition in a noisy rational expectations economy,” *Econometrica: Journal of the Econometric Society*, pp. 1415–1430.

- Wurgler, J., and E. Zhuravskaya, 2002, “Does arbitrage flatten demand curves for stocks?,” *The Journal of Business*, 75(4), 583–608.
- Zhang, X., 2006, “Information uncertainty and stock returns,” *The Journal of Finance*, 61(1), 105–137.

Figure 1: Cumulative Returns to AIP strategy



This figure shows the cumulative equal-weighted returns to the lowest and highest decile portfolios sorted on abnormal number of IPs searching for 10-K files in EDGAR system (AIP_10K). Grey line represents the cumulative returns to the top-minus-bottom portfolio formed on AIP_10K. The sample runs from January 2003 to December 2014.

Table 1: **Stock-Level Descriptive Statistics**

This table presents the descriptive statistics of our variables. Panel A reports the summary statistics for the full sample. Panel B reports the pairwise rank correlation among our variables where they overlap. Panel C reports the characteristics of portfolios sorted on abnormal number of IPs searching for 10-K files in SEC's EDGAR database (AIP_10K). IP_total is the total number of unique IP address searching all six types of EDGAR filings. IP_funtl is the total number of unique IP address searching 10-K, 10-Q and 8-K files. AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10K files in EDGAR database. Each month, we sort all stocks into deciles based on their AIP_10K. We first calculate the mean of each variable for each decile each month and then calculate the time-series average of cross-sectional means. LnME is the natural log of firm's market capitalization at the end of the June of each year in millions of US dollars. Coverage is log one plus analyst coverage. Turnover12 is the monthly turnover ratio averaged over the past 12 months. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. The overall sample period is from January 2003 to December 2014.

Panel A: Summary Statistics					
Variable	Mean	Median	STD	P25	P75
	<i>Number of IPs searching for EDGAR filings</i>				
IP_total	155	94	317	56	159
IP_funtl	107	64	213	37	111
IP_10K	60	32	135	17	60
IP_10Q	37	24	61	13	42
IP_8K	33	19	79	10	36
	<i>Firm-level characteristics</i>				
LnME	6.16	6.08	1.98	4.74	7.47
LnBM	-0.66	-0.56	0.84	-1.11	-0.12
Mom	16.67%	7.64%	57.57%	-12.06%	31.78%
Coverage	1.49	1.59	1.01	0.59	2.30
IVOL	0.02	0.02	0.02	0.01	0.03
Turnover12	0.17	0.12	0.19	0.05	0.21
IO	55.30%	59.15%	31.41%	28.92%	80.58%
	<i>Firm fundamentals and institutional trades</i>				
dROA (%)	0.032	-0.018	4.844	-0.684	0.599
FREV (%)	-0.106	-0.001	22.185	-0.070	0.052
Net buying by HFs (%)	0.102	-0.002	2.106	-0.475	0.582
Net buying by MFs (%)	0.266	0.093	2.649	-0.640	1.204

Table 1 Continued

Panel B: Rank Correlations										
	IP_total	IP_funtl	IP_10K	LnME	Cov	Turnover12	Ivol	LnBM	Mom	IO
IP_total	1.000									
IP_funtl	0.918	1.000								
IP_10K	0.812	0.897	1.000							
LnME	0.671	0.664	0.672	1.000						
Cov	0.594	0.605	0.603	0.832	1.000					
Turnover12	0.588	0.579	0.539	0.544	0.621	1.000				
Ivol	-0.134	-0.149	-0.212	-0.523	-0.360	-0.016	1.000			
LnBM	-0.239	-0.229	-0.224	-0.319	-0.326	-0.303	0.051	1.000		
Mom	0.031	0.023	0.044	0.112	0.051	0.049	-0.117	0.008	1.000	
IO	0.469	0.494	0.514	0.650	0.647	0.615	-0.306	-0.193	0.095	1.000

Table 1 Continued

Panel C: Descriptive statistics by AIP_10K deciles												
	Obs	AIP_10K	IP_total	IP_funtl	IP_10K	LnME	Cov	Turnover12	Ivol	LnBM	Mom	IO
1(Low)	330	-1.25	59	35	12	5.977	1.369	0.154	0.025	-0.590	0.150	45.53%
2	330	-0.60	76	51	22	6.074	1.513	0.163	0.024	-0.719	0.164	53.38%
3	330	-0.38	91	63	30	6.166	1.573	0.166	0.024	-0.742	0.163	57.21%
4	330	-0.21	104	72	36	6.248	1.611	0.170	0.024	-0.741	0.172	59.23%
5	330	-0.07	116	82	42	6.270	1.623	0.171	0.024	-0.711	0.176	60.20%
6	330	0.07	128	91	48	6.284	1.634	0.170	0.024	-0.700	0.173	60.79%
7	330	0.22	141	101	55	6.218	1.594	0.165	0.024	-0.662	0.174	60.19%
8	330	0.39	160	116	66	6.118	1.526	0.164	0.024	-0.623	0.164	58.91%
9	330	0.62	201	147	87	6.032	1.454	0.158	0.025	-0.563	0.162	56.09%
10(High)	330	1.14	464	342	226	6.257	1.483	0.163	0.025	-0.537	0.168	53.28%

Table 2: Cross-Sectional Determinants of Number of IPs Searching EDGAR Filings

This table presents the Fama-MacBeth regression of log number of IPs searching SEC Edgar files. In Panel A, the dependent variable is the log one plus number of unique IP addresses searching EDGAR filings in a month. In Panel B, the dependent variable is the log one plus number of unique IP addresses searching EDGAR 10-K, 10-Q and 8-K files in a month. In Panel C, the dependent variable is the log one plus number of unique IP addresses searching 10-K files in a month. LnME is the natural log of firm's market capitalization at the end of the June of each year in millions of US dollars. Coverage is log one plus analyst coverage. Turnover12 is the average monthly turnover ratio over the past 12 months. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. The overall sample period is from January 2003 to December 2014.

