

# Categorization Bias in the Stock Market\*

Philipp Krüger, Augustin Landier, and David Thesmar

First Version: January 2012

This Version: April 2012

## Abstract

This paper provides evidence of categorization bias in financial markets. Some investors perceive individual firms through the lenses of industries. Such categorical thinking generates mispricing and predictability in stock-returns. We measure the difference in returns between a firm's official SIC industry and its underlying fundamentals by constructing a basket of closely related firms, based on Hoberg and Phillips (2010a,b). We find that stocks comove strongly with their official industry in the short-term and then revert toward their basket of Hoberg and Phillips comparables. Long-short strategies exploiting mispricing due to industry categorization generate statistically significant and economically sizable risk-adjusted excess returns. We also provide evidence that financial analysts are biased by industry categorization: When a firm's official industry is less representative of its fundamentals, analysts sometimes tend to put excessive weight on information related to the official industry and thus make predictable forecast errors.

JEL-Classification: G10, G11, G14, M41

Keywords: Limited attention, behavioral finance, return predictability, analyst forecasts accuracy

---

\*Krueger is at the Geneva Finance Research Institute, Université de Genève. Landier is at the Toulouse School of Economics. Thesmar is at HEC Paris and CEPR. Philipp Krüger, philipp.krueger@unige.ch, Telephone: +41 (0)22 379 85 69; Augustin Landier, augustin.landier@tse-fr.eu, Telephone: +33 (0)5 61 12 86 88; David Thesmar, thesmar@hec.fr, Telephone: +33 (0)1 39 67 94 12. Landier acknowledges support from the Scor Chair at Fondation Jean-Jacques Laffont. Thesmar thanks the HEC Foundation for financial support.

# I Introduction

Even if the underlying reality is continuous, economic agents tend to see the world through coarse categories: a student is an A-student or is not, economies are in a recession or aren't, bonds are "investment grade" or "junk", firms belong to manufacturing or services... Categories exist because memory is limited. It would be inefficiently costly for agents to memorize all characteristics of individuals, events or firms: They use categories as "boxes" in which similar objects can be stored and then tend to consider them as similar when making decisions. The use of representative categories biases expectations as soon as agents forget whether an object is strongly or weakly representative of its category: An agent thinking through categories is biased in a predictable manner vis-a-vis a fully rational agent, making use of a wider information set. For instance, imagine that you grade students as A, B or C based on an underlying continuous score and memorize only the final grade. A student who got an A but was very close to the B threshold is categorized together with students that are on average better than him. Someone who remembers the actual score can thus predict that you tend to overestimate the level of this particular student. There is experimental evidence that such categorization bias is prevalent. For instance, temperature predicted by non-professional subjects jumps at the beginning of each month, even though forecasts are, on average over the month, unbiased. This is because each forecaster overweights the importance of the month in her intuitive statistical model of temperature ("July is hotter than June"). People therefore behave as if categories were better descriptions of reality than they actually are. Following Mullainathan (2002), we call this bias "categorization bias".

Categorization bias implies a combination of under- and overreaction to information (Mullainathan (2002); Hong, Stein, and Yu (2007)). If new information does not trigger a change in category, learning is abnormally slow vis-a-vis the rational Bayesian benchmark. But, as soon as new, though marginal, information forces a change in categories, agents' expectations over-shoot. For instance, as long as it keeps its Michelin star, a restaurant is

perceived as very good, even if the quality of its food has started declining; then, suddenly, it loses its star and falls strongly into disgrace with the gourmet community. From the point of view of an informed observer of the slow decline in quality of the restaurant, this discontinuous reaction is an error: there is under-reaction, followed by over-reaction of the public. Categorization bias seems to be everywhere in economics, with important consequences: In product markets (where it may explain advertisement spending, as in Mullainathan, Shleifer, and Schwartzstein (2008)), in the labor market (where it could lead to inefficient discrimination), in finance (where it could lead to return predictability, as in Mullainathan (2002)), etc. Yet, and surprisingly, the economics literature offers little evidence of this bias based on field data.

This paper shows that categorization bias has sizable effects on asset prices. Financial markets offer a natural testing ground of categorical thinking because of their wealth of data on expectations (analyst forecasts, stock price movements). Our premise is that some market participants tend to mentally group firms according to an industry classification (e.g. the SIC code reported in Compustat). However, such classifications provide only a noisy and sometimes misleading description of what firms actually do (Hoberg and Phillips (2010a)). Investors relying strongly on industrial categorization neglect some public information: For instance they could use more accurate comparables than simply looking at firms in the same SIC category. By using the heuristics that a firm is "like its SIC industry", such investors thus make predictable valuation mistakes, which in turn create mispricing and predictability in stock-returns. Think of a news shock affecting the "paper and allied products" industry. A fully rational market participant is aware that some members of this industry, like Rock-Tenn - which does paperboard food packaging - are well represented by the "paper" category, while others, like Schweitzer-Mauduit - which produces cigarette paper and therefore depends on cigarette sales - are not. If investors were fully rational, Schweitzer-Mauduit would react little to "paper" shocks, while Rock-Tenn would react more. If, however, investors are subject to categorization bias, Schweitzer-Mauduit is perceived to be more similar to mainstream paper firms

(such as Rock-Tenn) than it actually is. This leads to a testable implication: Firms (like Schweitzer-Mauduit as opposed to Rock-Tenn) should tend to over-react to their official industry shocks, and when such overreaction happens, they should subsequently revert toward their fundamental value.

To operationalize our empirical strategy, we first develop a method to detect whether a stock is subject to categorization bias at a given point in time. To do this, we rely on a measure of firms' industry fundamentals developed recently by Hoberg and Phillips (2010a,b), who compare firms based on the "business description" section of their 10K filing. We use their data to identify a firm's "fundamental" peers, i.e. firms whose business description is highly similar and thus likely to be a good proxy of fundamentals. For each firm, we then calculate the returns of its "official" peers, as defined by the standard SIC classification from Compustat, and those of its "fundamental" peers, as defined by the Hoberg and Phillips (2010a,b) similarity measure. We show that, at high frequencies (e.g. weekly), firms comove strongly with their official industry and weakly with their fundamental peers while the opposite pattern holds at lower frequencies. This sharp reversal is in line with the hypothesis that bounded rationality leads some agents to overweight official industry shocks in the short-term and that such mistakes are corrected over time. We test this hypothesis further by constructing portfolios based on the divergence between official and fundamental industry returns: If investors tend to mix up a firm with its official industry, a highly positive difference between official and fundamental returns should be a signal that predicts lower returns for the firm in the future. Going back to our prior example, if last month, Schweitzer-Mauduit's official peers (like Rock-Tenn) have outperformed their fundamental peers (like Alliance One, which grows tobacco), SM is likely to be overvalued, and its stock price will decrease subsequently. In a second step, we refine this signal by focusing on firms whose stock price has closely followed their "official industry" in a given month: This is when categorization bias is most likely to be occurring with high intensity. To return to the Schweitzer-Mauduit example, the overvaluation should be more severe if the stock has closely followed

the "paper" industry in the past month. We label such firms "social industry followers".

Consistent with the fact that investors put too much weight on the social industry classification, we find that "social industry followers", whose social peers have outperformed (resp. under-performed) fundamental peers, earn large negative (resp. positive) returns over the next month. A long-short strategy based on this insight yields a monthly abnormal return of 1.57 % (more than 18% annual) with a t-stat of 4.1. We also show that return predictability is more pronounced among smaller firms. This is consistent with bounded rationality: If there is a fixed or convex cost associated to evaluating an individual firm, more investors are likely to pay that cost for large firms (for which the stakes are higher) than for small firms. Small firms are thus more prone to categorization bias.

To complement the returns evidence, we then look at actual expectations data: We find that forecast errors of stock analysts follow a pattern consistent with categorization bias. To show this, we adapt the above methodology to earnings forecast data. We compute the difference between earnings forecasts of "social" peers and the forecasts of "fundamental" peers. This forecast discrepancy measures the extent to which a forecast based on the social classification is upward biased. For firms for which analysts tend to rely heavily on the consensus in the social industry, we find that analyst expectations tend to be upward biased when the forecast discrepancy is positive. This result is consistent with the hypothesis that some analysts are relying too much on comparisons with firms of the social industry, which makes their expectations biased when a firm is not well classified.

Our paper contributes to several literatures. The behavioral economics literature has produced models of non-Bayesian inference based on categorical thinking. These models rest on the fact that agents assign situations to categories, and assign a probability one to the most likely category. This non-Bayesian feature is what creates the categorization bias, and explains over- and underreaction to information (see Mullainathan (2002); Hong, Stein, and Yu (2007); Mullainathan, Shleifer, and Schwartzstein (2008) embed

categorization bias in a communication model). In this sense, our paper is related to well documented cognitive biases in psychology such as the salience or representativeness heuristics. This literature is, however, largely theoretical (its aim is to build theoretical frameworks able to account for available experimental evidence); our contribution is to provide direct evidence of categorization bias from field data. The paper also contributes to the large literature on stock habitats (see Barberis, Shleifer, and Wurgler (2005)). This literature suggests that stocks may comove "too much" when they are traded by similar investors (for instance by index funds when they belong to an index). Our paper suggests that categorization bias is a good candidate to generate excess comovement. In a sense, the official classification tends to generate stock habitat due to investor categorization biases. Last, this paper is related to the literature on investor attention (see Hong, Lim, and Stein (2000)). This literature shows that mispricing may occur as a by-product of investors' inability to quickly process all relevant information (see also Cohen and Frazzini (2008)). In another recent paper, Cohen and Lou (2012) show that investors have trouble assessing the effect of industry shocks on diversified firms, which generates momentum in their returns as information is slowly impounded into prices. Our paper contributes to this literature by isolating the impact of a well defined psychological bias.

Section II describes the data. Section III describes the results using returns data. Section IV describes our results on analyst forecasts.

## II The Data

### II.1 Text-based Network Industry Classifications

The main data used in this study are the Text-based Network Industry Classifications (TNIC)<sup>1</sup> developed in Hoberg and Phillips (2010a,b). The Hoberg and Phillips (HP) industry classifications are constructed by parsing product descriptions from 10K forms and forming word vectors for each firm to compute continuous measures of product simi-

---

<sup>1</sup><http://www.rhsmith.umd.edu/industrydata/index.html>

ilarity for every possible pair of firms. For any two firms  $i$  and  $j$  in year  $t$ , the data provide a real number in the interval of  $[0, 1]$  describing how similar the words used by firms  $i$  and  $j$  are. Unfortunately, the publicly distributed data do not contain this real number, but rather firm links in the form of (gvkey1-gvkey2 combinations) for which the product similarity score exceeds a certain threshold. Thus, for each firm, the data we use give a list of firms that use similar words to describe their products and therefore are likely to operate in the same product market.

## II.2 Return Sample

Our first set of tests requires stock-level returns data. In order to construct it, we start from all publicly traded securities in the CRSP universe between 1995 and 2009 excluding stocks with sharecodes other than 10 or 11. We match this stock level dataset with the Compustat annual files in order to obtain each firm's *official* SIC code. We chose the historical SIC code from Compustat (*SICH*) to define a firm's official industry (alternative choices of industry classification are described as robustness checks). Guenther and Rosman (1994); Kahle and Walkling (1996) give some details about the rules used by Compustat to produce firm-level SIC codes based on the firms' 10K filings. Next, we match this CRSP-Compustat merged dataset with the HP text-based network industry classifications. These classifications are available between 1996 and 2008. In order to allow for the HP industry classifications to be known to investors and to eliminate the possibility of look ahead bias, we lag each firm's network industry classification by one year. Finally, we obtain firm level analyst coverage from I/B/E/S by counting the number of analysts who have issued at least one fiscal year end earnings forecast for a given firm throughout a year. In line with the literature, we set the number of analysts to zero whenever a firm included in the CRSP/Compustat sample cannot be matched with I/B/E/S. We end up with a sample of about 48,000 firm-year observations. For each of these firm-year observations, we then obtain monthly stock prices, returns and shares outstanding. This procedure yields a final sample of 565,266 firm-month observations

between 1997 and 2009. We also consider weekly returns. Our weekly sample has about 2,600,000 firm-week observations.

We define a firm  $j$ 's "official" industry return  $r_{j,official,t}$  as the equally weighted average of the returns during period  $t$  of all firms belonging to the same SIC2 category as firm  $j$ . We then define  $r_{j,fundamental,t}$ , which captures the average return of firms who are fundamentally similar to firm  $j$ . It is calculated as the equally weighted return of a portfolio consisting of all firms that are linked to firm  $j$  in the Hoberg and Phillips sense in that period. In other words,  $r_{j,fundamental,t}$  could be interpreted as capturing shocks to economic fundamentals relevant to firm  $j$ . By contrast, a firm's official industry portfolio does not necessarily represent firms that are fundamentally linked to each other, but rather firms that are grouped together according to the official industry classification, which can suffer from plain mistakes (another official industry would be more appropriate) or simply coarseness (a lot of firms are put in the same bag). We consider returns at both weekly and monthly frequencies. Table I provides summary statistics for the main return related variables used in this study.

[Table I about here.]

### II.3 Analyst Sample

In the second part of the paper, we examine whether analysts are also subject to industry categorization biases when they issue earnings forecasts. To address this question, we first construct a dataset of forecasts at the analyst level. We start from the I/B/E/S Detail History file and restrict ourselves to fiscal year end EPS forecasts for US firms issued during the current fiscal year (Forecast Period Indicator (FPI)=1). We obtain realized earnings from the Actuals files. As pointed out in Diether, Malloy, and Scherbina (2002), we account for stock splits and make the forecast and actual earnings time series comparable. In doing so, we rely on the CRSP cumulative adjustment split factor extracted from the CRSP Daily files, which is the most reliable and accurate way of performing the



split adjustment (see Robinson and Glushkov (2006)). Next, we remove analyst forecasts which seem to be erroneous (e.g. forecasts which are issued after actual earnings have been announced, forecasts for which actual earnings are announced before the fiscal year end, forecasts which have been reconfirmed before they have been issued for the first time, etc.). We also trim the tails of the EPS forecast distribution at 1% to reduce the impact of statistical outliers.

During the fiscal year, analysts tend to adjust their forecast of year end earnings at several points in time. Hence, each unit of observation is an analyst-firm-date where the forecast of fiscal year end earnings is reported. In each quarter of the year, we keep only the most recently issued or revised analyst-firm observation. Consider for instance an analyst  $i$  who, in the same quarter, issues two forecasts for a firm  $j$  with fiscal year end in December, the first one on January 29, 1996 and the second one on February 15, 1996. We keep only the most recent forecast, i.e. the forecast issued on February 15. We then organize the forecast data by forecast horizon, which we define as the difference between the date at which a forecast is issued and the fiscal year end date. We dismiss forecasts issued or revised after the fiscal year ends. Given that we keep one single analyst-firm observation per quarter, we end up with a maximum of four forecasts per analyst-firm-year which differ regarding their forecast horizons. Forecasts issued in the first fiscal quarter have a horizon of  $T = 3$  quarters, forecasts issued in the second quarter have a horizon of  $T = 2$  quarters, forecasts issued in the third quarter have a horizon of  $T = 1$  quarter and forecast issued in the last fiscal quarter have a horizon of  $T = 0$  quarters. We denote analyst  $i$ 's  $T$  quarter ahead forecast of firm  $j$ 's fiscal year end earnings issued in quarter  $t$  as  $F(T)_{i,j,t}$ .

