

Global Business Networks

Christian Breitung* and Sebastian Müller⁺

July 10, 2023

Abstract

Applying large language models (GPT-3, Luminous, T5-XXL) to business descriptions of more than 79,000 firms, we construct business networks (*BNs*) of economically linked firms across the globe. We run multiple evaluation tasks and find that our networks are better suited to identify relevant competitors, suppliers, and customers across the globe than traditional industry classifications. We further demonstrate the usability of our networks by examining global lead-lag effects and M&A activity. Our results emphasize the importance of analyzing global economic links rather than exclusively focusing on domestic relations.

Keywords: Business Network, Textual Analysis, Natural Language Processing, Transformer, GPT-3, Large Language Models

We want to thank Vincent Bogousslavsky, Gerard Hoberg, Gordon Phillips, Jeffrey Pontiff, Ronnie Sadka, Milena Wittwer, and the participants of the seminars at the Carroll School of Management and the TUM School of Management for valuable comments and feedback.

*Technical University of Munich, TUM School of Management, Campus Heilbronn, Center for Digital Transformation, Am Bildungscampus 9, 74076 Heilbronn, Germany (christian.breitung@tum.de, Phone: +49 7131 264 18 822)

⁺Technical University of Munich, TUM School of Management, Campus Heilbronn, Center for Digital Transformation, Am Bildungscampus 9, 74076 Heilbronn, Germany (sebastian.mueller.hn@tum.de, Phone.: +49 713 126418806)

1. Introduction

In recent decades, global economies have experienced a growing trend of specialization. Companies in advanced economies frequently focus on producing highly specialized products, which leads to considerable heterogeneity within industries. For instance, two automobile manufacturers may offer comparable products with distinct features, such as electric or gasoline engines, or cater to different market segments, such as luxury or budget consumers. Additionally, these companies may vary in aspects such as digitization levels, supply chain resilience, and geographical locations. Previous research emphasizes that traditional sector and industry classifications, like SIC or NAICS, may not effectively represent this within-sector variety (Hoberg and Phillips, 2016). Instead, each company has a unique network of affiliated competitors, suppliers, and customers, interconnected through economic ties.

Identifying economically linked firms is essential in several scenarios. For example, acquirers can utilize publicly available data on a target’s peers to appraise its value (comparable company analysis), particularly for smaller or foreign targets, where the acquirer’s knowledge may be limited. Moreover, economic links can be used to explore research questions in economics, such as measuring competition intensity, identifying industry boundaries, and examining industrial organization. Simultaneously, investors may wish to evaluate the effect of news regarding economically linked firms on the focal company. Since competitors, suppliers, and customers may not necessarily operate in the same country, researchers and practitioners require information on international peers to study these questions.

In this paper, we identify economic links across the globe by applying state-of-the-art context-aware textual analysis methods to business descriptions of more than 79,000 publicly traded firms across 67 countries. We propose different global Business Networks (BNs) using the following pre-trained language models: A Sentence Transformer (T5-XXL), a Generative Pre-Trained Transformer (GPT-3) developed by OpenAI, the company behind the popular chatbot ChatGPT, and *LUMINOUS*, a multilingual model provided by the German startup *Aleph Alpha*.

Our results indicate that the *BNs* constructed using large language models (*LLMs*) exhibit improved precision compared to traditional word-based networks and networks based on well-known industry classifications. Next to competitors, our networks also uncover essential relationships with suppliers and customers.

We showcase the usability of our networks in two dimensions. First, we analyze the lead-

lag effect in global stock markets. We find that investors overlook important information when focusing on US links only, highlighting the unique contribution of our networks. A US long-short portfolio based on the past performance of global peers significantly outperforms US peer-based portfolios by a risk-adjusted 46 basis points (bps) per month. We further demonstrate that our *BNs* reveal economic links that are not discovered via shared analyst coverage, industry membership, similar firm characteristics, and past stock return correlations. Second, we show that acquiring firms are more likely to acquire firms with a high business description similarity. This effect persists even after controlling for industry, country, size, and several other firm characteristics.