Panel A: Dependent Variable is log(1+# of unique IP addresses searching all EDGAR filings)							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
LnME	0.2713*** (69.44)	0.2356*** (71.54)	0.2475*** (73.46)	0.2943*** (75.60)	0.2992*** (76.94)	0.3015*** (77.29)	0.3026*** (77.58)
Coverage		0.1310*** (32.65)	0.0422*** (14.39)	0.0382*** (14.36)	0.0321*** (12.17)	0.0332*** (12.56)	0.0360*** (14.17)
Turnover12			1.0083*** (30.21)	0.7934*** (29.08)	0.7862*** (30.04)	0.7912*** (29.75)	0.7877*** (30.52)
Ivol				9.1266*** (34.65)	9.0159*** (33.38)	9.0510*** (33.16)	9.0215*** (32.36)
Mom					-0.0518*** (-6.00)	-0.0529*** (-6.19)	-0.0507*** (-5.99)
LnBM						0.0171*** (8.19)	0.0158*** (7.25)
IO							-0.0299** (-1.99)
Constant	2.5352*** (39.20)	2.6342*** (40.68)	2.5357*** (40.19)	2.0730*** (33.45)	2.0483*** (33.37)	2.0408*** (33.32)	2.0449*** (32.62)
Ave.R-sq	0.404	0.483	0.520	0.554	0.558	0.559	0.563
N.of Obs.	610651	488129	488129	488123	488123	488123	484835

Table 2 Continued

Panel B: Dependent Variable is log(1+# of unique IP addresses searching 10-K, 10-Q and 8-K files)							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
LnME	0.2723*** (64.05)	0.2355*** (61.21)	0.2468*** (60.97)	0.2931*** (65.17)	0.2984*** (66.87)	0.3015*** (67.27)	0.3005*** (67.10)
Coverage		0.1405*** (35.59)	0.0530*** (15.60)	0.0492*** (15.87)	0.0421*** (13.86)	0.0436*** (14.48)	0.0369*** (15.57)
Turnover12			0.9833*** (29.18)	0.7702*** (26.62)	0.7708*** (27.23)	0.7787*** (26.95)	0.7560*** (27.32)
Ivol				9.0866*** (36.40)	8.9652*** (34.66)	9.0334*** (34.19)	9.0934*** (33.54)
Mom					-0.0684*** (-7.72)	-0.0698*** (-8.00)	-0.0685*** (-7.95)
LnBM						0.0251*** (10.23)	0.0223*** (9.01)
IO							0.0411*** (2.76)
Constant	2.2017*** (34.86)	2.2804*** (36.21)	2.1866*** (35.81)	1.7281*** (28.85)	1.7033*** (28.72)	1.6943*** (28.66)	1.6868*** (27.97)
Ave.R-sq	0.386	0.458	0.491	0.522	0.526	0.527	0.533
N.of Obs.	610651	488129	488129	488123	488123	488123	484835

Table 2 Continued

Panel C: Dependent Variable is log(1+# of unique IP addresses searching 10-K files)							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
LnME	0.2979*** (61.33)	0.2674*** (60.72)	0.2765*** (59.37)	0.3120*** (61.48)	0.3169*** (62.64)	0.3201*** (63.24)	0.3155*** (62.55)
Coverage		0.1453*** (35.85)	0.0729*** (23.41)	0.0698*** (23.28)	0.0637*** (21.42)	0.0649*** (21.49)	0.0431*** (16.48)
Turnover12			0.8122*** (30.68)	0.6461*** (28.59)	0.6415*** (28.59)	0.6522*** (28.74)	0.5924*** (28.38)
Ivol				6.9981*** (30.56)	6.9145*** (29.46)	7.0130*** (28.94)	7.2542*** (29.41)
Mom					-0.0484*** (-5.54)	-0.0510*** (-5.93)	-0.0521*** (-6.09)
LnBM						0.0267*** (9.03)	0.0213*** (7.48)
IO							0.1600*** (10.36)
Constant	1.3873*** (25.17)	1.4159*** (25.47)	1.3396*** (24.67)	0.9886*** (18.65)	0.9639*** (18.51)	0.9554*** (18.43)	0.9267*** (17.62)
Ave.R-sq	0.388	0.467	0.486	0.501	0.504	0.506	0.511
N.of Obs.	610651	488129	488129	488123	488123	488123	484835

Table 3: **Portfolio Excess Returns Sorted on Abnormal Number of IPs**

This table reports the monthly average excess returns for each of the decile portfolios, as well as the long-short portfolio (High-Low). AIP_total is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of filings in EDGAR database on a set of firm characteristics. Similarly, AIP_funtl (AIP_10K) is constructed using number of IPs searching for 10-K, 10-Q and 8-K files (10-K) in EDGAR database. At the end of each month, all stocks are sorted into deciles based on their abnormal number of IPs and a long-short portfolio is formed by buying the highest decile and shorting the lowest decile portfolio. Portfolio returns are computed over the next month. Panel A reports results for equally weighted portfolios and Panel B shows results for value-weighted portfolios. The sample runs from January 2003 to December 2014.