### II.3.1 Forecast Errors

We now compute the accuracy associated with analyst  $i$ 's  $T$  quarter ahead forecast of firm  $j$ 's earnings issued in quarter  $t$  as

$$ForecastError(T)_{i,j,t} = (F(T)_{i,j,t} - A_{j,t})/P_{j,t-4},$$

where  $A_{j,t}$  denotes the firm's realized earnings per share at fiscal year end and  $P_{j,t-4}$  is the price that prevailed 12 month prior to fiscal year end.

In our tests, we also use firm-level forecast errors, which are obtained by collapsing analyst-firm level forecast errors with identical horizon  $T$  at the firm-quarter level. Formally, we define the average bias across firm  $j$ 's  $T$  quarter ahead EPS forecasts in quarter  $t$  as

$$ForecastError(T)_{j,t} = \frac{1}{N(T)_{j,t}} \sum_{i=1}^{N(T)_{j,t}} ForecastError(T)_{i,j,t},$$

where  $N(T)_{j,t}$  refers to the number of issued or revised  $T$  quarter ahead forecasts for firm  $j$  in quarter  $t$ . For a given forecast horizon  $T$ , we thus end up with one yearly observation per firm  $j$ .

### II.3.2 Consensus Forecasts

Next, we match the stock price that prevailed at the beginning of the quarter in which a forecast is issued and define the  $T$  horizon forecast to price ratio for analyst  $i$ 's EPS forecast for firm  $j$  in month  $t$  as:

$$FP(T)_{i,j,t} = F(T)_{i,j,t}/P_{j,t-1}.$$

### II.3.3 Firm-level consensus forecast

We then calculate the average forecast to price ratio at the firm level by averaging across all issued or revised forecasts for firm  $j$  in quarter  $t$  with a horizon of  $T$  quarters, i.e.

$$FP(T)_{j,t} = \frac{1}{N(T)_{j,t}} \sum_{i=1}^{N(T)_{j,t}} FP(T)_{i,j,t}.$$

### II.3.4 Industry-level consensus forecasts

For each firm-quarter, we finally calculate forecast to price ratios of the firm's "official" and "fundamental" industries. To do so, we match the forecast data with historical SIC codes from COMPUSTAT, and calculate, for each firm-quarter-horizon triplet, the average forecast to price ratio across all firms belonging to the same SIC2 industry. We will refer to this variable as  $FP(T)_{j,official,t}$ , the consensus  $T$  quarter ahead forecast in firm  $j$ 's official industry at date  $t$ . Analogously, we use the Hoberg and Philippps classification and calculate, for each firm  $j$  at date  $t$ , the average  $T$  quarter ahead forecast to price ratio of all firms that are related to firm  $j$  in the Hoberg and Phillips sense. We denote this variable as  $FP(T)_{j,fundamental,t}$ .

[Table II about here.]

Table II shows summary statistics for the main analyst related variables. Panel A contains variables at the forecast (analyst-firm-year) level. Panel B shows variables at the firm-year level (e.g. industry-level forecast to price ratios). Consistent with the existing literature, the average and median forecast errors at longer horizons are small, but positive. In contrast, bias decreases at shorter forecast horizons (e.g.  $T = 0$ ). Finally, Panel C shows variables with variation at the analyst-year level: the typical analyst tracks about three SIC2 industries and 11 stocks and has been providing forecasts for about six years. All variables are trimmed at 1% and 99%.

## III Stock-price over-reaction to official industry returns

The idea we have in mind is that investors in the stock market overemphasize the representativeness of official industry classifications. If categorical thinking is at work, we thus expect the following: In the short run, stocks exhibit comovement with their official industry peers in excess of what is granted by fundamentals. With time, investors reverse

mistakes induced by such categorical thinking, which then leads prices to revert toward economic fundamentals at longer horizons. Consider, for instance, a firm XYZ for which the official industry is only weakly representative. Such limited representativeness could be due to the firm selling products different from those offered by the typical firm in its industry, or simply by the coarseness of the classification. If some investors put all firms of that official industry in the same box when analyzing information flow, they will create excess comovement between XYZ's returns and its official industry returns. This comovement will subsequently be followed by a reversal toward XYZ's fundamentals. Hence, comovement between imperfectly related firms arises only because investors treat firms that belong to the same official industry category as being highly similar. Thus, our hypothesis is that at high frequencies, firms comove strongly with their official industry and weakly with their fundamental peers while the opposite pattern should hold at lower frequencies.

We first provide a graphical test of our hypothesis using weekly return data. For each stock, we compute cumulative returns over different horizons. Let  $r_{j,t}^T$  denote the  $T$  week cumulative return of firm  $j$ , between time  $t$  and  $t + T$ . Correspondingly,  $r_{j,fundamental,t}^T$  and  $r_{j,official,t}^T$  denote the cumulative return between week  $t$  and week  $t + T$  of the firm's official industry and that of the basket of Hoberg-Phillips comparables. We regress a firm's  $T$  week cumulative return on the firm's fundamental and the firm's official  $T$  week cumulative industry returns. This is done in a pooled regression with double clustered standard errors (time and stock dimension). Formally, we estimate the following equation for horizons  $T$  running from 1 to 8 weeks:

$$r_{j,t}^T = a_T + b_T \times r_{j,official,t}^T + c_T \times r_{j,fundamental,t}^T + e_{j,t}^T \quad (1)$$

[Insert Figure 1 here.]

We plot the coefficient estimates  $b_T$  and  $c_T$  against the return horizon  $T$  in figure 1.

The dashed lines represent 95 % confidence intervals. A strikingly clear pattern emerges from the picture: in the short run (one week horizon,  $T=1$ ), firm level returns are more strongly related to their official than to their fundamental industry returns. This is equivalent to saying that at the one week horizon, returns are more sensitive to official returns than they are to fundamental returns (as proxied by the returns to a portfolio of firm's HP comparables). By contrast, and as expected, this relationship reverts at longer return horizons: Firm returns become significantly more sensitive to their "fundamental" than to their "official industry" for horizons of three or more weeks ( $T > 2$ ). The coefficient estimate  $c_T$ , i.e. the sensitivity of cumulative stock returns with respect to the fundamental industry return increases monotonically in  $T$ , while, simultaneously, the sensitivity with respect to the official industry ( $b_T$ ) decreases monotonically. The figure thus shows that in the short term, stock returns comove *excessively* with their official industries, while in the medium and longer term, they revert toward fundamentals, thus becoming more correlated with HP-peer returns. This second comovement test also validates the use of HP-peers as a proxy for fundamentals: the Hoberg and Philipps classification explains stock returns at a one month horizon better than the SIC2, which suggests that it contains more information.

## III.1 Portfolio Tests

### III.1.1 Basic Test

The graphical evidence presented in the previous section suggests a strong pattern of return predictability based on the idea of industry categorization bias, which we now test by forming portfolios. In this second step, our identification strategy focuses on instances in which returns of a firm's *fundamental* ( $r_{j,fundamental,t}$ ) and *official* industries ( $r_{j,official,t}$ ) diverge. Categorization bias implies that a strong positive divergence between official and fundamental industry shocks (large  $r_{j,official,t} - r_{j,fundamental,t}$ ) should signal strongly negative returns, as the firm's stock price subsequently reverts to fundamentals. Note

that to eliminate any possibility of look ahead bias, we lag the HP industry classification by one year when constructing each firm's "fundamental" industry portfolio, and use COMPUSTAT's historical SIC.

We build a simple portfolio strategy based on this insight by sorting stocks according to the industry return differential ( $r_{j,official,t} - r_{j,fundamental,t}$ ). We start with weekly data. At the beginning of each week, we sort firms according to the return differential in  $t - n$ . In total, we consider lags of up to 6 weeks ( $n = 6$ ). The first quintile contains stocks for which the return differential is the most negative: the firm's fundamental industry return strongly exceeds that of its official industry. In the fifth quintile, official industry peers strongly outperform fundamental ones. We expect stocks in the first quintile portfolio to revert positively to their fundamentally linked firms, while firms in the fifth quintile should revert negatively. We restrict all tests to stock for which the price at the beginning of the period exceeds \$5. In addition, we require the official and the fundamental industries to be populated by at least five different firms.

[Table III about here.]

In table III we report four factor alphas alongside loadings for the Fama and French (1993) and Carhart (1997) factors for the Q1 and the Q5 portfolios. We also show results for a long-short (Q1-Q5) portfolio, which buys firms in the first and sells firms in the fifth quintile of the industry return differential distribution. All portfolios are equally weighted. In Panel A we sort firms on the industry return differential at the beginning of the week ( $t-1$ ). In panel B, we use the industry return differential at the beginning of the previous week ( $t-2$ ) and so forth. Consistent with the hypothesis that investors overreact to shocks to a firm's official industry, we find strong evidence that prices revert to their fundamentally linked firms following sizable divergence between official industry returns and their fundamentals. The weekly four factor alpha of an equally weighted portfolio which is long in firms with the most negative prior industry return differential is 28 basis points. With a t-stat of 7.18, this risk adjusted excess

return is highly statistically significant. Similarly, firms that have experienced the most positive industry return differential over the previous week revert negatively. The Q5 portfolio yields a statistically significant weekly alpha of -9 basis points (t-stat: -2.33). The long-short portfolio Q1-Q5 yielding 37 basis points is also highly significant (t-stat: 5.66). It corresponds to an annualized return of about 19%.

Since we construct portfolios according to the industry return differential at different lags, we are able to examine the frequency at which the reversal occurs. The table reveals that the speed of reversal depends strongly on whether the industry return differential is positive or negative. Stocks which have been subject to unjustified upward price pressure (High fundamental{ocial differential; Q5 firms) revert fully within a week's time. Beyond that time, the signal fades out: The risk adjusted excess returns are no longer statistically different from zero after the first week. By contrast, the "long" signal of the strategy is much more persistent: "Q1 firms" take much longer to revert completely. The weekly equally weighted alpha decreases somewhat gradually from about 28 basis in the first, to about 16 basis points in the sixth week. The signal no longer produces significant risk adjusted excess returns for the long-short portfolio after three weeks, suggesting that the mispricing dissipates roughly within a month's time. Guided by this observation, we reproduce the portfolio analysis at the monthly level in the appendix (see table A.I) yielding identical conclusions, and we will focus on monthly returns in the remainder of the text.

[Table IV about here.]

### III.1.2 Sorting by Size

Bounded rationality theories predict that agents have the choice to "reduce their bias" when the benefit of doing so outweighs the cost (see, for instance, Gabaix (2010)). In the present context, we therefore expect that categorization bias is less pronounced for large caps, as the dollar value of mispricing generated by a similar cognitive mistake is higher for large than small firms. In table IV, we turn to examining risk adjusted excess returns

of value weighted portfolios. Following up on the prior insights regarding the speed of the reversal, we now focus on monthly returns. We first sort all firms according to their market capitalization in the prior month and group them into terciles. Within each size tercile, we then sort firms into five quintile portfolios based on their industry return differential over the previous month. In total, we obtain 15 size-industry differential and three long-short portfolios. Panel A, B and C report the results for all value weighted portfolios consisting of small, medium and large capitalization firms. The results show significant monthly portfolio alpha for long-short portfolios consisting of small (154 basis points;  $t=4.25$ ) and mid cap firms (119 basis points;  $t=2.39$ ). By contrast, value weighted portfolios containing large cap firms do not yield significant risk adjusted excess returns. All in all, categorization bias seems to be prevalent in the bottom two terciles of firm size, and rather absent in the top tercile, which is largely consistent with bounded rationality.

### III.1.3 Refining the strategy

If categorization bias explains the stock price anomaly we document, the pattern of reversal should be more pronounced when firms have indeed closely followed their "official" industry. This would be the case, for instance, if stock miscategorization fluctuates over time. Sometimes, the stock is properly categorized, and even though fundamental and official peers diverge, the stock follows its fundamental peers and therefore does not revert. Sometimes, the stock is miscategorized, and hence the peer-fundamental divergence signals future reversal.

To test this, we refine our strategy by the means of a double sort. We measure (1) the official-fundamental divergence as above, and (2) the extent to which the return of a stock has been abnormally close to its official industry return in a given month. Our second signal is computed as follows: For each firm-month pair, we define an *Official Industry Follower* measure as

$$\rho_{j,t} = \frac{r_{j,t} - r_{j,official,t}}{\sigma_{j,24}(r_{j,t} - r_{j,official,t})} ,$$



where  $r_{j,t}$  is firm  $j$ 's return in month  $t$ ,  $r_{j,official,t}$  denotes the equally weighted return of the firm's official industry and  $\sigma_{j,24}(r_{j,t} - r_{j,official,t})$  is the standard deviation of the difference between firm  $j$ 's return and the return of the firm's official industry calculated for rolling windows of 24 months. Low values of  $\rho$  indicate that firm  $j$ 's stock price change has closely followed the price changes in its official industry in month  $t$ , while large values indicate divergence. We choose this measure over the rolling correlation between  $r_{j,t}$  and  $r_{j,official,t}$  since we are ultimately interested in identifying cases in which stocks that typically do not follow their official industry-level returns have been abnormally close followers of their official industries. Relying on rolling correlations and selecting stocks for which the rolling correlation is high would identify firms that generally tend to follow their official industries, while we are interested in identifying cases in which official industry following is unusual (low values of  $\rho$ ).

We now build a portfolio strategy based on a combination of the *Official Industry Follower* measure  $\rho_{j,t}$  and the industry return differential  $r_{j,official,t} - r_{j,fundamental,t}$ . We first sort firms into three terciles of  $\rho$  over the previous month. The first tercile marks *Strong Official Industry Follower*, i.e. stocks that have moved unusually strongly in line with their official industries over the past month. In contrast, the third tercile of  $\rho$  indicates stocks which have not followed their official industry (*Weak Official Industry Follower*) over the previous month. Within each tercile of  $\rho$ , we then sort firms into quintiles of the previous month's industry return differential. Thus, the signal we use in order to construct the refined portfolio strategy is based on the following interaction term  $(r_{j,official,t-1} - r_{j,fundamental,t-1}) \times \rho_{j,t-1}$ . Figure 2 illustrates the idea of the interacted signal in a simple diagram.

[Figure 2 about here.]

We expect return predictability to be strongest among *Strong Official Industry Follower* stocks. This is because firms which have not followed their official industries (*Low Official Industry Follower*) are more likely of having moved in line with their fundamentally linked firms, reducing the scope for return reversal. Hence, the magnitude of the

return reversal should be increasing in the extent to which a stock has moved in line with its official industry peers.

[Table V about here.]

In table V, we report the risk adjusted returns and factor loadings for equally weighted portfolios, which are constructed according to a double sort on  $\rho_{j,t-1}$  and  $r_{j,official,t-1} - r_{j,fundamental,t-1}$ . The results show that return predictability, as evidenced by the risk adjusted excess returns of the long-short (Q1-Q5) portfolio, is monotonically declining in the *Official Industry Follower* measure: the effect is most pronounced among *Strong Official Industry Follower* stocks (158 basis points;  $t=4.10$ ; see Panel A) and weakest for *Weak Official Industry Follower* stocks (77 basis points;  $t=1.73$ ; see Panel C). A monotonically decreasing effect is consistent with the view that stock reversal to fundamentals is more pronounced if a stock has unusually closely tracked its official industry in the first place.