Our *BNs* share similarities with the *Text-based Network Industry Classifications* (TNIC) proposed by [Hoberg and Phillips \(2010, 2016\)](#), but there are also notable differences.

In contrast to *TNIC*, our networks also contain links with suppliers and customers in addition to firms that offer similar products (competitors). Moreover, the *TNIC* dataset is confined to US firms, relying on the business section of 10-K filings (Item 1). Our *BNs* encompass economic links across numerous countries, drawing upon business descriptions for global stocks provided by Refinitiv. We do not use international annual reports because significant variations exist across countries, firms, and time, as highlighted by [Breitung and Müller \(2022\)](#). This makes it challenging to identify business sections in international firm reports. Moreover, collecting international reports is both time and resource-intensive, as no Application Programming Interface (API) like *EDGAR*¹ exists for procuring international firm reports. In contrast, obtaining business descriptions from Refinitiv is a more streamlined process requiring less computational resources, ultimately increasing our analyses' replicability.

Furthermore, concerning the brevity of business descriptions, we do not advocate traditional word-based techniques, such as the *bag-of-words* (*BOW*) approach that [Hoberg and Phillips \(2016\)](#) rely on. We present anecdotal evidence suggesting that these methods are susceptible to inaccuracies when applied to concise business descriptions², given the limited number of business-specific words present. For instance, we discover that according to *BOW*, the oil and gas producer *BP* is among the most related firms in the case of the car manufacturer *Ford*, which we trace back to non-industry associated words such as "company", "activity" and "segment" rather than similarities in the business operations.

¹Firms that are subject to disclosure requirements in the US submit their files to EDGAR, the Electronic Data Gathering, Analysis, and Retrieval system of the Security Exchange Commission.

²The business descriptions obtained via Refinitiv do not exceed three hundred words.

Instead, we capitalize on the latest advancements in Natural Language Processing (NLP) and construct yearly updated global business networks by applying pre-trained language models to business descriptions. We use a state-of-the-art Sentence Transformer model provided by Ni et al. (2021) and access the *LLMs* of commercial startups, such as OpenAI and Aleph Alpha, through their APIs.

Using the embeddings (semantically sensitive vector representations) obtained from these models, we generate a cosine similarity matrix of dimensions $n \times n$, where n denotes the total number of descriptions of firms actively traded in the previous year. In the final step, we determine the 99th percentile of the cosine similarity distributions and classify firms with a cosine similarity above this threshold as economically related.

To pinpoint the optimal model that aligns with our objectives, we generate numerous *LLM* networks and assess their efficacy across a spectrum of performance metrics. First, we compare the share of firms whose peers operate in the same industry, country and NYSE size decile. We find evidence that supports our argument that context-aware networks are less likely to list spurious firm relations. Second, we calculate pairwise overlaps and correlations among all networks under examination and find that the context-aware networks tend to exhibit the highest overlap and strongest return correlations. Third, following Eisdorfer et al. (2021), we extract disclosed competitors from US annual reports and evaluate which network detects most of them. We find that context-aware networks substantially outperform a word-based network. With *LUMINOUS*, we correctly identify 53.98% of the disclosed US competitors, whereas a word-based network discovers only 15.47%. Using *TNIC* instead, we identify 43.98% of the competitors. The outperformance of *LUMINOUS* persists if we control for differences in the number of predictions per network, suggesting that our context-aware networks effectively capture relevant economic links. Finally, we follow the approach of Guo et al. (2023) and extract economic peers from merger and acquisition filings submitted to the SEC in 2022. We again find that our business networks achieve comparable performance to *TNIC*.