Panel A: Equal-weighted Decile Portfolio Excess Return						
	AIP_total	t-stat	AIP_funtl	t-stat	AIP_10K	t-stat
Low	0.46	1.20	0.50	1.29	0.47	1.22
2	0.78	1.78	0.76	1.73	0.63	1.40
3	0.81	1.83	0.80	1.79	0.75	1.68
4	1.08	2.33	1.04	2.27	0.85	1.81
5	1.00	2.15	1.00	2.13	0.93	1.99
6	1.07	2.24	0.99	2.07	1.02	2.11
7	1.19	2.40	1.14	2.34	1.11	2.28
8	1.12	2.19	1.06	2.05	1.26	2.51
9	1.14	2.21	1.24	2.35	1.32	2.54
High	1.18	2.29	1.29	2.55	1.48	2.98
High - Low	0.71	3.18	0.79	3.61	1.00	4.70

Panel B: Value-weighted Decile Portfolio Excess Return						
	AIP_total	t-stat	AIP_funtl	t-stat	AIP_10K	t-stat
Low	0.40	1.01	0.57	1.60	0.48	1.42
2	0.80	2.01	0.72	1.72	0.59	1.39
3	0.76	1.93	0.86	2.15	0.68	1.61
4	1.04	2.58	0.97	2.35	0.83	2.03
5	0.85	2.09	0.92	2.20	0.99	2.54
6	0.80	2.03	0.89	2.23	0.75	1.83
7	1.00	2.62	0.90	2.38	0.88	2.18
8	0.89	2.26	0.84	2.13	1.01	2.70
9	0.94	2.60	0.87	2.43	0.74	2.04
High	0.71	2.13	0.66	2.01	0.75	2.28
High - Low	0.31	1.23	0.09	0.44	0.26	1.32

Table 4: **Factor adjusted alphas of Portfolios Sorted on Abnormal Number of IPs**

This table reports the monthly Carhart (1997) four factor alphas for each of the 10 decile portfolios, as well as the long-short portfolio (High-Low). AIP_total is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of files in EDGAR database on a set of firm characteristics. Similarly, AIP_funtl (AIP_10K) is constructed using number of IPs searching for 10-K, 10-Q and 8-K filings (10-K) in EDGAR database. At the end of each month, all stocks are sorted into deciles based on their abnormal number of IPs and a long-short portfolio is formed by buying the highest decile and shorting the lowest decile portfolio. Portfolio returns are computed over the next month. Panel A reports results for equally weighted portfolios and Panel B shows results for value-weighted portfolios. The sample runs from January 2003 to December 2014.

Panel A: Equal-weighted Decile Portfolio 4-factor alpha						
	AIP_total	t-stat	AIP_funtl	t-stat	AIP_10K	t-stat
Low	-0.36	-3.40	-0.32	-2.98	-0.34	-2.84
2	-0.16	-1.77	-0.19	-2.13	-0.33	-3.68
3	-0.15	-2.02	-0.16	-1.79	-0.22	-2.33
4	0.07	0.74	0.04	0.44	-0.18	-2.13
5	-0.01	-0.15	-0.02	-0.20	-0.09	-1.05
6	0.04	0.46	-0.04	-0.50	-0.02	-0.21
7	0.14	1.33	0.11	0.89	0.08	0.63
8	0.06	0.47	-0.01	-0.10	0.20	1.80
9	0.08	0.56	0.16	1.20	0.27	1.85
High	0.16	0.94	0.29	1.87	0.48	3.30
High - Low	0.52	2.74	0.62	3.33	0.82	4.35

Panel B: Value-weighted Decile Portfolio 4-factor alpha						
	AIP_total	t-stat	AIP_funtl	t-stat	AIP_10K	t-stat
Low	-0.40	-2.15	-0.17	-1.05	-0.22	-1.38
2	-0.07	-0.52	-0.20	-1.57	-0.34	-2.69
3	-0.11	-0.96	-0.02	-0.16	-0.25	-2.10
4	0.14	1.27	0.05	0.44	-0.08	-0.72
5	-0.06	-0.52	0.01	0.07	0.13	1.21
6	-0.07	-0.68	0.00	0.03	-0.16	-1.57
7	0.16	1.71	0.07	0.70	-0.01	-0.12
8	0.05	0.43	-0.02	-0.14	0.20	2.40
9	0.15	1.69	0.10	1.19	-0.04	-0.41
High	0.02	0.20	-0.03	-0.32	0.05	0.57
High - Low	0.42	1.79	0.14	0.68	0.27	1.38

Table 5: **Returns and Alphas of Portfolios Sorted on Raw Number of IPs**

This table reports the monthly excess returns and Carhart (1997) four factor alphas for decile portfolios sorted on raw number of IPs searching for Edgar files. At the end of each month, all stocks are sorted into deciles based on their raw number of IPs and a long-short portfolio is formed by buying the highest decile and shorting the lowest decile portfolio. Portfolio returns are computed over the next month. Panel A reports results for equally weighted excess return and Panel B shows results Carhart (1997) four factor alphas. The sample runs from January 2003 to December 2014.

Panel A: Equal-weighted Decile Portfolio Excess Return						
	IP_total	t-stat	IP_funtl	t-stat	IP_10K	t-stat
Low	0.73	2.17	0.87	2.62	0.73	2.04
2	0.92	2.17	0.80	1.90	0.80	1.87
3	1.01	2.19	0.91	1.89	0.63	1.32
4	1.12	2.22	1.12	2.28	0.95	1.86
5	1.12	2.19	0.89	1.73	1.05	2.01
6	1.07	2.08	1.17	2.23	1.12	2.10
7	1.01	1.92	1.12	2.11	1.12	2.07
8	1.14	2.06	1.05	1.91	1.22	2.25
9	0.99	1.84	1.04	1.96	1.19	2.26
High	0.98	1.99	1.09	2.20	1.10	2.31
High - Low	0.26	1.19	0.22	0.68	0.37	1.58