[Table VI about here.]

We now restrict the analysis to *Strong Official Industry Follower* stocks (first tercile of  $\rho_{j,t-1}$ ), and form value weighted portfolios at the monthly level. At the beginning of each month, we first sort all *Strong Official Industry Follower* stocks into terciles of market capitalization in the previous period. Then, within a size tercile, we sort on the industry return differential. Comparing value weighted long-short portfolios restricted to *Strong Official Industry* followers with portfolios which are not conditioned on the *Official Industry Follower* measure (see Panels A, B and C of table IV) shows that, indeed, reversal is much more pronounced in both economic and statistical terms for *Strong Industry Follower* stocks. For mid cap stocks, for instance, the monthly alpha of the value weighted long-short portfolio increases from 119 ( $t=2.39$ ; see Panel B of table IV) to 160 basis points ( $t=3.05$ ; see Panel B of table VI). This increase is even more pronounced when looking at small capitalization firms. In the appendix (see Table A.II),

we evaluate equally weighted portfolios restricted to *Strong Official Industry Followers* at the weekly frequency. The same dramatic increase of both economic and statistical significance for official industry movers is evidenced at the higher frequency too: four factor alpha for the long-short portfolio in the first week increases from 37 (t=5.66; see Panel A of table III) to 48 (t=7.34; see Panel A of table A.II) weekly basis points, when we additionally condition on the *Official Industry Follower* measure over the previous week.

### III.2 Regression Tests

As an alternative to the calendar time portfolio approach, we now use cross sectional regression frameworks in order to test our hypothesis. Relying on regression tests allows us to control more thoroughly for other potential cross sectional determinants of stock returns, most notably the firm's lagged return  $r_{j,t-1}$ , the firm's past cumulative returns  $r_{j,t-2,t-12}$ , the log of the firm's market capitalization and a firm's book to market equity. We focus on monthly returns in the regression tests.

[Table VII about here.]

In column (1) of table VII, we start by regressing the monthly raw return on the industry return differential in the previous month. We double cluster standard errors in the month and stock dimension in the spirit of Petersen (2009) and also include month fixed effects in the equation. As expected, the coefficient estimate for  $r_{j,official,t-1} - r_{j,fundamental,t-1}$  has a negative sign and is statistically significant: when the official industry has outperformed the fundamental one, future returns are negative as the stock reverts to its fundamental value. Mirroring the portfolio analysis, we then refine the analysis in column (2) of table VII by interacting the return differential ( $r_{j,official,t-1} - r_{j,fundamental,t-1}$ ) with dummy variables indicating *Strong*, *Medium* and *Weak Official Industry Follower* stocks. Consistent with the portfolio tests, the return reversal effect is most pronounced among *Strong Official Industry Follower* stocks (the

reference category). The marginal effects for the reference category increases both in economic magnitude and statistical significance ( $t=-3.24$ ) when contrasted with the pooled estimate in column (1). In line with the evidence presented in the portfolio analysis, the effect of the industry return differential decreases monotonically in the three levels of *Official Industry Follower*. Equations (3)-(6) implement further controls, but to make results easier to read, we focus on *Strong Official Industry Follower* stocks, for which the mean reversion is the most pronounced. Column (3) shows an economically strong and statistically significant impact of the industry return differential even after controlling for commonly used cross sectional determinants of stock returns. Column (4) interacts the industry return differential with market capitalization and confirms that the categorization bias is weaker for large capitalization firms, which is consistent with bounded rationality as noted earlier. In column (5), we test whether information production, as proxied by analyst coverage, tends to attenuate the bias. We find no evidence in favor of this view.

One potential alternative interpretation of our results could be that return predictability is not due to the categorization bias we have in mind, but rather due to indexing by passive exchange traded funds. Assume for instance that industry-focused ETFs receive flows based on fundamental information about the representative firm of their industry. In this case, such ETFs would exert price pressure on the representative firm but the price would not revert as the flows were motivated by fundamental trading. Stocks that are misclassified in the industry, however, would first comove with the industry portfolio, and then return to their fundamental value. Misclassified stocks would comove too much with their official industry because of ETF flows. The problem would be less one of investor categorization, than one of financial market imperfection due to the low granularity of ETF definitions.

[Figure 3 about here.]

We provide a simple test in column (6). One implication of this alternative interpretation is that the effect should be stronger in recent periods. Figure 3 plots the monthly

dollar volume of the *SPDR S&P 500* (Ticker Symbol: *SPY*) as a fraction of total monthly dollar volume of the CRSP Universe and shows that growth in ETF trading has picked up remarkably after 2002. We code a dummy variable that marks all firm-month observations after 2002 and interact it with the industry return differential. The interaction effect is not statistically significant, suggesting that our results are not driven by indexing.

Furthermore, in the appendix (see table A.III), we re-estimate the first five specifications of the previous table in a Fama and MacBeth (1973) framework. In calculating the t-statistics, we allow for heteroscedasticity and serial correlation of up to 12 months in the spirit of Ponti (1996). The coefficient estimates resulting from the Fama-Macbeth (FMB) regressions are quite similar in terms of their economic magnitudes and statistical significance is, if anything, stronger in the FMB framework.

### III.3 A Placebo Test: Using Pseudo-Fundamental Peers

A potential concern in interpreting our results is that they might be exclusively driven by the returns of the official industry portfolio  $r_{j,official,t}$  rather than being due to divergence between the official and the fundamental industry portfolio returns  $r_{j,official,t} - r_{j,fundamental,t}$ . The observed reversal in the comovement intensity of stocks with their official and fundamental peers (Figure 1) already alleviates this concern somewhat. However, we also address this issue more directly within a portfolio analysis framework by constructing a pseudo Hoberg and Phillips (2010a,b) industry classification by the means of the following algorithm: For each firm  $j$  included in the TNIC data in year  $t$ , we randomly select  $l_{jt}$  different firms from the Compustat universe in year  $t$ , where  $l_{jt}$  is the number of Hoberg and Phillips links of firm  $j$  in year  $t$ . We then calculate the equally weighted return to a portfolio consisting of these pseudo Hoberg and Phillips firms and denote the pseudo Hoberg and Phillips portfolio return as  $r_{j,pseudo,t}$ . For each firm-month observation, the pseudo portfolio contains the same number of firms as the Hoberg and Phillips industry portfolio. Yet, the firms in the pseudo portfolio have no meaningful relation to firm  $j$  because they are randomly chosen. We now reproduce the portfolio analysis

for the pseudo return differential by first sorting stocks according to whether they track their official industries over the previous month (i.e. terciles of  $\rho_{j,t-1}$ ) and then sorting stocks within a given tercile of  $\rho_{j,t-1}$  by quintiles of  $r_{j,official,t-1} - r_{j,pseudo,t-1}$ . Hence, instead of sorting on the industry-return differential, we now sort stocks into quintiles of  $r_{j,official,t} - r_{j,pseudo,t}$ . All portfolios are equally weighted.

[Table VIII about here.]

Panels A, B and C of table VIII show risk adjusted excess returns and factor loadings for *Strong*, *Medium* and *Weak Industry Follower* stocks. The pseudo strategy goes long in stocks for which  $r_{j,fundamental,t-1} - r_{j,pseudo,t-1}$  is the most negative over the previous month and sells stocks for which the pseudo industry return difference is most positive. Assuming that  $r_{j,pseudo,t-1}$  is just random noise, such a strategy is akin to an inverse industry momentum strategy (buying short term losers and selling short term winners). Consistent with the idea of industry momentum (see Moskowitz and Grinblatt (1999)), we find that the long-short portfolio Q1{Q5 based on the pseudo return differential yields statistically significant negative four factor alpha. Negative alpha is in stark contrast to the positive risk adjusted returns we document for strategies based on  $r_{j,official,t-1} - r_{j,pseudo,t-1}$ . In unreported regression tests, we obtain a statistically significantly positive coefficient estimate for  $r_{j,official,t-1} - r_{j,pseudo,t-1}$ , which again is in stark contrast to the negative coefficient for  $r_{j,official,t-1} - r_{j,fundamental,t-1}$ . This evidence suggests that it is not  $r_{j,official,t-1}$ , but rather the relationship between the official and the fundamental industry returns, i.e.  $r_{j,official,t-1} - r_{j,fundamental,t-1}$  that is driving our results.

### III.4 Additional Robustness Checks

In untabulated analysis, we define a firm's official industry at the SIC3 level obtaining similar results. We also obtain qualitatively similar results by using the Global Industry Classification Standard (GICS) instead of the Standard Industrial Classification. In additional robustness checks, we test whether our results are sensitive to defining industry

level returns as medians, equally weighted or value weighted averages and find this not to be the case. In a last robustness check, we also find that using DGTW characteristics-adjusted stock returns (see Daniel, Grinblatt, Titman, and Wermers (1997)) leads to similar conclusions.

## IV Industry Categorization Bias in Analysts' Forecasts

The aim of this section is to provide evidence of categorization bias directly based on analyst expectations, instead of looking at prices. Here, we take advantage of the wealth of data on realized and expected analyst forecasts. The empirical approach will be similar to the one used in the stock returns section, except that we replace returns by true expectations.

Ideally, in the spirit of Table III, we would want to regress EPS forecast errors on the divergence between "official" and "fundamental" industry forecasts. The idea is that if expectations about "official" peers are, say, better than expectations about "fundamental" peers and *if* analysts suffer from categorization bias, analyst forecasts should be too optimistic: the average firm-level EPS forecast error should be positive. Put differently, divergence between official and fundamental industry forecast predicts systematic biases in expectations, an apparent violation of the rational expectation hypothesis. This naive version of the test is, however, difficult to implement as forecast errors are affected by many observable and unobservable determinants.

This is why we adopt the approach of Table V: we test if forecast errors are *more* affected by official{fundamental industry divergence, *when* analyst expectations tend to track the "official" industry consensus more closely. This conditional approach is akin to a difference-in-difference methodology and identical to the test we run on stock returns in Table V. The main advantage of this approach is that it refines the predictive power of the categorization hypothesis as it allows focusing on instances where (1) the scope for

categorization error is the biggest (official and fundamental industries diverge) and (2) analyst do focus a lot on what is going on in the official industry. In addition, this refined approach allows to better control for unobservable determinants of forecast errors.

First, we define our measure capturing how close forecasts for firm  $j$  are at time  $t$  to the average forecast prevailing in its the official industry as:

$$\rho(T)_{j,t} = \frac{|FP(T)_{j,t} - FP(T)_{j,official,t}|}{\sigma_t(FP(T)_{j,t} - FP(T)_{j,official,t})}.$$

$FP(T)_{j,t}$  is the firm's consensus forecast to price ratio and  $FP(T)_{j,official,t}$  denotes the average forecast to price ratio prevailing among a firm's official peers. Since we study different forecast horizons, this measure is also indexed by the horizon  $T$ .  $\sigma_t$  denotes the cross-sectional standard deviation.  $\rho(T)_{j,t}$  measures whether the consensus forecast for firm  $j$  is closer to the official industry forecast for firm  $j$  than for other firms. In other words, low values of  $\rho$  indicate that analysts tend to be close to the official industry level average forecast for firm  $j$ , while large values indicate that the consensus forecast for firm  $j$  diverges from its official industry-level consensus forecast. In contrast to the returns section, we calculate  $\rho(T)_{j,t}$  by normalizing by the cross sectional standard deviation. We do so because the lower frequency of analyst forecasts makes it difficult to use the time series of past forecasts as we did in our stock returns tests. Thus,  $\rho(T)_{j,t}$  closely mirrors the measure we adopted in our stock returns tests of Table V.

We then implement our test by running the following regression:

$$\begin{aligned} ForecastError(T)_{j,t} = & \alpha + \beta(FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}) + \gamma\rho(T)_{j,t} \\ & + \delta\rho(T)_{j,t} \times (FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}) + \lambda'X_{j,t} + \theta_t + \epsilon_{j,t}. \end{aligned}$$

where  $ForecastError(T)_{j,t}$  is the average  $T$  horizon firm-level forecast-error,  $FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}$  is the divergence between forecasts of official and fundamental peers, and  $\rho(T)_{j,t}$  is our measure of analysts' official industry focus. We run these regressions separately for different forecast horizons  $T = 0, 1, 2, 3$ . In the specifications, we also



control for a number of firm level characteristics (e.g. book-to-market, size, earnings volatility, etc.) that have been used in the literature to explain forecasts errors (see Hong and Kacperczyk (2010)) and include year fixed effects. The appendix contains a complete list of the firm-level controls used.

In the above specification, our hypothesis is that  $\delta < 0$ . If analysts have the tendency to closely track the official industry level forecast for a firm (small  $\rho(T)_{jt}$ ) and there is strong positive divergence between official and fundamental earnings prospects ( $FP(T)_{j,official,t} - FP(T)_{j,fundamental,t} > 0$ ), analysts will tend to be overly optimistic.

[Table IX about here.]

We report the results for several specifications and all four forecast horizons in table IX. In columns (1)-(4), we examine long horizon forecasts, i.e. forecasts issued in the first fiscal quarter (horizon  $T = 3$ ), while columns (5), (6) and (7) consider bias associated with shorter horizon forecasts ( $T = 2, 1, 0$  quarters respectively). In columns (1) and (2), we regress the three quarter ahead firm-level forecast error on the official-fundamental divergence. As noted before, the sign of this coefficient is difficult to interpret because of unobservable characteristics influencing forecast errors: whether we control for observable characteristics or not, it is statistically indistinguishable from zero. In columns (3)-(7), we interact the fundamental-official divergence with our firm-level measure of analysts' official industry focus  $\rho(T)_{jt}$ . Column (3) looks at a forecast horizon of 3 quarters ( $T = 3$ ): the coefficient estimate for the interaction term is significantly negative. This is consistent with what we expect: whenever the fundamental forecast is lower than the official one ( $FP(T)_{j,official,t} - FP(T)_{j,fundamental,t} > 0$ ), and analysts focus a lot on official peers ( $\rho(T)_{jt}$  small), then forecast errors are very positive. In these instances, analysts are too optimistic and their bias can be systematically predicted. In column (4) we add firm level control variables to the equation which leaves our conclusions unchanged. In columns (5) to (7) we repeat this exercise for shorter forecast horizons ( $T = 2, 1, 0$ ). The interaction term is significant for a forecast horizon of two quarters and becomes insignificant at

shorter horizons. The observation that categorization bias decays at shorter horizons is consistent with a bounded rationality argument: at shorter horizons, public information about future EPS is more abundant (in particular because of quarterly accounts). Hence, it is less costly for analyst to debias their expectations, so they become more rational.

We now study whether the tendency for analysts to be biased is stronger for specific kinds of analysts. Again, we try to isolate at the firm-level cases in which anchoring is particularly problematic (high official-fundamental divergence), and ask which analysts have, in these instances, a tendency to track the official industry in their forecasts. It should be noted that anchoring on a firm's official industry consensus is not bad per se. This is because it might be a good strategy to use the official industry average as an anchor whenever the firm in question is very similar to its official industry peers. By contrast, official industry emphasis is problematic at times when a firm's official and true industries differ substantially.