Motivated by these results, we demonstrate the potential of our *BNs* by studying informational spillover effects, such as the lead-lag effect in stock returns. The lead-lag effect is a well-documented anomaly for the US stock market in the financial literature. As summarized by Ali and Hirshleifer (2020), a vast body of literature establishes predictive links among firms grouped within the same industry (Moskowitz and Grinblatt, 1999; Hoberg and Phillips, 2018), sharing a similar geographic location (Parsons et al., 2020), related through supply

chains (Cohen and Frazzini, 2008; Menzly and Ozbas, 2010), or utilizing similar technologies (Lee et al., 2019). Moreover, Cohen and Lou (2012) find lead-lag effects from single-segment to multi-segment corporations from the same industry, and Müller (2019) identifies economic links between stocks with similar stock characteristics.

Nevertheless, there has been relatively limited research on lead-lag effects in international stock markets. In this regard, Huang (2015) finds that past industry-level returns from foreign countries can predict the returns of US multinational stocks. Grüner et al. (2018) examine whether lead-lag effects are observable and exploitable in the stock markets of the G7 countries, but they do not investigate spillovers from economically linked firms. Instead, they study lead-lag effects from large to small stocks and between stocks with high and low analyst coverage and institutional ownership.

Our study differs from these two studies in at least two critical dimensions. On the one hand, we examine the entire universe of global stocks rather than solely focusing on US multinationals. On the other hand, we demonstrate how sophisticated textual analysis techniques can be utilized to identify firm-specific peers across the globe. This is in contrast to Huang (2015) who studies spillovers using industry classifications. By going long (short) in stocks whose economically linked performed best (worst) in the past month with monthly rebalancing from 1996 to 2021, and relying on US stocks only, we find evidence that aligns with the literature on a lead-lag relationship. We obtain a significant monthly seven-factor alpha of 63 bps using a word-based network (BN_{BOW}). If we use context-aware BN s instead, we observe substantially larger monthly seven-factor alphas between 89 and 103 bps with a t-statistic of up to 9.41. Networks based on the four-digit SIC classification generate a lower alpha of 79 bps. Since a portfolio based on a time-invariant version of the TNIC dataset yields a 103 bps alpha, we conclude that our networks are well suited for capturing the lead-lag effect.

If we construct US portfolios based on the past performance of economically linked global stocks instead, we observe a substantial and statistically significant increase in alpha up to 149 bps across the different networks, suggesting that value-relevant information from foreign markets might be less efficiently priced.

In our analysis of global financial markets, we observe even stronger lead-lag effects. Here, long-short portfolios based on context-aware networks generate monthly seven-factor alphas of up to 235 bps, substantially outperforming word-based networks. Our findings further suggest that our context-aware networks reveal unpriced economic links with firms from different indus-

tries, as a network based on firms within the same four-digit SIC code generates a substantially lower alpha of 142 bps.

To test whether limits to arbitrage may explain the extraordinarily high alphas, we construct long-short portfolios based on stocks below and above the fifth NYSE size decile. We indeed observe a highly significant size effect, suggesting that limits to arbitrage explain a substantial fraction of the lead-lag effect. While we obtain seven-factor alphas up to 175 (254) bps for smaller stocks in the US (globally), the alphas are substantially lower for larger stocks but remain highly significantly at the 1% level with 90 (125) bps.

We also examine to what extent our networks detect economic links that other approaches suggested in the literature do not capture. We therefore conduct Fama MacBeth regressions (Fama and MacBeth, 1973) that account for shared analyst coverage (Ali and Hirshleifer, 2020), industry membership, similar firm characteristics (Müller, 2019) and return correlations of stocks within the same industry (Gatev et al., 2006; De Franco et al., 2011). Our results indicate that our *BNs* identify relevant economic links that may not be identified otherwise.

Two types of potential look-ahead biases potentially affect our analyses. First, our reliance on recent over historical descriptions might inflate our estimates of the lead-lag effect. However, by replicating our analyses with historical business descriptions from Refinitiv SDC Platinum, we gather evidence that suggests our analyses are largely unaffected by this bias.

Second, the language models, having been trained on extensive datasets that coincide with our evaluation period, might induce a more nuanced bias. Mitigating such bias presents a greater challenge, as no existing language model has been exclusively trained on data preceding our evaluation period, and computational constraints deter us from crafting such a model.