Panel B: Equal-weighted Decile Portfolio 4-factor alpha						
	IP_total	t-stat	IP_funtl	t-stat	IP_10K	t-stat
Low	0.05	0.30	0.18	1.18	0.04	0.23
2	0.06	0.44	-0.05	-0.39	-0.12	-0.78
3	0.08	0.68	-0.07	-0.54	-0.26	-1.96
4	0.00	0.00	0.05	0.35	-0.08	-0.59
5	0.01	0.08	-0.11	-0.94	-0.08	-0.70
6	-0.09	-0.83	-0.02	-0.17	-0.01	-0.12
7	-0.09	-0.94	-0.08	-0.96	0.01	0.15
8	-0.11	-1.17	-0.08	-0.73	0.06	0.77
9	-0.13	-1.33	-0.05	-0.49	0.05	0.49
High	-0.05	-0.50	-0.02	-0.20	0.13	1.49
High - Low	-0.09	-0.56	-0.20	-1.15	0.09	0.47

Table 6: **Two-way sorts on Firm Size and Abnormal Number of IPs**

This table reports monthly Carhart (1997) 4-factor alphas (in percentages) sorted on stock's market capitalization and abnormal number of IPs searching 10-K files (AIP_10K). AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in EDGAR database on a set of firm characteristics. At the end of each month, all the stocks are sorted into quintiles based on NYSE size breakpoints, and within each quintile the stocks are further sorted into quintiles based on their AIP_10K. We also report, for each size quintile, the high-AIP minus high-AIP portfolio alpha. Panel A reports results on an equal-weighted basis and panel B on a value-weighted basis. The sample runs from January 2003 to December 2014.

Panel A: Equal-weighted 4 factor alpha					
	Small firms	2	3	4	Large firms
Low AIP	-0.48	-0.22	-0.18	-0.06	-0.22
2	-0.29	-0.14	-0.19	0.02	0.06
3	-0.11	-0.10	-0.10	-0.01	0.02
4	0.21	0.07	0.16	0.15	0.28
High AIP	0.54	-0.12	0.19	0.12	0.23
High-Low	1.03	0.10	0.37	0.19	0.45
t-stat	5.19	0.40	1.95	0.87	2.29
Panel B: Value-weighted 4 factor alpha					
	Small firms	2	3	4	Large firms
Low AIP	-0.57	-0.22	-0.19	-0.08	-0.24
2	-0.35	-0.09	-0.17	0.00	-0.06
3	-0.16	-0.11	-0.05	-0.01	-0.02
4	-0.02	0.08	0.17	0.13	0.05
High AIP	0.36	-0.13	0.22	0.14	0.07
High-Low	0.93	0.09	0.41	0.22	0.31
t-stat	4.66	0.36	2.18	1.01	1.89

Table 7: **Limits to Arbitrage**

This table reports results on limits to arbitrage. We sort stocks into tercile based on each limits-to-arbitrage variable X, including idiosyncratic volatility (IVOL), institutional ownership (IO) and analyst coverage (Coverage). We then independently sort stocks into quintiles based on abnormal number of IPs searching 10-K files (AIP_10K). AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in Edgar database on a set of firm characteristics. We report the Carhart (1997) four-factor alpha of the lowest and highest AIP portfolio among the lowest and highest X group. The "High-Low" column reports Carhart (1997) four-factor alpha of the high-AIP minus low-AIP portfolios. The sample runs from January 2003 to December 2014.

	Low AIP_10K	High AIP_10K	High-Low
High IVOL	-0.76 (-3.27)	0.48 (1.95)	1.24 (4.44)
Low IVOL	0.03 (0.30)	0.27 (3.34)	0.23 (1.76)
High IO	-0.17 (-1.61)	0.23 (1.75)	0.40 (2.36)
Low IO	-0.56 (-3.53)	0.48 (1.91)	1.03 (4.41)
High Coverage	-0.33 (-3.08)	0.18 (1.54)	0.51 (3.07)
Low Coverage	-0.41 (-2.59)	0.68 (3.23)	1.10 (5.77)

Table 8: **Complexity of Financial Filings**

This table reports return predictability results on variation in the complexity of financial filings. We sort stocks into tercile based on the size or number of words of its most recent 10-K filing. We then independently sort stocks into quintiles based on abnormal number of IPs searching 10-K files (AIP_10K). AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in Edgar database on a set of firm characteristics. We report the Carhart (1997) four-factor alpha of the lowest and highest AIP portfolio among the lowest and highest information cost group. The "High-Low" column reports Carhart (1997) four-factor alpha of the high-AIP minus low-AIP portfolios. The sample runs from January 2003 to December 2014.

	Low AIP_10K	High AIP_10K	High-Low
Large File Size	-0.48 (-3.98)	0.44 (2.86)	0.92 (4.46)
Small File Size	-0.29 (-2.13)	0.36 (2.64)	0.65 (3.51)
More word count	-0.48 (-4.08)	0.49 (3.18)	0.97 (5.06)
Lesser word count	-0.36 (-3.05)	0.32 (2.48)	0.68 (4.21)