To implement this test, we first need to calculate the extent to which analyst  $i$ 's forecast for firm  $j$ 's year end earnings is close to the official industry by adapting our  $\rho$  measure to individual analyst forecasts:

$$\rho(T)_{i,j,t} = \frac{|FP(T)_{i,j,t} - FP(T)_{j,official,t}|}{\sigma_{jt}(FP(T)_{i,j,t} - FP(T)_{j,official,t})}.$$

In calculating  $\rho(T)$  at the analyst level, we now choose to normalize by the cross-sectional standard deviation at the firm-year level ( $\sigma_{jt}$ ). We do so because we are ultimately interested in capturing whether an analyst  $i$  compared to other analysts forecasting earnings for the same firm  $j$  tends to anchor more strongly on the official industry.

We then regress  $\rho(T)_{i,j,t}$  on an interaction term between the absolute fundamental-official divergence  $|FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}|$  and observable analyst characteristics. We expect more rational analysts to deviate more from the official industry consensus, when such deviations are justified, i.e. the official category provides a bad approximation (high divergence  $|FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}|$ ). In other words, more

rational analysts should have higher  $\rho(T)_{i,j,t}$  when  $|FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}|$  is relatively high. In all regressions, in addition to the *industry divergence*  $\times$  *analyst characteristics* interaction term, we include as controls several firm and other analyst characteristics (e.g. number of stocks or industries covered by the analyst, etc.) as well as analyst and year fixed effects.

[Table X about here.]

Taken together, our results suggest that more informed analysts tend to be more rational (i.e. deviate from the consensus) when such deviations are required (high fundamental-official divergence). We report results in Table X, and focus on long horizon forecasts ( $T=3$ ), where biases seem to be the strongest according to results from Table IX. We first regress  $\rho(3)_{i,j,t}$  on the industry difference and year fixed effects. The coefficient estimate is significantly positive, suggesting that on average analysts deviate from the official consensus (see column (1)). This suggests that, on average, when official and fundamental signals diverge, analyst forecasts tend to track the official industry less closely ( $\rho(3)_{i,j,t}$  is bigger). On average, analysts trust official signals less when they are less relevant. As we have seen in the previous Table, this correction is not large enough to prevent the consensus from being biased, but analysts may differ in their ability to correct for miscategorization. This is what we investigate in columns (2)-(7). In column (2), we look at analyst experience. We measure experience as the number of years an analyst appears in the I/B/E/S database. The interaction term is positive and significant at 10%; This is weakly consistent with the view that more experienced analysts rely less on the "official" category when official and fundamental consensus diverge. In column (3), we interact the absolute industry difference with firm-level experience of the analyst defined as the number of years analyst  $i$  has been issuing forecasts for firm  $j$ . Again, we find a statistically significant positive sign. These first two results suggest that experience reduces biases.

Next, we look at analyst's breadth of information. The effect of following a broad set of stocks is ambiguous. On the one hand, it could create cognitive overload, making

it more difficult to correct categorization bias. On the other hand, it could provide the analyst with adequate benchmarks to form more rational expectations. In column (4), we interact official-fundamental divergence with the number of stocks covered by the analyst. A broader view seems to be beneficial: We find a statistically significant positive sign for the interaction effect between the absolute industry difference and the number of stocks covered by analyst  $i$  in year  $t$ . Similarly, analysts who cover stocks from a higher number of official industries also tend to be more rational when a firm appears to be only limitedly representative of its official sector (see column (5)). Covering more official industries indicates to some extent that an analyst's understanding of comparable firms depends to a lesser extent on official industry classifications. In column (6), we seek to capture whether an analyst is more focused on the official or fundamental peers of firm  $j$ . We expect that, when an analyst follows many fundamental peers of a firm, she will form more rational expectations, in particular when it is important to do so. To implement this, we calculate the number of official peers  $\# Official_{i,j,t}$ , as the number of firms belonging to the same SIC2 industry as firm  $j$  that are also covered by analyst  $i$  in year  $t$ . Similarly, the number of fundamental peers  $\# Fundamental_{i,j,t}$  is the number of firms covered by analyst  $i$  that are linked to firm  $j$  in the Hoberg and Phillips sense. We normalize the difference between official and fundamental peers by the number of stocks covered by the analyst. We label this measure  $Coverage_{i,j,t} = (\# Official_{i,j,t} - \# Fundamental_{i,j,t}) / Stocks_{i,j,t}$ . We expect that whenever an analyst's stock universe is tilted more toward official peers (positive values for  $Coverage$ ), she will not deviate enough from official industry averages when in fact such deviation would be rational (high absolute industry difference). In line with this argument we find a negative coefficient which is highly statistically significant for the interaction between  $Coverage$  and  $|FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}|$  (see column (6)).

Finally, we examine whether the size of the brokerage house at which the analyst is employed has an impact on excessive industry anchoring. We find no evidence of such an effect (see column (7)). Altogether, our results suggest that analysts with more informa-

tion (more experience, following a broader set stock, following more relevant comparables) tend to be less subject to categorization bias.

## V Conclusion

This paper is part of the growing literature on the "economics of inattention" that studies the real effects induced by cognitive costs of information processing. In the first part of the paper, we explore the stock-price effects of industry categorization bias. If some investors mentally group firms belonging to the same official industry category when processing new information, one should expect over-reaction to official industry shocks and under-reaction to fundamental shocks that are not captured by the official industry classification. To test this hypothesis, we compute the difference between returns of a firm's official industry and a portfolio of its fundamentals peers, which we define according to whether firms operate in similar product markets (see Hoberg and Phillips (2010a,b)). In line with the existence of industry categorization bias, we find that, in the short run firms comove strongly with their official industry and weakly with their fundamental peers, whereas this pattern reverses at longer horizons. Furthermore, by constructing portfolios based on this insight, we show that divergence between a firm's official industry returns and those of its fundamental peers strongly predicts a firm's subsequent reversal toward fundamentals. The long-short strategy based on this signal generates highly significant risk adjusted excess returns, which become even larger when conditioning on whether a firm's stock price has moved closely in line with its official industry in the previous month. We thus conclude that the overemphasis of official industry classifications by boundedly rational investors creates large and predictable short-term deviations of prices from fundamentals, resulting in cross sectional return predictability.

Second, by studying analyst forecasts errors, we explore the extent to which financial analysts are themselves subject to industry categorization biases. We provide evidence that analysts tend to overemphasize official industry categories, resulting in cross sectional

predictability of earnings forecast errors for firms which are poorly representative of their official industry.

## References

- Barberis, N., A. Shleifer, and J. Wurgler, 2005, Comovement, *Journal of Financial Economics* 75, 283{317.
- Carhart, Marc M., 1997, On Persistence in Mutual Fund Performance, *Journal of Finance* 52, 57{82.
- Cen, L., G. Hilary, and K.C. Wei, 2011, The role of anchoring bias in the equity market: Evidence from analysts' earnings forecasts and stock returns, forthcoming in *Journal of Financial and Quantitative Analysis*.
- Chan, L.K.C., J. Lakonishok, and B. Swaminathan, 2007, Industry classifications and return comovement, *Financial Analysts Journal* pp. 56{70.
- Chan, W.S., 2003, Stock price reaction to news and no-news: drift and reversal after headlines, *Journal of Financial Economics* 70, 223{260.
- Clement, M.B., 1999, Analyst forecast accuracy: Do ability, resources, and portfolio complexity matter?, *Journal of Accounting and Economics* 27, 285{303.
- Cohen, Lauren, and Andrea Frazzini, 2008, Economic Links and Predictable Returns, *The Journal of Finance* 63, 1977{2011.
- Cohen, Lauren, and Dong Lou, 2010, Complicated firms, forthcoming in the *Journal of Financial Economics*.
- Da, Z., J. Engelberg, and P. Gao, 2011, In search of attention, *The Journal of Finance* 66, 1461{1499.

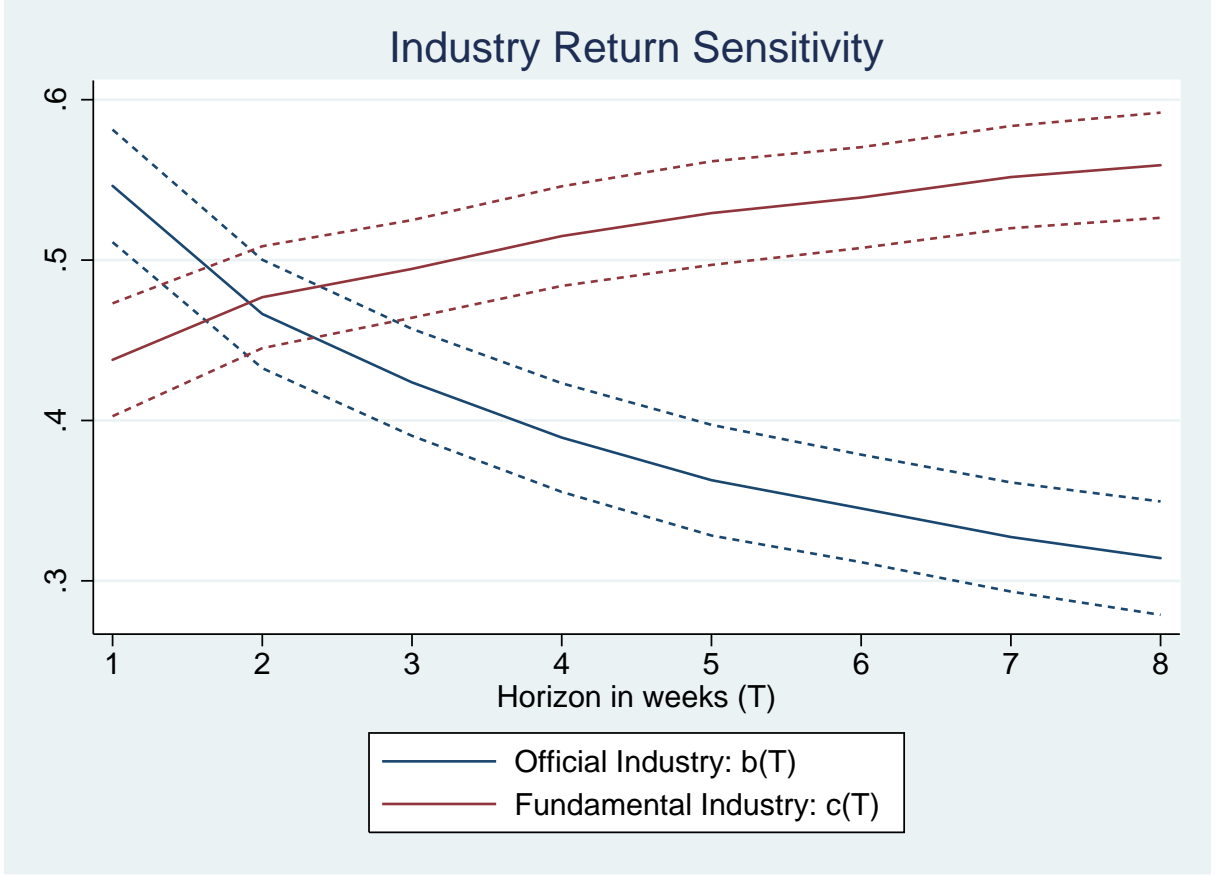
- Daniel, K., M. Grinblatt, S. Titman, and R. Wermers, 1997, Measuring mutual fund performance with characteristic-based benchmarks, *Journal of Finance* pp. 1035{1058.
- Diether, K.B., C.J. Malloy, and A. Scherbina, 2002, Differences of opinion and the cross section of stock returns, *The Journal of Finance* 57, 2113{2141.
- Fama, E.F., and J.D. MacBeth, 1973, Risk, return, and equilibrium: Empirical tests, *The Journal of Political Economy* pp. 607{636.
- Fama, Eugene F., and Kenneth French, 1993, Common risk factors in the returns on bonds and stocks, *Journal of Financial Economics* 33, 3{53.
- Foster, G., 1981, Intra-industry information transfers associated with earnings releases, *Journal of Accounting and Economics* 3, 201{232.
- Gabaix, Xavier, 2010, A Sparsity-Based Model of Bounded Rationality, Working Paper, Stern School of Business, New York University.
- Guenther, D.A., and A.J. Rosman, 1994, Differences between compustat and crsp sic codes and related effects on research, *Journal of Accounting and Economics* 18, 115{128.
- Hoberg, G., and G. Phillips, 2010a, Product market synergies and competition in mergers and acquisitions: A text-based analysis, *Review of Financial Studies* 23, 3773.
- Hoberg, Gerard, and Gordon Phillips, 2010b, Text-Based Network Industries and Endogenous Product Differentiation, University of Maryland Working Paper.
- Hong, H., and M. Kacperczyk, 2010, Competition and bias, *The Quarterly Journal of Economics* 125, 1683.
- Hong, H., and J.D. Kubik, 2003, Analyzing the analysts: Career concerns and biased earnings forecasts, *The Journal of Finance* 58, 313{351.

- 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, 265{295.
- Hong, H., J. Stein, and J. Yu, 2007, Simple forecasts and paradigm shifts, *Journal of Finance*.
- Hong, H., W. Torous, and R. Valkanov, 2007, Do industries lead stock markets?, *Journal of Financial Economics* 83, 367{396.
- Hou, K., 2007, Industry information diffusion and the lead-lag effect in stock returns, *Review of Financial Studies* 20, 1113.
- , and T.J. Moskowitz, 2005, Market frictions, price delay, and the cross-section of expected returns, *Review of Financial Studies* 18, 981{1020.
- Ivkovic, Z., and N. Jegadeesh, 2004, The timing and value of forecast and recommendation revisions, *Journal of Financial Economics* 73, 433{463.
- Jegadeesh, N., 1990, Evidence of predictable behavior of security returns, *Journal of Finance* 45, 881{98.
- Kahle, K.M., and R.A. Walkling, 1996, The impact of industry classifications on financial research, *Journal of Financial and Quantitative Analysis* 31, 309{335.
- Lo, A.W., and A.C. MacKinlay, 1990, When are contrarian profits due to stock market overreaction?, *Review of Financial studies* 3, 175.
- Malloy, C.J., 2005, The geography of equity analysis, *The Journal of Finance* 60, 719{755.
- Menzly, L., and O. Ozbas, 2006, Cross-industry momentum, manuscript, University of Southern California.
- Moskowitz, T.J., and M. Grinblatt, 1999, Do industries explain momentum?, *The Journal of Finance* 54, 1249{1290.