Nevertheless, we propose a methodology to investigate whether our networks are affected by such a look-ahead bias. Initially, we ascertain the similarities in business descriptions among stocks within specific sectors. For this, we utilize a transformer model introduced in 2019, as well as an OpenAI model that has undergone training up to September 2021. Next, we scrutinize whether the correlation of the computed cosine similarities is smaller for sectors that endured substantial impacts from the *COVID-19* pandemic, such as the *Biotech and Pharmaceuticals* sector, compared to others. Interestingly, our findings challenge our original conjecture. Contrary to our initial hypothesis, we find that the correlation is higher for the *Biotech and Pharmaceuticals* sector, which implies that employing a model trained with more recent data is unlikely to introduce a significant bias to our business networks.

Additionally, our study demonstrates that our business networks can be utilized to identify target firms in M&A deals effectively. Specifically, we discover that approximately 45% of the target firms are ranked among the top 100 firms with the highest business description similarity to the acquiring firm when utilizing *LLM* networks. In contrast, when employing bootstrapping to select 100 firms from the same industry, on average, only 9% of the firms were eventually acquired.

Inspired by these promising findings, we conduct a logistic regression analysis, controlling for industry, country, size, profitability, and other relevant variables. Our results indicate that, on average, target firms exhibit a higher degree of similarity in their business descriptions compared to non-target firms. These findings suggest that business description similarities may help uncover potential M&A targets.

Our study makes several notable contributions to the finance and economics literature. Firstly, by presenting our approach to identifying economic links among firms in a global setting, we aim to provide researchers with a valuable tool for exploring previously inaccessible research questions. Secondly, we contribute to the growing literature on textual analysis in finance. Although prior research employs textual analysis to identify competitors (Eisdorfer et al., 2021), measure competition intensity (Li et al., 2013) and derive a time-varying industry classification (Hoberg and Phillips, 2016) with traditional textual analysis methods, our work provides compelling evidence for the use of advanced NLP tools in finance and economics. We extend the work of Breitung and Müller (2022), who show that context-aware similarity measures are superior to word-based measures in the context of new information detection, by studying whether *LLMs* such as GPT-3 are well suited to detect similarities in business operations. Finally, our results highlight the importance of considering global economic links rather than domestic ones only.

The structure of this paper is outlined as follows. Section 2 describes how we construct the business networks. Section 3 provides an overview of the business descriptions and stock data obtained from Refinitiv. In Section 4, we assess the performance of various business networks on multiple dimensions. In Section 5 and 6, we study two applications for our business networks and discuss how the networks may be modified to serve other applications in section 7. Finally, we conclude our findings in Section 8.

2. Methodology

2.1. Business Networks

To fully understand global economic links, it is crucial to identify firms with related business operations. In theory, various approaches can achieve this. The most straightforward approach is to define economically linked firms as those operating in the same industry. However, this assumes that all firms operating in the same industry are economically linked, which does not necessarily have to be the case. At the same time, firms may also be linked with companies from other industries, such as suppliers or customers that operate in different industries.

An alternative approach could be to extract competitor, supplier, and customer information from disclosures of publicly traded firms. For example, [Eisdorfer et al. \(2021\)](#) extract competitors from the business section of US annual reports. However, since firms have some flexibility in disclosing business relations and extracting firm names from text is error-prone, the resulting business network might lack essential links. Moreover, as indicated by [Breitung and Müller \(2022\)](#), international annual reports lack a harmonized structure which complicates the extraction of relevant information, even if it is available. This problem would also occur if we followed the approach of [Hoberg and Phillips \(2016\)](#), who identify similar firms by comparing the business sections of US annual reports (10-K filings). Internationally, many firms do not disclose information on its business operations in a separate section.