Table 9: **Fama-MacBeth Regression: Baseline**

This table reports the results from the Fama and MacBeth (1973) regression of monthly stock returns on abnormal number of IPs searching Edgar files (AIP). AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of files in EDGAR database on a set of firm characteristics. Column (1)-(3) use IPs searchings all types of EDGAR filings. Column (4)-(6) uses IPs searching 10-K, 10-Q and 8-K files. Column (7)-(9) looks at IPs searching for 10-K files. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	All EDGAR Filings			10-K, 10-Q and 8K			10-K		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
AIP	0.0060*** (2.68)	0.0053*** (2.64)	0.0050*** (2.88)	0.0047*** (2.70)	0.0041*** (2.78)	0.0042*** (2.94)	0.0051*** (3.73)	0.0046*** (3.81)	0.0044*** (3.74)
Rev		-0.0247*** (-3.18)	-0.0283*** (-3.74)		-0.0245*** (-3.16)	-0.0281*** (-3.72)		-0.0247*** (-3.19)	-0.0284*** (-3.75)
LnME		-0.0006 (-0.89)	-0.0014** (-2.59)		-0.0006 (-0.92)	-0.0014** (-2.60)		-0.0006 (-0.93)	-0.0014** (-2.58)
LnBM		0.0019 (1.64)	0.0014 (1.29)		0.0019 (1.59)	0.0013 (1.24)		0.0019 (1.58)	0.0013 (1.24)
Mom		-0.0058 (-0.95)	-0.0048 (-0.88)		-0.0057 (-0.94)	-0.0047 (-0.86)		-0.0058 (-0.94)	-0.0048 (-0.86)
Ivol			-0.0015 (-0.02)			-0.0025 (-0.04)			-0.0007 (-0.01)
Turnover12			-0.0094 (-1.37)			-0.0091 (-1.32)			-0.0089 (-1.28)
IO			0.0122*** (4.00)			0.0119*** (3.94)			0.0114*** (3.86)
Constant	0.0123** (2.18)	0.0122 (1.65)	0.0119** (2.33)	0.0122** (2.18)	0.0122* (1.66)	0.0120** (2.36)	0.0122** (2.18)	0.0123* (1.67)	0.0119** (2.35)
Ave.R-sq	0.003	0.030	0.046	0.003	0.030	0.046	0.003	0.030	0.046
N.of Obs.	483667	483667	480793	483667	483667	480793	483667	483667	480793

Table 10: **Predicting Long-horizon Returns**

This table reports the results from the Fama and MacBeth (1973) regression of cumulative returns from month $t + j$ to $t + k$ on abnormal number of IPs searching 10-K files in EDGAR database (AIP_10K) at month t . AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in EDGAR database on a set of firm characteristics. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month $t-12$ to $t-2$. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	Ret(2,4)	Ret(5,7)	Ret(8,13)	Ret(14,25)
AIP_10K	0.0102*** (2.95)	0.0068** (2.05)	0.0150 (1.57)	0.0175 (0.64)
Rev	-0.0072 (-0.53)	0.0037 (0.21)	0.0033 (0.11)	-0.0451 (-0.93)
LnME	-0.0023 (-1.64)	-0.0013 (-1.03)	-0.0015 (-0.61)	-0.0048 (-1.11)
LnBM	0.0046* (1.72)	0.0041 (1.57)	0.0118** (2.36)	0.0197* (1.79)
Mom	-0.0193 (-1.24)	-0.0117 (-0.88)	-0.0300* (-1.75)	-0.0421 (-1.26)
Ivol	0.0407 (0.20)	-0.0184 (-0.10)	0.2652 (0.73)	0.5759 (0.84)
Turnover12	-0.0165 (-0.92)	-0.0312* (-1.95)	-0.0451 (-1.53)	-0.0488 (-1.08)
IO	0.0116 (1.63)	0.0152** (2.18)	0.0414** (2.42)	0.0956** (2.47)
Constant	0.0370** (2.41)	0.0281* (1.72)	0.0451 (1.53)	0.0947 (1.51)
Ave.R-sq	0.051	0.044	0.036	0.035
N.of Obs.	469185	456068	425505	360584

Table 11: **Which Types of EDGAR Files?**

This table reports the results from the Fama and MacBeth (1973) regression of monthly stock returns on abnormal number of IPs searching EDGAR filings (AIP). AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of files in Edgar database on a set of firm characteristics. AIP_total use IPs searchings all types of Edgar files. AIP_fundl uses IPs searching 10-K, 10-Q and 8-K files. AIP_10K looks at IPs searching for 10-K files. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	(1)	(2)
AIP_total	-0.0014 (-0.63)	-0.0003 (-0.17)
AIP_fundl	0.0022 (1.11)	0.0012 (0.70)
AIP_10K	0.0049*** (3.96)	0.0043*** (4.02)
Rev		-0.0287*** (-3.80)
LnME		-0.0014** (-2.52)
LnBM		0.0013 (1.24)
Mom		-0.0048 (-0.88)
Ivol		-0.0027 (-0.04)
Turnover12		-0.0088 (-1.27)
IO		0.0112*** (3.84)
Constant	0.0122** (2.18)	0.0120** (2.34)
Ave.R-sq	0.005	0.048
N.of Obs.	483667	480793

Table 12: Number of IPs or Number of Searches?