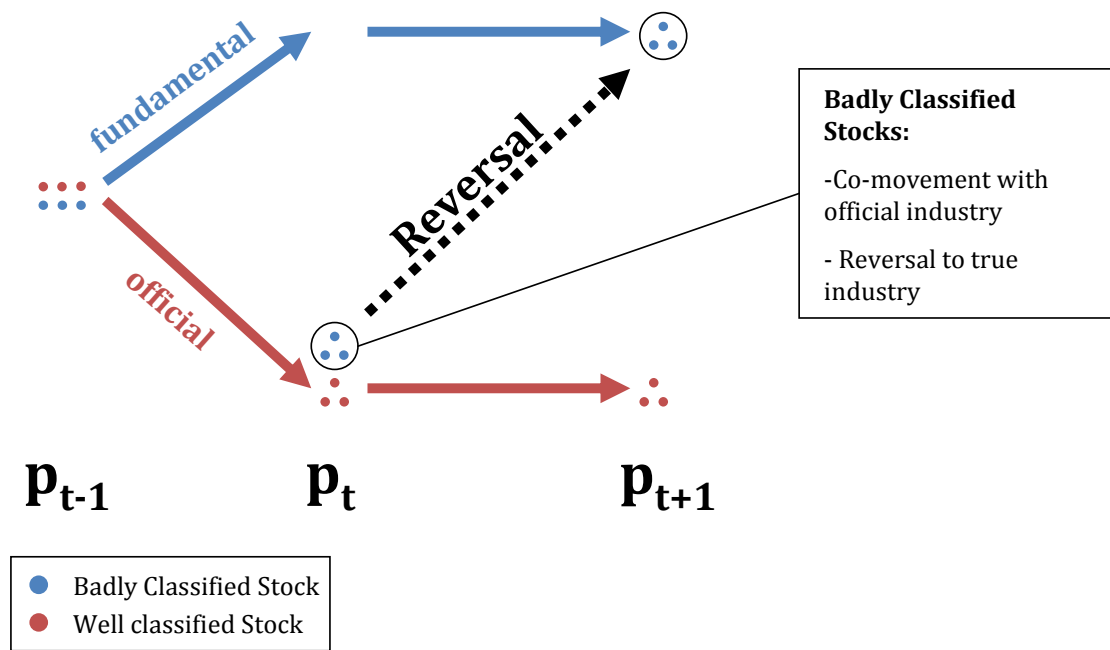


- Mullainathan, S., 2002, Thinking through categories, *mimeo MIT*.
- , A. Shleifer, and J. Schwartzstein, 2008, Coarse thinking and persuasion, *Quarterly Journal of Economics* pp. 577{619.
- Payne, J.L., and W.B. Thomas, 2003, The implications of using stock-split adjusted i/b/e/s data in empirical research, *Accounting Review* pp. 1049{1067.
- Peng, L., and W. Xiong, 2006, Investor attention, overconfidence and category learning, *Journal of Financial Economics* 80, 563{602.
- Petersen, M.A., 2009, Estimating standard errors in finance panel data sets: Comparing approaches, *Review of financial studies* 22, 435{480.
- Ponti, J., 1996, Costly arbitrage: Evidence from closed-end funds, *The Quarterly Journal of Economics* 111, 1135.
- Ramnath, S., 2002, Investor and analyst reactions to earnings announcements of related firms: An empirical analysis, *Journal of Accounting Research* 40, 1351{1376.
- Robinson, D., and D. Glushkov, 2006, A note on ibes unadjusted data, Wharton Research Data Services.

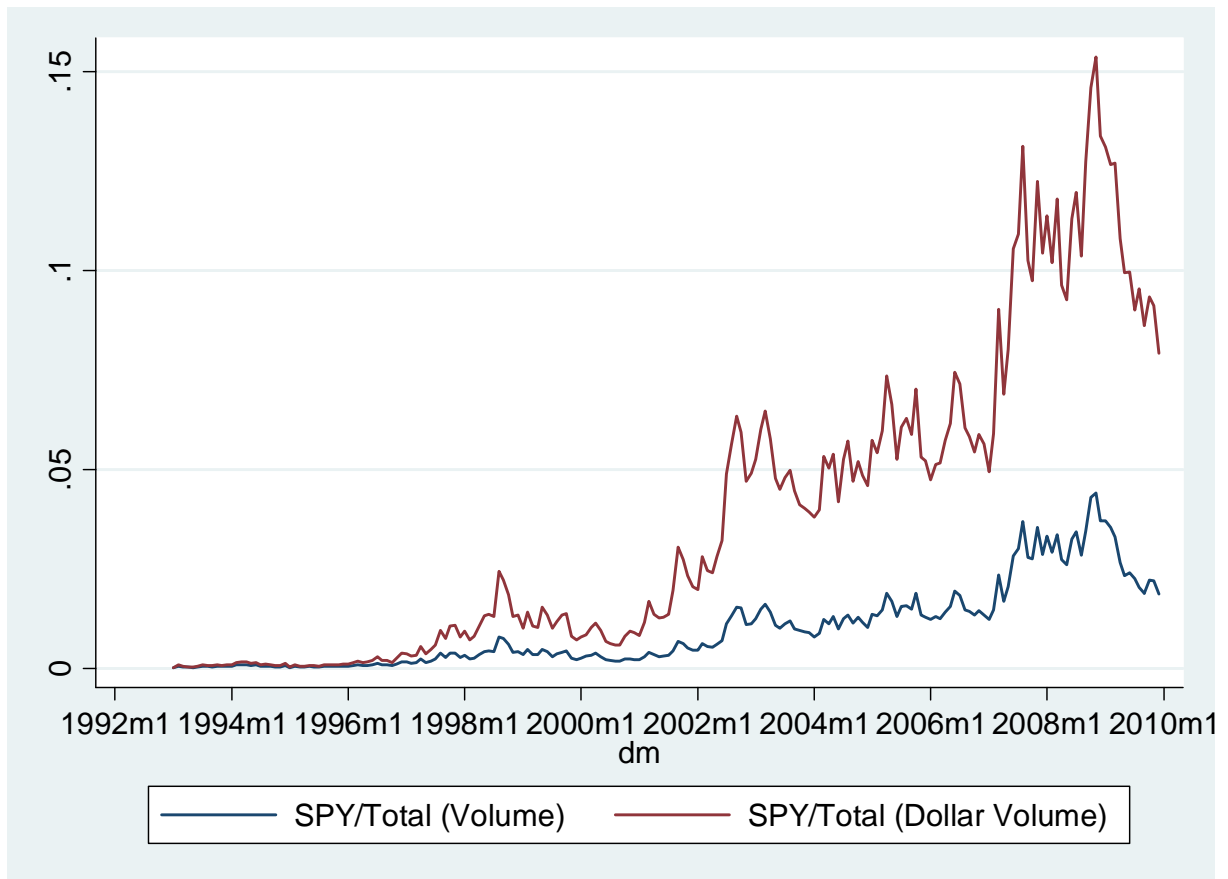
# Figures



**Figure 1.** This figure shows how cumulative stock returns covary with official and fundamental cumulative industry returns. A firm's fundamental industry return is defined as the equally weighted return of a portfolio consisting of all firms that are linked to firm  $j$  in the Hoberg and Phillips (2010a,b) sense in the previous calendar year. A firm's official industry is defined as the historical two digit SIC code reported by Compustat for the previous calendar year. We regress a firm's  $T$  week cumulative return on the firm's fundamental and the firm's official  $T$  week cumulative industry returns in a pooled cross section with double clustered standard errors (time and stock dimension). Formally, we estimate the following equation  $r_{j,t}^T = a_T + b_T \times r_{j,official;t}^T + c_T \times r_{j,fundamental;t}^T + e_{j,t}^T$  for varying  $T$ .  $r_{j,t}^T$  denotes firm  $j$ 's  $T$  week cumulative return.  $r_{j,fundamental;t}^T$  and  $r_{j,official;t}^T$  are the  $T$  week cumulative returns of baskets consisting of firms that are fundamentally and officially linked to firm  $j$ . We restrict the regressions to stocks for which the price at the beginning of the period exceeds \$5. In addition, we require the official and the fundamental industries to be populated by at least five different firms. We then plot the resulting coefficient estimates  $b_T$  and  $c_T$  against the time horizon  $T$  of the cumulative returns. The dotted lines show 95% confidence intervals for both coefficient estimates.



**Figure 2.** This diagram illustrates the idea behind the refined categorization bias strategy, which is based on a signal combining the official-fundamental industry return differential and a measure indicating whether a stock has moved closely in line with its official industry over the previous month.



**Figure 3.** This figure plots monthly (dollar) volume of the exchange traded fund SPDR S&P 500 (Ticker: SPY) normalized by total (dollar) volume of the CRSP universe.

# Tables

**Table I**  
**Summary Statistics: Returns**

This table reports summary statistics for the returns sample, which runs from 1997-2009. All variables are defined in the Appendix. All variables that rely on Compustat data are trimmed at the first and 99th percentile.

Panel A: Annual Variables						
	1					
	Mean	Median	SD	P25	P75	N
SIC2 firms	224.786	140.000	188.681	66.000	372.000	48641
Hoberg and Phillips firms	83.278	34.000	118.682	11.000	98.000	47935
$\ln(\text{BE}/\text{ME})$	-0.723	-0.644	0.854	-1.200	-0.166	47218
Analysts	6.416	4.000	7.879	0.000	10.000	48641
NYSE Breakpoint	3.442	2.000	2.879	1.000	5.000	48641
Panel B: Monthly Variables						
	1					
	Mean	Median	SD	P25	P75	N
$r_{j;t}$	0.012	0.000	0.201	-0.078	0.079	565266
$r_{j;\text{official};t}$	0.011	0.012	0.084	-0.035	0.053	565266
$r_{j;\text{fundamental};t}$	0.012	0.010	0.104	-0.039	0.057	564540
$\rho_{j;t}$	0.810	0.626	0.699	0.288	1.135	565261
$r_{j;\text{official};t} - r_{j;\text{fundamental};t}$	-0.001	0.001	0.068	-0.024	0.025	564540
$r_{j;t-2;t-12}$	0.089	0.004	0.597	-0.272	0.302	561824
$\ln(\text{M})_{j;t-1}$	12.535	12.472	2.156	10.942	13.989	565266
Panel C: Weekly Variables						
	Mean	Median	SD	P25	P75	N
$r_{j;t}$	0.003	-0.000	0.103	-0.038	0.036	2606105
$r_{j;\text{official};t}$	0.003	0.004	0.036	-0.014	0.021	2606105
$r_{j;\text{fundamental};t}$	0.003	0.003	0.043	-0.017	0.022	2606105
$\rho_{j;t}$	0.758	0.564	0.711	0.256	1.037	2379113
$r_{j;\text{official};t} - r_{j;\text{fundamental};t}$	0.000	0.000	0.024	-0.009	0.010	2606105

**Table II**  
**Summary Statistics: Analyst Forecasts**

This table reports summary statistics of all employed firm, analyst and analyst-firm level variables. The sample period runs from 1997-2009. All variables are defined in the Appendix. All variables are trimmed at the first and 99th percentile.

Panel A: Analyst-Firm Level						
	Mean	Median	SD	P25	P75	N
$ForecastError(3)_{i,j;t}$	0.002	0.000	0.021	-0.006	0.009	114874
$FP(3)_{i,j;t}$	0.052	0.054	0.037	0.036	0.072	114874
$\rho(3)_{i,j;t}$	8.129	4.611	9.518	1.820	10.478	114874
# Official $_{i,j;t}$	4.433	3.000	4.737	1.000	6.000	108129
# Fundamental $_{i,j;t}$	3.569	2.000	3.935	1.000	5.000	108129
Coverage $_{i,j;t}$	0.058	0.000	0.191	0.000	0.143	108129
Firm Experience $_{i,j;t}$	3.572	2.000	3.756	1.000	5.000	114874

Panel B: Firm Level						
	Mean	Median	SD	P25	P75	N
$ForecastError(0)_{j;t}$	-0.000	-0.000	0.005	-0.002	0.001	18751
$ForecastError(1)_{j;t}$	0.001	-0.000	0.009	-0.003	0.004	20487
$ForecastError(2)_{j;t}$	0.002	0.000	0.016	-0.004	0.008	20200
$ForecastError(3)_{j;t}$	0.004	0.001	0.022	-0.005	0.011	18550
$FP(0)_{j;t}$	0.050	0.053	0.044	0.033	0.071	18751
$FP(1)_{j;t}$	0.049	0.053	0.042	0.034	0.071	20487
$FP(2)_{j;t}$	0.053	0.056	0.043	0.036	0.074	20200
$FP(3)_{j;t}$	0.054	0.057	0.041	0.038	0.074	18550
$FP(0)_{j;official;t}$	0.048	0.048	0.025	0.030	0.064	18751
$FP(1)_{j;official;t}$	0.047	0.048	0.025	0.029	0.065	20487
$FP(2)_{j;official;t}$	0.050	0.051	0.026	0.032	0.066	20200
$FP(3)_{j;official;t}$	0.051	0.053	0.025	0.033	0.068	18550
$FP(0)_{j;fundamental;t}$	0.046	0.051	0.032	0.027	0.066	18751
$FP(1)_{j;fundamental;t}$	0.046	0.051	0.031	0.028	0.067	20487
$FP(2)_{j;fundamental;t}$	0.049	0.054	0.031	0.032	0.069	20200
$FP(3)_{j;fundamental;t}$	0.050	0.056	0.030	0.033	0.070	18550
$FP(0)_{j;official;t} - FP(0)_{j;fundamental;t}$	0.002	0.001	0.022	-0.006	0.012	18751
$FP(1)_{j;official;t} - FP(1)_{j;fundamental;t}$	0.001	0.000	0.021	-0.007	0.010	20487
$FP(2)_{j;official;t} - FP(2)_{j;fundamental;t}$	0.001	0.000	0.022	-0.007	0.010	20200
$FP(3)_{j;official;t} - FP(3)_{j;fundamental;t}$	0.001	-0.000	0.021	-0.007	0.009	18550
$ FP(3)_{j;official;t} - FP(3)_{j;fundamental;t} $	0.014	0.008	0.015	0.003	0.020	18550
$\rho(0)_{j;t}$	0.549	0.364	0.608	0.157	0.716	18751
$\rho(1)_{j;t}$	0.538	0.360	0.601	0.154	0.703	20487
$\rho(2)_{j;t}$	0.540	0.359	0.596	0.153	0.706	20200
$\rho(3)_{j;t}$	0.542	0.364	0.601	0.156	0.707	18550
ln(Size)	13.728	13.570	1.714	12.509	14.778	22839
VolRet	0.122	0.105	0.073	0.071	0.153	22839
Ret	0.015	0.013	0.041	-0.007	0.034	22839
ln(BM)	-0.764	-0.706	0.693	-1.185	-0.299	22839
VolRoe	0.072	0.046	0.075	0.024	0.089	22839
Roe	0.181	0.186	0.209	0.100	0.277	22839
SP500	0.201	0.000	0.401	0.000	0.000	22839
ln(P)	3.098	3.135	0.735	2.588	3.606	22839
Analysts	6.811	5.000	6.156	2.000	9.000	22839

Panel C: Analyst Level						
	Mean	Median	SD	P25	P75	N
Experience $_{i;t}$	5.986	4.000	5.317	2.000	9.000	27657
Industries $_{i;t}$	3.017	3.000	2.140	1.000	4.000	27657
Stocks $_{i;t}$	10.820	10.000	6.259	7.000	14.000	27657
Broker Size $_{i;t}$	53.297	38.000	46.796	16.000	81.000	27657

**Table III**  
**Weekly Categorization Bias Strategy (EW)**

This table shows alphas and factor loadings from weekly performance regressions of the Categorization Bias Strategy. At the beginning of each week, stocks are sorted into five quintiles according to the difference between the stock's official and fundamental industry return during week  $t - n$ , i.e.  $(r_{j:official;t-n} - r_{j:fundamental;t-n})$ . A firm's official industry return is defined as the equally weighted return of all its SIC2 peers. A firm's fundamental industry return is defined as the equally weighted return of a portfolio consisting of all its Hoberg and Phillips peers. We consider a maximum horizon of six weeks ( $n = 6$ ). Portfolios in the Q1 portfolio have experienced the most negative difference between official and fundamental industry returns. Q1-Q5 is a portfolio that is long in stocks for which the return difference between the official and the fundamental industry portfolios has been the most negative (Q1, bottom 20%) and short in stocks for which the return differential has been the most positive (Q5, top 20%). All portfolios are equally weighted and are rebalanced on a weekly basis.