A different approach involves the comparison of past stock returns. The underlying assumption is that the most similar firms should have the highest return co-movements due to their similar risk exposures. We could thus identify related firms by comparing historical correlations of daily stock return data. [Gatev et al. \(2006\)](#) show that this approach can be used to construct a profitable pairs-trading investment strategy. Based on the idea that firms which comoved in the past should also comove in the future, they find that shorting the winner and buying the loser may generate excess returns up to 11%.

However, a disadvantage is that firms could randomly comove and thus appear related even though there are not. While we could mitigate this effect by identifying the most similar firms within an industry, similar to [De Franco et al. \(2011\)](#), who correlate past earnings within industries to identify firms with similar financial statement information, we would not be able to detect economic links with firms outside the own industry.

Within this paper, we construct business networks (BNs) based on widely available business descriptions from the commercial data vendor Refinitiv. To do so, we first collect business

descriptions for more than 79,000 firms worldwide. Then, we obtain a vector representation for each business description by either applying word-based (*BOW*) or context-aware (via language models) similarity measures. In the case of *BOW*, we identify a set of unique nouns from the universe of business descriptions and construct word-frequency vectors, similar to the implementation of [Hoberg and Phillips \(2016\)](#). In the case of the *language model* approach, we apply a language model to all business descriptions. If a description exceeds the maximum number of tokens the model supports, we split the description into sentences, obtain embeddings on the sentence level, and pool them into a single embedding to ensure that we have one vector representation for each description in our dataset. We provide a more detailed description of the language models and the process of obtaining embeddings in the upcoming section.

Considering that not all 79,000 firms in our dataset were active throughout the sample, we construct yearly networks based on firms that were actively traded in the previous year. Using the embeddings of the business descriptions, we create an $n \times n$ cosine similarity matrix, where n represents the number of actively traded firms. This matrix allows us to measure the similarity between all business descriptions. We then rank firms based on their cosine similarity in descending order.

We need to set a cosine similarity threshold to isolate firms with a sufficiently high business description similarity. Defining a fixed number of highest-ranked firms as economically related is not ideal, given that we observe differences in the number of economic links a firm may have. Firms that offer specialized products and internalized a large share of the supply chain should have less economic links than firms who offer generic products and outsourced a large share of their supply chain. Therefore, we argue in favor of a relative cosine similarity threshold as this allows for differences in the number of related firms. We also control for differences in the distribution of cosine similarities across different models using a percentile value rather than the same cosine similarity threshold across various models. More precisely, we consider firms to be economically linked if their business description similarity is in the top 1% of all values in the respective model's cosine similarity matrix.³ It's worth noting that our network might suggest that there are no sufficiently similar companies. In this case, we do not include the firm in our *BNs*.

³Note that we determine the 99th percentile based on a random subset of 1000 firms. This approach helps us reduce the computer memory (RAM) requirements while still obtaining accurate results.

2.2. Similarity measures

In the realm of identifying firms with similar business operations using textual data, choosing a suitable similarity measure is essential. One commonly adopted approach is *bag-of-words*, where text is encoded as a vector of its constituent words. However, the computational complexity of this approach can vary greatly, depending on the chosen implementation. For instance, [Cohen et al. \(2020\)](#) calculate pairwise similarities using a list of relevant words determined on a per-pair basis, resulting in a challenging parallelization problem due to the lack of a fixed vector size. Assuming 40,000 active firms on average, implementing [Cohen et al. \(2020\)](#) would necessitate the creation of 800 million unique word vectors per year.⁴ An alternative implementation that is more feasible for our purpose is presented by [Hoberg and Phillips \(2016\)](#), who first identify a set of relevant words and then construct binary high-dimensional vectors based on the presence of words in the text. By doing so, a parallelization of the cosine similarity calculations becomes feasible as all vectors possess a common dimensionality.