This table reports the results from the Fama and MacBeth (1973) regression. Asearch is the residual from a monthly regression of log one plus total number of EDGAR requests for SEC filings. AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for EDGAR files on a set of firm characteristics. Column (1)-(2) looks at searching activities for all types of EDGAR files. Column (3)-(4) look at searching activities for 10-K, 10-Q and 8-K files. Column (5)-(6) looks at searching activities for 10-K files. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	All EDGAR Files		10K, 10Q, 8K		10K	
Asearch	0.0014 (1.54)	-0.0004 (-0.42)	0.0020* (1.90)	-0.0024 (-1.49)	0.0033*** (3.93)	-0.0039 (-1.57)
AIP		0.0055** (2.45)		0.0062*** (2.83)		0.0084*** (2.90)
Rev	-0.0283*** (-3.73)	-0.0284*** (-3.76)	-0.0283*** (-3.74)	-0.0284*** (-3.77)	-0.0284*** (-3.75)	-0.0289*** (-3.75)
LnME	-0.0014** (-2.59)	-0.0014*** (-2.63)	-0.0014** (-2.61)	-0.0014** (-2.52)	-0.0014*** (-2.64)	-0.0013*** (-3.11)
LnBM	0.0013 (1.26)	0.0014 (1.31)	0.0014 (1.34)	0.0014 (1.36)	0.0012 (1.13)	0.0015* (1.71)
Mom	-0.0049 (-0.89)	-0.0048 (-0.88)	-0.0048 (-0.87)	-0.0049 (-0.89)	-0.0048 (-0.86)	-0.0049 (-1.15)
Ivol	0.0048 (0.07)	-0.0014 (-0.02)	0.0065 (0.09)	-0.0033 (-0.05)	0.0039 (0.05)	-0.0021 (-0.03)
Turnover12	-0.0100 (-1.46)	-0.0096 (-1.39)	-0.0095 (-1.38)	-0.0091 (-1.33)	-0.0095 (-1.37)	-0.0088 (-1.33)
IO	0.0127*** (4.10)	0.0123*** (4.04)	0.0122*** (4.06)	0.0115*** (3.86)	0.0120*** (4.03)	0.0109*** (3.57)
Constant	0.0115** (2.26)	0.0120** (2.35)	0.0116** (2.29)	0.0119** (2.33)	0.0117** (2.32)	0.0120*** (3.19)
Ave.R-sq	0.046	0.047	0.046	0.048	0.046	0.049
N.of Obs.	480793	480793	480793	480793	480793	480793

Table 13: **Predicting Fundamental Performance**

This table reports the results from panel regression of future fundamental performance measure on abnormal number of IPs searching 10K files in Edgar database at month t . AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for EDGAR filings on a set of firm characteristics. The dependent variable in Column (1)-(3) are change of quarterly Return-on-Assets from four quarters ago. The dependent variable in Column (4)-(6) are monthly revision of analysts consensus annual EPS forecast. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month $t-12$ to $t-2$. Coverage is log one plus analyst coverage. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). We control year-quarter fixed effect in Column (1)-(3) and year-month fixed effect in Column (4)-(6). Turnover12 is the monthly turnover ratio averaged over past 12 months. Standard errors are double clustered at both firm and time level. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	Change of ROA			Forecast Revision		
	AIP_total	AIP_fundl	AIP_10K	AIP_total	AIP_fundl	AIP_10K
AIP	0.0017*	0.0026**	0.0028***	0.0007***	0.0016***	0.0019***
	(1.96)	(2.51)	(2.92)	(2.78)	(6.19)	(5.28)
LROA	-0.3425***	-0.3428***	-0.3430***			
	(-4.71)	(-4.73)	(-4.74)			
LnME	0.0008	0.0008	0.0008	-0.0005	-0.0005	-0.0005
	(1.27)	(1.31)	(1.33)	(-1.51)	(-1.55)	(-1.63)
LnBM	-0.0013	-0.0012	-0.0012	-0.0008**	-0.0008**	-0.0009**
	(-0.87)	(-0.84)	(-0.83)	(-2.37)	(-2.39)	(-2.47)
Mom	0.0100***	0.0099***	0.0100***	0.0025***	0.0025***	0.0025***
	(3.55)	(3.56)	(3.57)	(5.20)	(5.14)	(5.18)
Coverage	0.0004	0.0004	0.0005	0.0021***	0.0021***	0.0021***
	(0.29)	(0.31)	(0.32)	(3.30)	(3.29)	(3.28)
Turnover12	-0.0118**	-0.0117**	-0.0117**	-0.0082***	-0.0082***	-0.0082***
	(-2.43)	(-2.42)	(-2.43)	(-3.15)	(-3.16)	(-3.17)
IO	-0.0010	-0.0011	-0.0013	0.0049***	0.0051***	0.0052***
	(-0.48)	(-0.51)	(-0.61)	(5.47)	(5.61)	(5.73)
Ivol	-0.0777	-0.0773	-0.0775	-0.1111**	-0.1115**	-0.1138**
	(-1.41)	(-1.41)	(-1.42)	(-2.35)	(-2.37)	(-2.41)
Time FE	yes	yes	yes	yes	yes	yes
Adj.R-sq	0.056	0.056	0.056	0.002	0.002	0.002
N.of Obs.	128504	128504	128504	348130	348130	348130

Table 14: **Mutual Fund Outflows Induced Mispricing and Abnormal Number of IPs**

This table reports the results from the Fama and MacBeth (1973) regression of quarterly change of abnormal number of IPs searching EDGAR files on quarterly mutual fund outflows. Outflows is calculated following Edmans, Goldstein, and Jiang (2012). AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for EDGAR filings on a set of firm characteristics. dAIP equals the within-firm change in AIP in the quarter around mutual fund outflows. LnME is the natural log of firm's market capitalization at the end of the June of each year in millions of US dollars. Coverage is log one plus analyst coverage. Turnover12 is the monthly turnover ratio averaged over the past 12 months. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	dAIP_total		dAIP_funtl		dAIP_10K	
	(1)	(2)	(3)	(4)	(5)	(6)
Outflows	-2.4242***	-1.7256***	-1.9145***	-1.3527***	-1.9303**	-1.5459**
	(-4.02)	(-4.92)	(-3.36)	(-3.27)	(-2.06)	(-2.31)
LnME		-0.0091***		-0.0094***		-0.0093***
		(-6.03)		(-5.68)		(-5.81)
LnBM		0.0013		-0.0014		-0.0017
		(0.56)		(-0.57)		(-0.75)
Coverage		0.0080***		0.0076***		0.0087***
		(4.50)		(4.28)		(3.70)
Ivol		-1.8233***		-1.9963***		-1.8354***
		(-6.48)		(-7.68)		(-6.19)
Turnover12		-0.0015		0.0158		0.0203
		(-0.09)		(1.13)		(1.56)
IO		-0.0023		-0.0141**		-0.0143**
		(-0.36)		(-2.54)		(-2.28)
Mom		-0.0336***		-0.0370***		-0.0398***
		(-5.17)		(-5.70)		(-7.68)
Constant	0.0007	0.0901***	0.0050**	0.1036***	0.0049**	0.0967***
	(0.29)	(7.79)	(2.09)	(8.54)	(2.06)	(6.54)
Ave.R-sq	0.001	0.031	0.001	0.034	0.001	0.026
N.of Obs.	131863	131041	131863	131041	131863	131041