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Signal (t-1)</b>							
Q1	0.281*** (7.18)	0.935*** (61.84)	0.790*** (29.37)	0.166*** (6.39)	-0.112*** (-7.13)	0.903	675
Q5	-0.0918** (-2.33)	1.037*** (68.16)	0.717*** (26.47)	0.0663** (2.54)	-0.0852*** (-5.41)	0.914	675
Q1-Q5	0.373*** (5.66)	-0.101*** (-3.98)	0.0739 (1.63)	0.0995** (2.28)	-0.0265 (-1.01)	0.0401	675
<b>Panel B: Signal (t-2)</b>							
Q1	0.173*** (4.30)	0.948*** (61.01)	0.741*** (26.82)	0.202*** (7.57)	-0.0981*** (-6.10)	0.898	674
Q5	0.00627 (0.15)	1.047*** (66.38)	0.741*** (26.41)	0.0591** (2.19)	-0.107*** (-6.55)	0.911	674
Q1-Q5	0.167** (2.45)	-0.0989*** (-3.76)	0.000349 (0.01)	0.143*** (3.16)	0.00885 (0.32)	0.0493	674
<b>Panel C: Signal (t-3)</b>							
Q1	0.210*** (5.32)	0.978*** (64.33)	0.789*** (29.16)	0.176*** (6.77)	-0.0704*** (-4.47)	0.907	673
Q5	-0.000617 (-0.01)	1.005*** (63.04)	0.689*** (24.29)	0.0335 (1.23)	-0.156*** (-9.45)	0.905	673
Q1-Q5	0.210*** (3.07)	-0.0278 (-1.05)	0.0993** (2.11)	0.143*** (3.15)	0.0858*** (3.13)	0.0347	673
<b>Panel D: Signal (t-4)</b>							
Q1	0.143*** (3.55)	0.984*** (63.35)	0.780*** (28.21)	0.199*** (7.47)	-0.0768*** (-4.77)	0.904	672
Q5	0.0457 (1.10)	1.004*** (62.75)	0.735*** (25.80)	0.0575** (2.10)	-0.157*** (-9.46)	0.906	672
Q1-Q5	0.0971 (1.40)	-0.0202 (-0.76)	0.0451 (0.95)	0.141*** (3.09)	0.0800*** (2.90)	0.0271	672
<b>Panel E: Signal (t-5)</b>							
Q1	0.112*** (3.10)	1.012*** (72.27)	0.899*** (36.08)	0.147*** (6.14)	-0.0281* (-1.94)	0.927	671
Q5	0.0600 (1.52)	0.971*** (63.55)	0.678*** (24.92)	0.0944*** (3.61)	-0.204*** (-12.90)	0.909	671
Q1-Q5	0.0522 (0.84)	0.0401* (1.67)	0.221*** (5.19)	0.0526 (1.29)	0.176*** (7.11)	0.119	671
<b>Panel F: Signal (t-6)</b>							
Q1	0.165*** (4.28)	1.029*** (69.03)	0.834*** (31.42)	0.101*** (3.98)	-0.00889 (-0.58)	0.918	670
Q5	0.0293 (0.70)	0.960*** (59.20)	0.669*** (23.16)	0.112*** (4.02)	-0.208*** (-12.41)	0.897	670
Q1-Q5	0.136** (2.03)	0.0687*** (2.66)	0.165*** (3.58)	-0.0102 (-0.23)	0.199*** (7.45)	0.119	670

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table IV

**Monthly Categorization Bias Strategy (VW): By Market Capitalization**

At the beginning of each month, stocks are ranked according to their market capitalization and separated into brackets of small, mid and large capitalization stocks (three terciles). Within each tercile, stocks are then sorted into five quintiles according to the difference between an equally weighted portfolio of stocks belonging to the firm's official and an equally weighted portfolio consisting of stocks belonging to a firm's fundamental industry ( $r_{j;official;t-1} - r_{j;fundamental;t-1}$ ). Within each size tercile, the ranked stocks are assigned to one of five quintile portfolios. Q1 contains firms for which the return differential between the official and the fundamental industry is the lowest. All stocks are value weighted within a quintile portfolio and rebalanced every month. Monthly alphas as well as Fama and French (1993) and Carhart (1997) factor loadings are in percent. Panels A, B and C show the results for value weighted portfolios consisting of small, mid and large stocks respectively.

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Small Cap Stocks (First Tercile)</b>							
Q1	0.915*** (3.69)	0.808*** (14.41)	1.157*** (17.32)	0.211*** (2.91)	-0.116*** (-2.73)	0.835	155
Q2	0.0683 (0.36)	0.839*** (19.28)	0.894*** (17.23)	0.229*** (4.06)	-0.0994*** (-3.02)	0.871	155
Q3	-0.0535 (-0.24)	0.827*** (16.32)	0.815*** (13.50)	0.416*** (6.34)	-0.135*** (-3.53)	0.817	155
Q4	-0.420* (-1.92)	0.979*** (19.82)	0.820*** (13.93)	0.416*** (6.50)	-0.112*** (-2.99)	0.852	155
Q5	-0.629*** (-2.58)	1.025*** (18.57)	0.903*** (13.73)	0.328*** (4.59)	-0.144*** (-3.46)	0.844	155
Q1-Q5	1.544*** (4.25)	-0.217*** (-2.64)	0.254*** (2.60)	-0.117 (-1.10)	0.0284 (0.46)	0.103	155
<b>Panel B: Mid Cap Stocks (Second Tercile)</b>							
Q1	0.652** (2.21)	0.955*** (14.30)	1.133*** (14.23)	0.0854 (0.99)	-0.0520 (-1.03)	0.807	155
Q2	0.0130 (0.08)	0.937*** (24.33)	0.720*** (15.67)	0.326*** (6.52)	-0.0270 (-0.93)	0.887	155
Q3	-0.241 (-1.15)	0.941*** (19.80)	0.603*** (10.65)	0.508*** (8.25)	-0.0293 (-0.82)	0.822	155
Q4	-0.350* (-1.84)	1.014*** (23.56)	0.702*** (13.69)	0.312*** (5.59)	-0.0692** (-2.13)	0.878	155
Q5	-0.536** (-1.98)	1.202*** (19.66)	0.749*** (10.28)	0.225*** (2.84)	-0.0760* (-1.65)	0.831	155
Q1-Q5	1.187** (2.39)	-0.246** (-2.19)	0.384*** (2.87)	-0.139 (-0.95)	0.0240 (0.28)	0.0946	155
<b>Panel C: Large Cap Stocks (Third Tercile)</b>							
Q1	0.255 (1.17)	0.903*** (18.28)	0.0334 (0.57)	-0.153** (-2.39)	-0.0190 (-0.51)	0.768	155
Q2	0.144 (1.03)	0.872*** (27.55)	-0.164*** (-4.34)	0.0519 (1.26)	0.0305 (1.28)	0.858	155
Q3	-0.103 (-0.58)	0.961*** (24.03)	-0.0923* (-1.94)	0.277*** (5.35)	0.0296 (0.98)	0.818	155
Q4	-0.216 (-1.45)	0.968*** (28.58)	-0.179*** (-4.44)	0.0236 (0.54)	0.0220 (0.86)	0.870	155
Q5	0.00746 (0.03)	1.060*** (21.41)	-0.124** (-2.10)	-0.164** (-2.56)	-0.0573 (-1.53)	0.814	155
Q1-Q5	0.247 (0.66)	-0.157* (-1.86)	0.158 (1.57)	0.0112 (0.10)	0.0383 (0.60)	0.0467	155

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$



Table V

**Monthly Categorization Bias Strategy (EW): By Industry Follower**

At the beginning of each month, stocks are ranked according to whether their returns tracked the returns of their official industries over the previous month. To do so, we calculate the absolute value of the difference between the firm's return and the firm's official industry return normalized by the rolling standard deviation of the difference. We denote this measure as  $\rho_{j;t-1} = |r_{j;t-1} - r_{j;official;t-1}| / \sigma_{j;24}(r_{j;t-1} - r_{j;official;t-1})$ . Low values of  $\rho_{j;t-1}$  indicate that over the prior month, the return of firm  $j$  has been abnormally close to the return of the official industry, while high values indicate divergence. Stocks in the first tercile of the distribution of  $\rho$  are labeled as *Strong Official Industry Follower* stocks. In each *Official Industry Follower* tercile, we then sort stocks according to the difference between an equally weighted portfolio of stocks belonging to the firm's official and an equally weighted portfolio consisting of stocks that belong to the firm's fundamental industry ( $r_{j;official;t-1} - r_{j;fundamental;t-1}$ ). Q1 contains firms for which the return differential between the official and the fundamental industry is the lowest. Q1-Q5 is a portfolio that is long in stocks for which the return difference between the official and the fundamental industry portfolios has been the most negative (Q1, bottom 20%) and short in stocks for which the return differential has been the most positive (Q5, top 20%). All stocks are equally weighted within a quintile portfolio and rebalanced every month. Monthly alphas as well as Fama and French (1993) and Carhart (1997) factor loadings are in percent. Panels A, B and C show the results for quintile portfolios restricted to *Strong*, *Medium* and *Weak Official Industry Follower* stocks (three terciles of  $\rho$ ).

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Strong Official Industry Follower (First Tercile)</b>							
Q1	0.957*** (4.18)	0.891*** (17.20)	0.902*** (14.61)	0.0857 (1.28)	-0.0776** (-1.99)	0.842	155
Q2	0.0738 (0.48)	0.904*** (25.72)	0.564*** (13.46)	0.253*** (5.55)	-0.0464* (-1.75)	0.890	155
Q3	-0.401** (-2.54)	0.913*** (25.60)	0.495*** (11.65)	0.381*** (8.24)	-0.0227 (-0.84)	0.878	155
Q4	-0.424** (-2.32)	0.987*** (23.80)	0.574*** (11.61)	0.249*** (4.63)	-0.0364 (-1.16)	0.870	155
Q5	-0.617*** (-2.81)	1.060*** (21.28)	0.649*** (10.94)	0.0382 (0.59)	-0.0967** (-2.57)	0.860	155
Q1-Q5	1.575*** (4.10)	-0.169* (-1.95)	0.253** (2.44)	0.0475 (0.42)	0.0191 (0.29)	0.0615	155
<b>Panel B: Medium Official Industry Follower (Second Tercile)</b>							
Q1	0.540** (2.57)	0.921*** (19.37)	0.835*** (14.74)	0.112* (1.82)	-0.0752** (-2.09)	0.860	155
Q2	0.320** (2.13)	0.878*** (25.82)	0.545*** (13.45)	0.221*** (5.02)	-0.0459* (-1.79)	0.892	155
Q3	0.0326 (0.18)	0.872*** (21.69)	0.468*** (9.77)	0.392*** (7.52)	-0.0392 (-1.29)	0.839	155
Q4	-0.452*** (-2.64)	0.972*** (25.10)	0.548*** (11.86)	0.307*** (6.11)	-0.0437 (-1.50)	0.879	155
Q5	-0.423* (-1.81)	1.096*** (20.69)	0.642*** (10.17)	0.180*** (2.62)	-0.105*** (-2.64)	0.845	155
Q1-Q5	0.963** (2.51)	-0.175** (-2.02)	0.193* (1.86)	-0.0676 (-0.60)	0.0303 (0.46)	0.0586	155
<b>Panel C: Weak Official Industry Follower (Third Tercile)</b>							
Q1	0.469 (1.56)	0.820*** (12.04)	1.041*** (12.83)	0.0645 (0.73)	-0.0218 (-0.42)	0.760	155
Q2	0.0326 (0.21)	0.889*** (25.61)	0.602*** (14.55)	0.268*** (5.97)	-0.0475* (-1.81)	0.893	155
Q3	-0.217 (-1.28)	0.873*** (22.67)	0.494*** (10.77)	0.521*** (10.46)	-0.0679** (-2.34)	0.857	155
Q4	-0.162 (-1.07)	0.973*** (28.31)	0.516*** (12.62)	0.412*** (9.25)	-0.0919*** (-3.54)	0.903	155
Q5	-0.300 (-1.41)	1.070*** (22.27)	0.562*** (9.82)	0.323*** (5.19)	-0.130*** (-3.58)	0.857	155
Q1-Q5	0.769* (1.73)	-0.250** (-2.49)	0.479*** (3.99)	-0.259** (-1.98)	0.108 (1.42)	0.204	155

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table VI  
**Monthly Categorization Bias Strategy (VW): Strong Industry Follower; By Market Capitalization**

In this table, we restrict the portfolios to *Strong Official Industry Follower Stocks* only (low values of  $\rho$ ). See the previous table for a formal definition of  $\rho$ . In each month, we first sort the set of *Strong Official Industry Follower Stocks* into three size brackets (small, mid and large caps). Within each size tercile, we then sort the stocks according to the difference between an equally weighted portfolio of stocks belonging to the stock's official and an equally weighted portfolio consisting of stocks that belong to the stock's fundamental industry ( $r_{j:official;t-1} - r_{j:fundamental;t-1}$ ). The ranked stocks are assigned to one of five quintile portfolios. Q1 contains firms for which the return differential between the official and the fundamental industry is the lowest. Q1-Q5 is a portfolio that is long in stocks for which the return difference between the official and the fundamental industry portfolios has been the most negative (Q1, bottom 20%) and short in stocks for which the return differential has been the most positive (Q5, top 20%). All stocks are value weighted within a quintile portfolio and rebalanced every month. Monthly alphas as well as Fama and French (1993) and Carhart (1997) factor loadings are in percent. Panels A, B and C show alphas and loadings for value weighted portfolios consisting of small, mid and large stocks respectively.