Despite its widespread use, the *BOW* approach has several limitations. First, measuring document similarity by counting common words may not be an accurate proxy, as two firms with dissimilar businesses may still use similar words in distinct contexts. For instance, the word *security* could relate to *cyber security*, *production security*, or *health security*. Without considering the context, its actual meaning cannot be identified. This is only exacerbated in our case, as we deal with large amounts of short text, where the amount of informative words is typically low. Differentiating between informative and uninformative words requires specialized language and domain knowledge, which might be available for widely-used languages like English but could be lacking for less spoken languages. Second, we cannot use *BOW* to compare text that is written in different languages. While one could obtain translations to tackle this problem, this process often induces an information loss. Finally, the *BOW* approach does not control for synonyms. For example, firms might describe their business as "selling cars" or "selling automobiles". Consequently, the *BOW* approach might overlook similarities in text.

Researchers may circumvent the problems mentioned above by leveraging the latest advances in NLP. These advances were sparked with the invention of the transformer architecture ([Vaswani et al., 2017](#)), an architecture that substantially improved the training speed and performance of deep learning models. Researchers quickly realized the potential behind this improvement, as it allows machine learning models to gain knowledge from text in an unsuper-

⁴The number of unique word vectors is obtained by taking the square of 40,000 and dividing it by two.

vised manner. When Google released the pre-trained transformer model *BERT* (Devlin et al., 2018), other tech firms were soon to follow. Facebook presented an improved version, RoBERTa (Liu et al., 2019) in 2019. OpenAI, a startup founded by Sam Altman, Elon Musk, and Peter Thiel, received great attention with its GPT-3 release in 2020 (Brown et al., 2020). Aleph Alpha, a German startup, presented their *luminous* model family and showed that the largest version of this model might achieve comparable results even though it is substantially smaller than GPT-3 (Aleph-Alpha, 2023). Most recently, OpenAI presented GPT-4 (OpenAI, 2023), an even more powerful language model that might even process "sparks of Artificial General Intelligence" (Bubeck et al., 2023).

What these language models have in common is the transformer architecture, which allows them to determine the semantic similarity of text when fine-tuned likewise. Fine-tuned models may vary in size, speed, and the number of tokens they can process. For instance, the *Sentence Transformer* models proposed by Reimers and Gurevych (2019) may process up to 512 tokens, making them suitable for obtaining embeddings for sentences or small paragraphs instead of entire documents. Currently, there are 38 pre-trained Sentence Transformer models available online, each trained on different datasets and methods. Some models perform well in semantic similarity, while others excel in semantic search tasks. The *sentence-t5-xxl* (*T5-XXL*) model, which is specifically trained on sentence similarity, yields the best average performance on a set of 14 diverse sentence similarity tasks. Meanwhile, the *multi-qa-mpnet-base-cos-v1* model shows the best average performance on six semantic search tasks. The *all-mpnet-base-v2* model provides strong results in both tasks. Comparing the speed and size of these models, *T5-XXL* is significantly larger and thus has a 56 times slower inference time than the other two models.⁵

Next to these open-source models, commercial solutions are available. OpenAI offers pre-trained models for sentence similarity via an API. The *text-embedding-ada-002* model (*ADA-002*), a GPT-3 derivative, is capable of generating vector representations for up to 8192 tokens. According to OpenAI, this model has achieved state-of-the-art performance on the SentEval dataset, an evaluation toolkit for universal sentence representations (Conneau and Kiela, 2018). Next to OpenAI, the German startup *Aleph Alpha* also offers access to their *luminous* model via an API.

In this paper, we consider three different pre-trained language models: A Sentence Trans-

⁵For more information on the available models and their performance, see: https://www.sbert.net/docs/pretrained_models.html

former model (*T5-XXL*), the GPT-3 derivative *ADA-002*, as well as the *luminous-base* model from Aleph Alpha. This allows us to investigate whether our approach to identify economic links is robust to different language models. Furthermore, we may observe whether there exist significant differences in the accuracy of open-source and commercial models.

3. Data

We collect English business descriptions for more than 79,000 global stocks from Refinitiv at the beginning of August 2022. These descriptions belong to stocks that were actively traded in 2022 as well as firms that delisted or filed bankruptcy within the last decades. The descriptions do not contain more than 300 words and we only consider those comprising at least ten words to ensure sufficient information on the business of a company. Although these brief descriptions offer a less comprehensive view of a company’s business compared to the 10-K filings, they cover the most crucial aspects.