Table 15: **Information Acquisition and Institutional Trading**

This table reports Fama and MacBeth (1973) regression of institutional trading at quarter q on abnormal number of IPs searching EDGAR filings (AIP) at quarter $q - 1$. The dependent variable in Column (1)-(3) is net buying by hedge fund in that quarter. The dependent variable in Column (4)-(6) is net buying by mutual funds in that quarter. AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of files in EDGAR database on a set of firm characteristics. Column (1) and (4) use IPs searchings all types of EDGAR filings. Column (2) and (5) uses IPs searching 10-K, 10-Q and 8-K files. Column (3) and (6) looks at IPs searching for 10-K files. Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month $t-12$ to $t-2$. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	Net Buying by Hedge Funds			Net Buying by Mutual Funds		
	AIP_total	AIP_funtl	AIP_10K	AIP_total	AIP_funtl	AIP_10K
AIP	0.0053** (2.45)	0.0043** (2.43)	0.0034** (2.27)	0.0030 (0.52)	0.0046 (0.77)	0.0065 (1.00)
Lagged Holding	-0.1245*** (-7.62)	-0.1234*** (-7.72)	-0.1218*** (-9.24)	-0.1557*** (-3.91)	-0.1569*** (-3.83)	-0.1573*** (-3.77)
LnME	-0.0011*** (-3.37)	-0.0011*** (-3.39)	-0.0011*** (-4.12)	-0.0003 (-1.07)	-0.0003 (-1.08)	-0.0003 (-1.06)
LnBM	-0.0002* (-1.88)	-0.0002* (-1.74)	-0.0002 (-1.64)	0.0003 (0.32)	0.0004 (0.38)	0.0003 (0.34)
Cov	-0.0006 (-1.40)	-0.0006 (-1.33)	-0.0005 (-1.51)	-0.0008 (-0.30)	-0.0007 (-0.26)	-0.0005 (-0.22)
Ivol	0.0145 (0.76)	0.0153 (0.80)	0.0164 (0.97)	-0.1965* (-1.96)	-0.1975* (-1.94)	-0.2012* (-1.95)
Turnover12	0.0041* (1.76)	0.0042* (1.78)	0.0044** (2.23)	-0.0069** (-2.04)	-0.0068** (-2.22)	-0.0064** (-2.38)
IO	0.0150*** (3.09)	0.0145*** (3.14)	0.0140*** (3.87)	0.0739*** (2.86)	0.0736*** (2.90)	0.0728*** (2.92)
Mom	0.0003 (1.14)	0.0004 (1.22)	0.0002 (0.97)	0.0066*** (6.44)	0.0067*** (6.29)	0.0065*** (7.74)
Constant	0.0070*** (4.06)	0.0070*** (4.01)	0.0069*** (4.68)	0.0048* (1.95)	0.0050* (1.91)	0.0053* (1.96)
Ave.R-sq	0.091	0.091	0.091	0.113	0.113	0.113
N.of Obs.	131795	131795	131795	131795	131795	131795

Appendices

Table A1: **Robustness of Decile Portfolio Sorts**

This table reports several robustness tests for a long/short portfolio based on abnormal number of IPs searching for 10-K files in EDGAR database (AIP_10K). AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in EDGAR database on a set of firm characteristics. In the first set of robustness tests, we report the gross return-weighted portfolio returns in which the weights are 1 + the stock's lagged monthly return, following Asparouhova, Bessembinder, and Kalcheva (2013). The second robustness test show the portfolio returns adjusted using DGTW method. The third set of robustness tests shows the Fama-French 48 industry-adjusted excess return. The fourth row shows the alpha using Pástor and Stambaugh (2003) liquidity factor augmented with the Fama-French factors and the momentum factor. In the fifth set of tests, we report the alphas using the Fama and French (2016) Five Factor model. In the sixth and seventh set of tests, we report the alphas using the Stambaugh and Yuan (2016) Mispricing Factors model and the Hou, Xue, and Zhang (2015) Q-factor model. In the eighth set of analyses, we exclude stocks whose market capitalization are in the bottom quintile based on NYSE size breakpoints. In the ninth panel, we skip six months between the moment at abnormal IP is constructed and the moment at which we start measuring returns. In the tenth and eleventh row, we report the 4-factor alpha for two sub-sample periods, one from 2003 to 2008 and another from 2009 to 2014.