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Small Cap Stocks (First Tercile)</b>							
Q1	1.338*** (4.15)	0.834*** (11.43)	1.214*** (13.96)	0.240** (2.54)	-0.103* (-1.88)	0.762	155
Q2	-0.0419 (-0.15)	0.856*** (13.37)	0.778*** (10.20)	0.175** (2.11)	-0.0764 (-1.58)	0.741	155
Q3	-0.407 (-1.47)	0.837*** (13.35)	0.803*** (10.75)	0.337*** (4.15)	-0.139*** (-2.93)	0.749	155
Q4	-0.544** (-1.98)	0.993*** (15.96)	0.816*** (11.00)	0.322*** (4.00)	-0.0724 (-1.54)	0.785	155
Q5	-0.731** (-2.33)	0.948*** (13.36)	0.983*** (11.63)	0.196** (2.14)	-0.192*** (-3.59)	0.772	155
Q1-Q5	2.069*** (4.72)	-0.113 (-1.14)	0.231* (1.95)	0.0434 (0.34)	0.0888 (1.18)	0.0504	155
<b>Panel B: Mid Cap Stocks (Second Tercile)</b>							
Q1	0.760** (2.37)	0.947*** (13.07)	1.101*** (12.75)	0.110 (1.17)	-0.0637 (-1.16)	0.774	155
Q2	-0.235 (-1.15)	0.936*** (20.20)	0.771*** (13.97)	0.385*** (6.42)	-0.0346 (-0.99)	0.849	155
Q3	-0.189 (-0.78)	0.968*** (17.73)	0.621*** (9.55)	0.398*** (5.63)	-0.0438 (-1.06)	0.791	155
Q4	-0.393 (-1.54)	0.984*** (17.03)	0.726*** (10.55)	0.258*** (3.45)	-0.0294 (-0.67)	0.794	155
Q5	-0.849*** (-2.82)	1.209*** (17.72)	0.749*** (9.21)	0.0616 (0.70)	-0.0577 (-1.12)	0.805	155
Q1-Q5	1.609*** (3.05)	-0.262** (-2.20)	0.352** (2.48)	0.0480 (0.31)	-0.00605 (-0.07)	0.0638	155
<b>Panel C: Large Cap Stocks (Third Tercile)</b>							
Q1	0.724** (2.36)	0.961*** (13.83)	0.191** (2.31)	-0.0403 (-0.45)	0.00996 (0.19)	0.653	155
Q2	0.206 (0.82)	0.955*** (16.85)	-0.185*** (-2.74)	0.124* (1.69)	0.0897** (2.10)	0.679	155
Q3	-0.606*** (-2.78)	0.935*** (18.93)	-0.148** (-2.51)	0.275*** (4.30)	0.0224 (0.60)	0.737	155
Q4	-0.0231 (-0.11)	1.018*** (20.74)	-0.0758 (-1.30)	-0.0679 (-1.07)	0.0447 (1.20)	0.786	155
Q5	0.0884 (0.31)	1.036*** (15.86)	-0.0454 (-0.58)	-0.225*** (-2.65)	0.00817 (0.17)	0.707	155
Q1-Q5	0.636 (1.43)	-0.0759 (-0.76)	0.237** (1.98)	0.184 (1.42)	0.00179 (0.02)	0.0340	155

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table VII

Pooled Cross-Sectional Regressions (Monthly Returns)

This table shows results from pooled cross sectional regressions with double clustered standard errors (month-firm). The dependent variable is the monthly stock return of firm  $j$ . We restrict the regressions to stock-month observations for which the stock price at the end of the preceding month is greater or equal to \$5. All specifications include month dummies. Weak, Med and Strong Ind Follower $_{j,t-1}$  indicate terciles of the Official Industry Follower measure  $\rho$ . Equations (3), (4), (5) and (6) are restricted to Strong Official Industry Follower stocks. Small, Mid and Large Cap $_{j,t-1}$  are dummy variables indicating size terciles based on market capitalization in each month. Med and High Analyst $_{j,t-1}$  indicate whether the stock receives medium or high analyst coverage. After 2002 is a dummy equal to one for firm-month observations after 2002. All other variables are defined in the appendix.

	(1)	(2)	(3)	(4)	(5)	(6)
--	-----	-----	-----	-----	-----	-----

$$r_{j,official;t-1} - r_{j,fundamental;t}$$

**Table VIII**  
**Monthly Pseudo Industry Classification Strategy (EW)**

In this table, we construct a portfolio strategy based on a Pseudo Hoberg and Phillips Industry classification (see section III.3 for details). We first sort firms into terciles of the previously defined Industry Follower Measure  $\rho_{j;t-1}$  (see the appendix for a formal definition of  $\rho_{j;t-1}$ ). Then, within each tercile of  $\rho - j, t - 1$ , we sort firms according to  $r_{j;official;t-1} - r_{j;pseudo;t-1}$ .  $r_{j;pseudo;t-1}$  is the return on an equally weighted portfolio consisting of random firms. For firm  $j$ , the number of random firms contained in the pseudo portfolio is equal to the number of firm  $j$ 's Hoberg and Phillips links in that month. As in preceding portfolio sorts,  $r_{j;official;t-1}$  is an equally weighted portfolio made up of all firms belonging to firm  $j$ 's two digit SIC code. Q1 contains portfolios for which the firm's official-pseudo industry return differential is most negative. Panels A, B and C show the results for quintile portfolios restricted to *Strong, Medium and Weak Official Industry Follower* stocks (three terciles of  $\rho$ ).

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Strong Official Industry Follower (First Tercile)</b>							
Q1	-0.687*** (-2.68)	1.041*** (17.92)	0.511*** (7.38)	0.253*** (3.37)	-0.127*** (-2.89)	0.793	155
Q2	-0.397* (-1.93)	1.053*** (22.67)	0.523*** (9.45)	0.214*** (3.56)	-0.0603* (-1.72)	0.855	155
Q3	-0.0149 (-0.10)	0.920*** (27.62)	0.682*** (17.18)	0.129*** (2.99)	-0.0510** (-2.03)	0.915	155
Q4	0.201 (1.17)	0.873*** (22.42)	0.707*** (15.24)	0.172*** (3.40)	-0.0630** (-2.15)	0.881	155
Q5	0.483 (1.64)	0.852*** (12.73)	0.889*** (11.15)	0.118 (1.36)	0.0288 (0.57)	0.738	155
Q1-Q5	-1.170** (-2.35)	0.189* (1.68)	-0.378*** (-2.81)	0.135 (0.93)	-0.155* (-1.82)	0.125	155
<b>Panel B: Medium Official Industry Follower (Second Tercile)</b>							
Q1	-0.511** (-2.33)	1.046*** (21.07)	0.484*** (8.18)	0.356*** (5.53)	-0.130*** (-3.47)	0.837	155
Q2	-0.0863 (-0.41)	1.023*** (21.66)	0.568*** (10.10)	0.188*** (3.07)	-0.0799** (-2.24)	0.851	155
Q3	0.0743 (0.52)	0.931*** (28.50)	0.621*** (15.95)	0.230*** (5.42)	-0.0378 (-1.54)	0.912	155
Q4	0.264* (1.72)	0.864*** (24.77)	0.647*** (15.58)	0.189*** (4.18)	-0.0597** (-2.27)	0.895	155
Q5	0.318 (1.19)	0.854*** (14.10)	0.790*** (10.95)	0.222*** (2.82)	0.000292 (0.01)	0.752	155
Q1-Q5	-0.829* (-1.93)	0.192** (1.97)	-0.306*** (-2.64)	0.134 (1.06)	-0.130* (-1.78)	0.127	155
<b>Panel C: Weak Official Industry Follower (Third Tercile)</b>							
Q1	-0.622*** (-3.36)	1.022*** (24.41)	0.528*** (10.59)	0.459*** (8.46)	-0.112*** (-3.56)	0.874	155
Q2	-0.201 (-1.10)	0.964*** (23.31)	0.557*** (11.29)	0.268*** (5.00)	-0.0999*** (-3.20)	0.870	155
Q3	0.150 (1.08)	0.925*** (29.32)	0.597*** (15.89)	0.309*** (7.55)	-0.0843*** (-3.54)	0.916	155
Q4	0.248 (1.47)	0.829*** (21.76)	0.636*** (13.99)	0.306*** (6.20)	-0.0833*** (-2.90)	0.868	155
Q5	0.238 (1.00)	0.911*** (16.93)	0.724*** (11.29)	0.374*** (5.37)	0.00643 (0.16)	0.790	155
Q1-Q5	-0.860** (-2.35)	0.110 (1.33)	-0.196** (-1.99)	0.0845 (0.79)	-0.119* (-1.90)	0.0913	155

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table IX

**Firm-Level Forecast Error**

The dependent variable in columns (1)–(4) is the average firm-level forecast error associated with forecasts issued in the first fiscal quarter of firm  $j$  (i.e. forecasts with a horizon of three quarters). In calculating the average firm level forecast error, we use the last forecast in quarter  $t$  that analyst  $i$  issues or revises for firm  $j$ 's year end earnings. The dependent variables in columns (5), (6) and (7) are the average firm-level forecast errors associated with the last forecast issued or revised in the second, third and fourth quarter of the fiscal year respectively. Thus, the forecast horizons in these regressions are two, one or zero quarters.  $FP(T)_{j,official;t}$  is the average forecast to price ratio prevailing in firm  $j$ 's two digit SIC industry in quarter  $t$ .  $FP(T)_{j,fundamental;t}$  is the average forecast to price ratio prevailing among firm  $j$ 's Hoberg and Phillips peers.  $\rho(T)_{j,t}$  measures whether firm  $j$ 's average forecast to price ratio in quarter  $t$  is close to that of its official industry and is formally defined as  $\rho(T)_{j,t} = |FP(T)_{j,t} - FP(T)_{j,official;t}| / \sigma_t(FP(T)_{j,t} - FP(T)_{j,official;t})$  where  $FP(T)_{j,t}$  is the average forecast to price ratio for firm  $j$  in year  $t$  and  $\sigma_t$  denotes the cross-sectional standard deviation in year  $t$ . All other variables are defined in the appendix. All regressions include year fixed effects. We restrict all regressions to forecasts for which the stock price of the firm at the beginning of the quarter exceeds \$5 and for which the industry portfolios are populated by at least five different firms. Standard errors are clustered at the firm level.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
$FP(T)_{j,official;t} - FP(T)_{j,fundamental;t}$	0.0034 (0.43)	-0.0137 (-1.62)	0.0437*** (3.70)	0.0292** (2.26)	0.0013 (0.15)	-0.0070 (-1.39)	-0.0011 (-0.42)
$\rho(T)_{j,t}$			0.0026*** (7.28)	0.0008* (1.79)	0.0004 (1.37)	0.0003** (1.98)	0.0004*** (4.09)
$\rho(T)_{j,t} \times (FP(T)_{j,official;t} - FP(T)_{j,fundamental;t})$			-0.0489*** (-4.13)	-0.0493*** (-3.48)	-0.0226** (-2.26)	-0.0042 (-0.72)	-0.0013 (-0.47)
ln(Size)		-0.0000 (-0.02)		0.0000 (0.03)	0.0004*** (2.79)	0.0004*** (4.27)	0.0002*** (3.55)
VolRet		0.0246*** (6.46)		0.0232*** (6.08)	0.0111*** (3.98)	0.0009 (0.55)	-0.0021** (-2.49)
Ret		-0.0573*** (-9.42)		-0.0567*** (-9.24)	-0.0235*** (-5.54)	-0.0088*** (-3.56)	-0.0030** (-2.35)
ln(BM)		0.0003 (0.73)		0.0001 (0.23)	0.0001 (0.40)	0.0002 (1.46)	-0.0003*** (-3.50)
VolRoe		-0.0016 (-0.52)		-0.0024 (-0.76)	-0.0000 (-0.00)	-0.0021* (-1.79)	-0.0030*** (-4.63)
Roe		-0.0012 (-0.96)		-0.0019 (-1.45)	-0.0008 (-0.90)	0.0001 (0.16)	0.0003 (1.07)
SP500		0.0003 (0.61)		0.0003 (0.62)	0.0001 (0.28)	-0.0002 (-0.81)	-0.0002* (-1.90)
ln(P)		-0.0025*** (-5.97)		-0.0025*** (-5.89)	-0.0023*** (-7.29)	-0.0013*** (-6.17)	-0.0005*** (-5.46)
Analysts		-0.0002*** (-5.89)		-0.0002*** (-5.89)	-0.0002*** (-7.51)	-0.0001*** (-6.43)	-0.0000*** (-3.88)
Constant	0.0042*** (23.92)	0.0119*** (4.64)	0.0027*** (12.11)	0.0114*** (4.47)	0.0043** (2.32)	0.0005 (0.41)	-0.0007 (-1.11)
Observations	22737	18550	22737	18550	20200	20487	18751
$R^2$	0.026	0.059	0.032	0.060	0.053	0.041	0.022

t statistics in parentheses. \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table X

## Analyst-Firm Level

The dependent variable is  $\rho(3)_{i,j,t} = |FP(3)_{i,j,t} - FP(3)_{j,official,t}| / \sigma_{j,t}(FP(3)_{i,j,t} - FP(3)_{j,official,t})$ , which measures whether an analyst has the tendency to more strongly track the firm's official industry level forecast than all other analysts following the same firm  $j$ . Experience is the number of years the analyst has been issuing forecasts. Firm Experience is the number of years analyst  $i$  has been forecasting earnings for firm  $j$ . Stocks is the number of stocks for which analyst  $i$  has issued at least one forecast in year  $t$ . Industries is the number of different SIC2 industries covered by analyst  $i$  in year  $t$ .  $Coverage_{i,j,t}$  is defined as  $(\# Official_{i,j,t} - \# Fundamental_{i,j,t}) / Stocks_{i,j,t}$  where  $\# Official_{i,j,t}$  is the number of firms covered by analyst  $i$  in year  $t$  which belong to firm  $j$ 's official industry.  $\# Fundamental_{i,j,t}$  is the number of Hoberg and Phillips peers of firm  $j$  covered by analyst  $i$  in year  $t$ . Broker Size is the number of analysts employed at analyst  $i$ 's brokerage house in year  $t$ . All regressions include the previously used firm-level control variables and year fixed effects. We restrict all regressions to forecasts for which the stock price of the firm at the beginning of the quarter exceeds \$5 and for which the industry portfolios are populated by at least five different firms. Standard errors are clustered at the analyst level.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t} $	51.5470*** (14.01)	43.0378*** (7.95)	44.8418*** (10.34)	39.8575*** (6.35)	-9.4689 (-1.43)	54.2078*** (13.97)	53.3291*** (9.84)
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Experience_{i,t}$		1.1799* (1.95)					
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Firm\ Experience_{i,j,t}$			1.8540** (2.39)				
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Stocks_{i,j,t}$				0.8468** (2.18)			
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Industries_{i,t}$					15.3842*** (10.35)		
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Coverage_{i,j,t}$						-42.3956*** (-3.12)	
$ FP(T)_{j,official,t} - FP(T)_{j,fundamental,t}  \times Broker\ Size_{i,t}$							-0.0384 (-0.53)
Coverage <sub>i,j,t</sub>						0.5621** (2.05)	
Experience <sub>i,t</sub>		-0.0043 (-0.19)	0.0116 (0.57)	0.0114 (0.56)	0.0060 (0.30)	0.0093 (0.44)	0.0100 (0.49)
Firm Experience <sub>i,j,t</sub>		-0.0789*** (-5.60)	-0.1041*** (-6.27)	-0.0792*** (-5.64)	-0.0779*** (-5.54)	-0.0772*** (-5.30)	-0.0784*** (-5.55)
Stocks <sub>i,t</sub>		0.0089 (1.01)	0.0087 (0.98)	-0.0000 (-0.00)	0.0083 (0.92)	0.0087 (0.93)	0.0082 (0.91)
Industries <sub>i,t</sub>		-0.0128 (-0.33)	-0.0112 (-0.29)	-0.0171 (-0.44)	-0.2171*** (-4.75)	-0.0069 (-0.16)	0.0032 (0.08)
Broker Size <sub>i,t</sub>		-0.0027 (-1.47)	-0.0026 (-1.46)	-0.0026 (-1.45)	-0.0025 (-1.41)	-0.0028 (-1.50)	-0.0021 (-1.01)
Observations	114874	114874	114874	114874	113914	107183	113914
$R^2$	0.255	0.256	0.256	0.256	0.259	0.259	0.257
Analyst Fixed Effects	No	Yes	Yes	Yes	Yes	Yes	Yes

† statistics in parentheses. \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

# Appendix

## A Definition of return related variables

$SIC2\ firms$  is the number of different firms belonging to firm  $j$ 's two digit SIC industry in year  $t - 1$ .