By manually evaluating randomly chosen descriptions, we find that most business descriptions adhere to the same structure. The first sentence typically describes the core business of a firm. It is then followed by an enumeration of the company’s segments, which are then explained in greater detail. Finally, many descriptions also list the most important products offered by the firm at the end. To ensure that the descriptions also contain sufficient information for smaller firms, we illustrate the distribution of the number of words contained in a business descriptions across different NYSE size deciles.

[Figure 1 about here.]

According to the boxplot provided in Figure 1, we find that the largest firms exhibit marginally higher word counts than those in the lowest decile. Nonetheless, over 95% of descriptions for firms in the lowest decile exceed forty words.

Following [Ibriyamova et al. \(2019\)](#), we evaluate how accurate the business descriptions are by comparing them to the Item 1 section of the firm’s 10-K filing in 2022. More precisely, we obtain embeddings for all US business descriptions and Item 1 sections as of 2022.⁶ Next to the embeddings from OpenAI that may represent up to 8196 tokens, we test the similarity

⁶We receive pre-processed Item 1 sections from *Qannual*, a company specializing in the information extraction of annual reports.

based on the *T5-XXL* model that may process no more than 256 tokens.⁷ If the length of the Item 1 section exceeds this threshold, we consider the first 8196 (256) tokens. We then rank all business descriptions based on their cosine similarity to each Item 1 and count for how many firms the corresponding business description ranks among the most similar descriptions.

[Figure 2 about here.]

According to Figure 2, we find that in 84% of all cases, the correct Refinitiv business description ranks highest if we consider the embeddings provided by OpenAI. If we consider the three highest ranked Refinitiv business descriptions, the accuracy increases to 90%. Even though smaller, we observe a similar effect for the other embedding model (*T5-XXL*). Given the high similarity between the Item 1 sections and Refinitiv business descriptions, we thus conclude that the quality of the business descriptions should be adequate to describe a firm’s business activities.

[Table 1 about here.]

Table 1 summarizes the number of business descriptions we collect. The dataset encompasses business descriptions from 67 countries. Significant disparities exist in the number of descriptions per country due to differences in the number of publicly traded firms in these countries. The dataset includes over 22,000 US firms, accounting for around 28% of all descriptions, while Japan, Canada, China, and the UK make up an additional 30%.

It is important to note that this dataset only contains business descriptions as of August 2022. This lack of historical descriptions could introduce bias when analyzing historical data. Firms change their business operations over time, which means that firms that appear economically linked today may not have operated in similar segments in the past.

While we cannot access historical descriptions via *Refinitiv Workspace*, we collect historical descriptions for a large share of the firms in our dataset from *Refinitiv SDC Platinum*. *SDC Platinum* provides an extensive overview of international M&A deals over the last decades, including business descriptions of the acquiring and targeted firm of the time of the deal. For instance, we observe 128 events where the automotive company *Ford* was involved.⁸

⁷We refrain from calculating the similarity according to the *LUMINOUS* model, as it is considerably more costly than the embedding model provided by OpenAI.

⁸These events also include stock repurchases and recapitalizations, which explains the large number of events.

[Table 2 about here.]

Table 2 shows the most recent as well as the historical descriptions of *Ford* and its competitor *General Motors*. It seems that the business descriptions got more informative over time. While the description of *Ford* allocated to a deal in 1983 contains only three sentences, the most recent description obtained from Refinitiv Workspace counts seven sentences. Historical descriptions thus seem to contain less information that may be used to identify related firms.

Note that we do not observe variations in the historical descriptions of Ford in the SDC Platinum database. In contrast, we observe various business description versions of *General Motors*. This discrepancy can be attributed to Refinitiv’s policy, which states that business descriptions are updated exclusively in case of a significant transformation in a firm’s operations. In total, we collect around 112,000 unique business descriptions in combination with the announcement date and other M&A data of the respective M&A deals for more than 70,000 firms.