	EW	VW
Gross return-weighted portfolio	1.096 (5.16)	NA
DGTW adjusted	0.910 (4.51)	0.410 (2.22)
FF48 Industry-adjusted	0.739 (3.26)	0.155 (1.16)
FF + Cahart + PS Factor	0.800 (4.23)	0.348 (1.78)
FF five factor (2015)	0.685 (3.36)	0.248 (1.19)
Mispricing factors (Stambaugh and Yuan 2017)	0.892 (4.42)	0.276 (1.35)
Q-factor (Hou, Xue and Zhang 2015)	0.897 (4.66)	0.183 (0.87)
Remove microcap stocks	0.518 (2.58)	0.276 (1.35)
Skip six months	0.532 (2.23)	0.266 (1.28)
2003-2008	0.620 (2.41)	0.261 (0.89)
2009-2014	1.073 (3.74)	0.121 (0.45)

Table A2: **Alternative Implementations of AIP**

This table reports several alternative implementations of AIP_10K when calculating long/short portfolio Carhart (1997) 4-factor alpha. AIP_10K is the residual from a monthly regression of log one plus total number of unique IP addresses searching for 10-K files in EDGAR database on a set of firm characteristics. In the first row, we calculate AIP_10K using model (7) of equation 2. In the second row, we sort portfolios based on changes in AIP_10K relative to its 12-month moving average. In the third row, we also include the square term of the four firm characteristics when calculating AIP. In the fourth row, we include lagged number of IPs in expected IP regression. Column (1) reports the results for equal-weighted portfolio, and Column (2) reports for the value-weighted portfolio. The sample runs from January 2003 to December 2014.

	EW	VW
Model (7) of Expected IP Regression	0.658 (3.95)	0.156 (0.82)
Change in AIP relative to 12 months average	0.883 (4.82)	0.388 (1.44)
Nonlinear functional form of Expected IP Regression	0.689 (4.30)	0.552 (2.39)
Control for lagged # of IPs in Expected IP Regression	0.698 (5.44)	0.508 (2.03)

Table A3: Fama-MacBeth Regression: Controlling for Earnings Surprise, Earnings Announcement Premium and Change of Breadth of Ownership

This table reports the results from the Fama and MacBeth (1973) regression of monthly stock returns on abnormal number of IPs searching EDGAR filings (AIP). AIP is the residual from a monthly regression of log one plus total number of unique IP addresses searching for all type of files in EDGAR site on a set of firm characteristics. Column (1) and (4) use IPs searchings all types of EDGAR filings. Column (2) and (5) uses IPs searching 10-K, 10-Q and 8-K files. Column (3) and (6) looks at IPs searching for 10-K files. SUE is a firm's standardized unexplained earnings, defined as the realized earnings per share (EPS) minus EPS from four quarters prior, divided by the standard deviation of this difference over the prior eight quarters. EAM is a dummy variable that equals one when a given firm announces earnings in the month. dBreadth is the percentage change of breadth of 13F institutional ownership, following Chen, Hong, and Stein (2002). Size (LnME) is the natural log of a firm's market capitalization at the end of the June of each year. Book-to-market (LnBM) is the natural log of the book-to-market ratio. The cases with negative book value are deleted. Momentum (MOM) is defined as the cumulative returns from month t-12 to t-2. The short term reversal measure (REV) is the lagged monthly return. Institutional ownership (IO) is the sum of shares held by institutions from 13F filings in each quarter divided by the total shares outstanding. IVOL is the idiosyncratic volatility, calculated following Ang, Hodrick, Xing, and Zhang (2006). Turnover12 is the monthly turnover ratio averaged over past 12 months. All t-statistics are Newey-West adjusted to control for heteroskedasticity and autocorrelation. ***, **, and * represent significance levels of 1%, 5%, and 10%, respectively.

	(1)	(2)	(3)	(4)	(5)	(6)
	AIP_total	AIP_fundl	AIP_10K	AIP_total	AIP_fundl	AIP_10K
AIP	0.0057*** (3.31)	0.0047*** (3.32)	0.0044*** (4.02)	0.0056*** (3.29)	0.0046*** (3.27)	0.0044*** (3.97)
Rev	-0.0300*** (-4.08)	-0.0298*** (-4.06)	-0.0300*** (-4.09)	-0.0299*** (-4.10)	-0.0297*** (-4.08)	-0.0299*** (-4.10)
LnME	-0.0017*** (-3.13)	-0.0017*** (-3.16)	-0.0016*** (-3.14)	-0.0017*** (-3.13)	-0.0017*** (-3.16)	-0.0016*** (-3.15)
LnBM	0.0016 (1.56)	0.0016 (1.51)	0.0016 (1.53)	0.0016 (1.57)	0.0016 (1.51)	0.0016 (1.54)
Mom	-0.0065 (-1.13)	-0.0064 (-1.12)	-0.0064 (-1.11)	-0.0068 (-1.18)	-0.0067 (-1.17)	-0.0067 (-1.17)
Ivol	0.0116 (0.16)	0.0104 (0.15)	0.0121 (0.17)	0.0077 (0.11)	0.0063 (0.09)	0.0083 (0.12)
Turnover12	-0.0097 (-1.38)	-0.0093 (-1.32)	-0.0092 (-1.29)	-0.0093 (-1.33)	-0.0089 (-1.27)	-0.0088 (-1.24)
IO	0.0118*** (3.42)	0.0114*** (3.35)	0.0109*** (3.26)	0.0118*** (3.43)	0.0114*** (3.36)	0.0109*** (3.27)
SUE	0.0028*** (8.36)	0.0028*** (8.39)	0.0027*** (8.46)	0.0027*** (8.29)	0.0027*** (8.31)	0.0027*** (8.37)
EAM	0.0032** (2.51)	0.0034** (2.59)	0.0027** (2.23)	0.0032** (2.48)	0.0033** (2.56)	0.0027** (2.19)
dBreadth				0.0710 (0.95)	0.0789 (1.05)	0.0842 (1.13)
Constant	0.0121** (2.46)	0.0122** (2.49)	0.0122** (2.50)	0.0121** (2.42)	0.0123** (2.46)	0.0123** (2.47)
Ave.R-sq	0.051	0.051	0.051	0.052	0.052	0.052
N.of Obs.	443261	443261	443261	442794	442794	442794