$Hoberg\ and\ Phillips\ firms$  is the number of different firms that are linked to firm  $j$  in the Hoberg and Phillips (2010a,b) sense in year  $t-1$ .

$\ln(BE/ME)$  is the log of the firm's book to market equity. Book equity is defined as  $(at - lt - pstkl + txditc + dcvt)$ , market equity is calculated as  $prc\_f \times csho$ . In line with the asset pricing literature, the  $BE/ME$  ratio is updated at the end of each June in order to ensure a gap of at least 6 month between fiscal year ends and stock returns.

$Analysts_{j,t}$  is the number of different analysts having issued at least one fiscal year end EPS forecast in a given year. If there is no match between the CRSP/Compustat and the I/B/E/S datasets, it is assumed that the corresponding firm does not have any analyst coverage.

$r_{j,t}$  is the firm's monthly or weekly return from CRSP.

$r_{j,official,t}$  is the monthly or weekly return to an equally weighted portfolio consisting of all firms with the same SIC2 code.

$r_{j,fundamental,t}$  is the monthly or weekly return to an equally weighted portfolio consisting of all firms that are linked to firm  $j$  in the Hoberg and Phillips (2010a,b) sense. This variable varies at the firm-month level.

$\rho_{j,t} = |(r_{j,t} - r_{j,official,t})/\sigma_{j,24}(r_{j,t} - r_{j,official,t})|$  is a measure that indicates to what extent the return of firm  $j$  has followed the return of its official SIC2 industry.  $\sigma_{j,24}(r_{j,t} - r_{j,official,t})$  denotes the standard deviation of the return differential calculated for rolling windows of 24 months.

*Low, Medium and High Official Industry Comovement* indicate the first, second and third tercile of  $\rho_{j,t-1}$  in each month.

$r_{j,t-2,t-12}$  is the firm's cumulative return between month  $t - 12$  and  $t - 2$ .

$\ln(M)_{j,t-1}$  is the natural log of the firm's market capitalization at the end of the previous month.

$Coverage_{i,j,t}$  is defined as  $(\# Official_{i,j,t} - \# Fundamental_{i,j,t})/Stocks_{i,j,t}$ .

Firm Experience $_{i,j,t}$  is the number of years analyst  $i$  has been forecasting earnings for firm  $j$  in year  $t$ .



## B Definition of analyst related variables

### B.1 Firm Level Variables

#### Main Variables

$ForecastError(T)_{j,t}$  is the average firm-level forecast error for all  $T$  horizon forecasts that were issued for firm  $j$ 's fiscal year end earnings in quarter  $t$ .

$FP(T)_{j,t}$  is the average  $T$  horizon forecast to price ratio at the firm level which is calculated by averaging across all outstanding  $T$  horizon forecasts for firm  $j$  in quarter  $t$ .

$FP(T)_{j,official,t}$  is the average  $T$  horizon forecast to price ratio prevailing in firm  $j$ 's official industry in quarter  $t$ . The official industry is defined as the firm's two digit SIC industry. It is calculated using all outstanding forecasts for firms in the official sector. We require the official sector to be populated by at least five different firms.

$FP(T)_{j,fundamental,t}$  is the average forecast to price ratio prevailing among firm  $j$ 's Hoberg and Phillips (2010a,b) siblings in the first month of the quarter in which a forecast was issued. It is calculated using all outstanding forecasts for firms in the official sector. We require the firm to be linked to at least five different firms in the Hoberg and Phillips (2010a,b) sense.

$\rho(T)_{j,t}$  is  $|FP(T)_{j,t} - FP(T)_{j,official,t}| / \sigma_t(FP(T)_{j,t} - FP(T)_{j,official,t})$ .  $\sigma_t$  is the cross sectional standard deviation in period  $t$ . The measure gives an indication about whether analysts as a whole tend to base their earnings forecasts for firm  $j$  on the average earnings forecast prevailing in firm  $j$ 's official industry category.

## Control Variables

In all regressions of the second part of the paper we control for the following firm level characteristics.

*Analysts* is the number of analysts with at least one year end EPS forecast throughout the current year.

*Ret* is the average monthly return over the previous calendar year.

*VolRet* is the standard deviation of monthly stock returns over the previous calendar year.

$\ln(Size)$  is the log of market capitalization at the end of the previous calendar year.

*Roe* is the return on equity in previous year  $(oibdp - dp) / (at - lt - pstkl + txditc + dcvt)$ .

*Volroe* is the firm level standard deviation of return on equity between year  $t - 1$  and  $t - 5$ .

$\ln(BM)$  is the log of the firm's book to market equity at the end of the previous calendar year. Book equity is defined as  $(at - lt - pstkl + txditc + dcvt)$ , market equity is calculated as  $prc\_c \times csho$ .

*SP500* is a dummy indicating whether the company belongs to the SP500 index in the previous calendar year.

$\ln(P)$  is the log of the firm's stock price at the beginning of quarter  $t$ .

## B.2 Analyst and Analyst-Firm Level Variables

### Analyst-Firm Level

The following variables are calculated at the forecast level.

$ForecastError(T)_{i,j,t}$  is the forecast error associated with the last  $T$  horizon forecast analyst  $i$  issues or revises in quarter  $t$  for year end earnings of firm  $j$ . The horizon  $T$  is measured in quarters.

$FP(3)_{i,j,t}$  is the forecast to price ratio in quarter  $t$  for analyst  $i$ 's three quarter ahead forecast of firm  $j$ 's fiscal year end earnings. It is calculated by normalizing the actual three quarter ahead forecast  $F(3)_{i,j,t}$  by the stock price prevailing for firm  $j$  at the beginning of quarter  $t$ .

$\rho(3)_{i,j,t}$  is defined as  $(|FP(3)_{i,j,t} - FP(3)_{j,official,t}|) / (\sigma_{jt}(FP(3)_{i,j,t} - FP(3)_{j,official,t}))$

$\# Official_{i,j,t}$  is the number of SIC2 peers of company  $j$  covered by analyst  $i$  in year  $t$ .

$\# Official_{i,j,t}$  is the number of Hoberg and Phillips (2010a,b) peers of company  $j$  covered by analyst  $i$  in year  $t$ .

### Analyst Level

The following variables are calculated at the analyst level.

$Stocks_{i,t}$  is the number of different stocks covered by analyst  $i$  in year  $t$ .

$Industries_{i,t}$  is the number of different industries (at the SIC2 level) covered by ana-

lyst  $i$  in year  $t$ .

$Experience_{i,t}$  is the number of years analyst  $i$  appears in the I/B/E/S database.

$BrokerSize_{i,t}$  is the number of analysts with at least one EPS forecast working at the same broker house as analyst  $i$  in year  $t$ .

Table A.I  
Monthly Categorization Bias Strategy (EW)

At the beginning of each month, stocks are ranked according to the difference between an equally weighted portfolio of stocks belonging to the firm's official and an equally weighted portfolio consisting of stocks that belong to the firm's fundamental industry ( $r_{j,official;t-1} - r_{j,fundamental;t-1}$ ). A firm's official industry is defined as the two digit SIC industry to which the firm belongs. A firm's fundamental industry is composed of all companies to which the firm has economic links in the Hoberg and Phillips (2010a,b) sense. The ranked stocks are assigned to one of five quintile portfolios. Q1 contains firms for which the return differential between the official and the fundamental industry is the lowest. Q1-Q5 is a portfolio that is long in stocks for which the return difference between the official and the fundamental industry portfolios has been the most negative (Q1, bottom 20%) and short in stocks for which the return differential has been the most positive (Q5, top 20%). All stocks are equally weighted within a quintile portfolio and re-balanced every month. Monthly alphas as well as Fama and French (1993) and Carhart (1997) factor loadings are in percent.

	alpha	mktrf	smb	hml	mom	R2	N
Q1	0.646*** (2.85)	0.871*** (16.98)	0.929*** (15.19)	0.0825 (1.24)	-0.0606 (-1.56)	0.844	155
Q2	0.140 (1.10)	0.891*** (30.89)	0.566*** (16.46)	0.249*** (6.66)	-0.0444** (-2.04)	0.922	155
Q3	-0.189 (-1.23)	0.884*** (25.37)	0.486*** (11.71)	0.430*** (9.51)	-0.0444* (-1.69)	0.878	155
Q4	-0.348** (-2.31)	0.977*** (28.65)	0.547*** (13.46)	0.321*** (7.26)	-0.0582** (-2.26)	0.905	155
Q5	-0.452** (-2.15)	1.083*** (22.78)	0.610*** (10.76)	0.194*** (3.15)	-0.107*** (-2.98)	0.866	155
Q1-Q5	1.098*** (2.82)	-0.212** (-2.40)	0.319*** (3.03)	-0.112 (-0.98)	0.0462 (0.69)	0.115	155

*t* statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table A.II

**Weekly Categorization Bias Strategy (EW): High Official Industry Mover**

In this table we restrict the portfolios to *Strong Official Industry Follower* stocks (first tercile of  $\rho$ ; see the appendix for a formal definition). We sort stocks that have abnormally closely followed their official industries in week  $t - n$  (low  $\rho$ ) into five quintiles of the difference between the stock's official and fundamental industry return during week  $t - n$ , i.e.  $(r_{j:official;t-n} - r_{j:fundamental;t-n})$ . All portfolios are equally weighted and rebalanced every week.

	alpha	mktrf	smb	hml	mom	R2	N
<b>Panel A: Signal (t-1)</b>							
Q1	0.308*** (7.54)	0.939*** (59.57)	0.847*** (30.20)	0.172*** (6.37)	-0.0922*** (-5.64)	0.898	675
Q5	-0.176*** (-4.25)	1.000*** (62.40)	0.770*** (27.01)	0.0704** (2.56)	-0.101*** (-6.07)	0.903	675
Q1-Q5	0.484*** (7.34)	-0.0609** (-2.39)	0.0771* (1.70)	0.102** (2.33)	0.00863 (0.33)	0.0230	675
<b>Panel B: Signal (t-2)</b>							
Q1	0.174*** (4.20)	0.939*** (58.74)	0.770*** (27.08)	0.194*** (7.08)	-0.0995*** (-6.01)	0.892	674
Q5	-0.0205 (-0.50)	1.050*** (66.63)	0.791*** (28.22)	0.0677** (2.51)	-0.0817*** (-5.01)	0.913	674
Q1-Q5	0.194*** (2.94)	-0.110*** (-4.33)	-0.0207 (-0.46)	0.126*** (2.89)	-0.0178 (-0.67)	0.0566	674
<b>Panel C: Signal (t-3)</b>							
Q1	0.216*** (5.33)	0.996*** (63.60)	0.822*** (29.52)	0.190*** (7.10)	-0.0640*** (-3.95)	0.905	673
Q5	-0.0210 (-0.49)	0.999*** (60.01)	0.691*** (23.33)	0.0312 (1.09)	-0.127*** (-7.34)	0.895	673
Q1-Q5	0.237*** (3.46)	-0.00365 (-0.14)	0.131*** (2.78)	0.159*** (3.52)	0.0627** (2.29)	0.0302	673
<b>Panel D: Signal (t-4)</b>							
Q1	0.134*** (3.27)	1.013*** (64.33)	0.795*** (28.35)	0.244*** (9.03)	-0.0621*** (-3.81)	0.905	672
Q5	0.00915 (0.23)	0.994*** (63.66)	0.744*** (26.78)	0.0856*** (3.20)	-0.133*** (-8.25)	0.907	672
Q1-Q5	0.124* (1.95)	0.0195 (0.79)	0.0508 (1.16)	0.158*** (3.74)	0.0712*** (2.78)	0.0256	672
<b>Panel E: Signal (t-5)</b>							
Q1	0.0963** (2.42)	1.001*** (65.21)	0.893*** (32.70)	0.119*** (4.54)	-0.0277* (-1.75)	0.912	671
Q5	0.0630 (1.52)	0.977*** (61.27)	0.712*** (25.07)	0.104*** (3.82)	-0.184*** (-11.16)	0.903	671
Q1-Q5	0.0333 (0.53)	0.0233 (0.97)	0.181*** (4.22)	0.0148 (0.36)	0.157*** (6.26)	0.0951	671
<b>Panel F: Signal (t-6)</b>							
Q1	0.134*** (3.33)	1.042*** (67.23)	0.840*** (30.40)	0.100*** (3.79)	0.00481 (0.30)	0.913	670
Q5	0.0228 (0.53)	0.960*** (58.12)	0.680*** (23.09)	0.120*** (4.24)	-0.198*** (-11.55)	0.893	670
Q1-Q5	0.111* (1.72)	0.0817*** (3.28)	0.160*** (3.60)	-0.0193 (-0.45)	0.202*** (7.85)	0.133	670

$t$  statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table A.III

**Fama-MacBeth (1973) Forecasting Regressions (Monthly Returns)**

This table shows Fama and Macbeth forecasting (1973) regressions. The dependent variable is the monthly stock return of firm  $j$ . Equations (3), (4) and (5) are restricted to Strong Official Industry Follower stocks. We restrict the regressions to stock-month observations for which the stock price at the end of the preceding month is greater or equal to \$5. All variables are defined as in table VII.

	(1)	(2)	(3)	(4)	(5)
$r_{j,official;t-1} - r_{j,fundamental;t-1}$	-0.07 (-4.99)	-0.10 (-5.02)	-0.11 (-5.33)	-0.14 (-5.92)	-0.12 (-5.84)
Med Ind Follower $_{j,t-1}$		0.00 (1.90)			
Weak Ind Follower $_{j,t-1}$		0.00 (1.90)			
Med Ind Follower $_{j,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$		0.03 (2.58)			
Med Ind Follower $_{j,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$		0.04 (3.50)			
Mid Cap $_{j,t-1}$				0.00 (0.44)	
Large Cap $_{j,t-1}$				0.00 (0.49)	
Mid Cap $_{j,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$				0.03 (0.89)	
Large Cap $_{j,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$				0.06 (2.22)	
Med Analysts $_{j,t-1}$					0.00 (2.15)
High Analysts $_{j,t-1}$					0.01 (2.18)
Med Analyst $_{jt,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$					0.01 (0.33)
High Analyst $_{jt,t-1} \times (r_{j,official;t-1} - r_{j,fundamental;t-1})$					0.00 (-0.00)
$r_{j,t-1}$	-0.02 (-3.42)	-0.02 (-3.44)	0.05 (2.20)	0.05 (2.28)	0.05 (2.49)
$r_{j,t-2;t-12}$	0.00 (0.57)	0.00 (0.58)	0.00 (0.27)	0.00 (0.28)	0.00 (0.34)
$\ln(M)_{j,t-1}$	0.00 (-0.45)	0.00 (-0.46)	0.00 (0.28)		
$\ln(BE/ME)$	0.00 (0.51)	0.00 (0.50)	0.00 (0.95)	0.00 (0.94)	0.00 (1.23)
HAC adjusted t statistics in parentheses.					