In addition, we obtain monthly stock return data from Refinitiv. While many researchers rely on CRSP as the data provider for US stock reports, we advocate using a single provider for US and international stocks. To ensure the reproducibility of our results with stock return data from CRSP, we only consider US stocks that we can identify in Refinitiv and CRSP (11979).

4. Evaluation of the Business Networks

Evaluating the accuracy of *BNs* presents a challenge due to the lack of a definite ground truth to compare against. We therefore rely on alternative evaluations described in the following.

4.1. Anecdotal evidence

We examine the competitors of *Ford* in 2021 through various networks, presenting a selection of anecdotal evidence. Since our networks include competitors, suppliers, and customers, we rank firms with the highest similarity in terms of past year market capitalization first. The reason is that the most relevant competitors should have a similar size. Furthermore, we prioritize domestic firms if two competitors belong to the same NYSE size decile. This is because domestic firms are more likely to be direct competitors.

[Table 3 about here.]

Table 3 displays the closest competitors of the car producer *Ford* according to the different BNs. According to networks based on large language models, *General Motors* emerges as the closest counterpart, consistent with the TNIC dataset. In contrast, a *BOW* network ranks the German car manufacturer *BMW* highest. While being a car manufacturer itself, *BMW* does not offer pickup trucks, in contrast to *Ford* and *General Motors*. At the same time, *BMW* also produces motorcycles, a product that is not offered by *Ford*, suggesting that *BMW* might be less similar than *General Motors*.

Furthermore, a *BOW* network suggests that the oil and gas producer *BP* is among the five most related firms. To determine the cause of this likely inaccurate classification, we delve into a detailed analysis of the lemmatized nouns in both company descriptions to comprehend the source of similarity. We find that both descriptions contain words like "Africa", "Europe", "Asia", and "America" that are not related to the business operations but to the area of operation instead. We further find that both descriptions contain words like "company", "activity", and "segment" that offer limited insights into a firm's business, potentially leading to erroneous conclusions regarding similarities in business descriptions.

Given its limitation to US firms, the TNIC dataset does not mention competitors like *Toyota*, *Hyundai*, or *Volkswagen*, which emphasizes the importance of considering a global network to effectively identify the most important competitors.

4.2. Comparison of the BNs

We aim to comprehensively and systematically understand the differences between the constructed BNs by analyzing multiple dimensions. Firstly, we investigate to what extent the models identify relations between firms in the same industry, domicile, and with similar size. Secondly, we calculate the "overlap percentage" of a focal network's links that can also be found in other networks relative to the total number of relations in the focal network. We present the results in a matrix covering all network pairs. Finally, we estimate the correlation between the various networks based on the average return of economically linked firms in 2021 to determine how similar the identified economic peers of the different networks perform.

[Table 4 about here.]

Panel A in Table 5 provides an overview of each firm's average number of identified relations and how many recognized links include firms from the same industry, country, and NYSE size decile. We observe substantial differences in the number of relations. As expected, the TNIC

network comprises fewer firms than our networks given its restriction to US firms. However, substantial differences exist across the three *LLM* networks. On average, the *T5-XXL* network comprises the highest number of relations with a mean value of 351 and a median of 255. This may seem counter-intuitive, since we use the same percentile for each network. However, because we determine the percentile over the entire universe of cosine similarities, the concentration of high cosine similarities may vary across the networks. For some networks, the 1% highest cosine similarities may spread across many firms, whereas they are more concentrated in other networks. Consequently, some firms may drop out from some networks and remain in others, leading to differences in the average number of links per firm.

Moreover, we observe substantial differences across the share of relations comprising firms of the same industry. It is substantially higher, at around 18%, for the three context-aware networks, whereas only 10.47% of the relations in a word-based network include firms from the same four-digit SIC code. This effect persists if we consider broader industry classifications and is in line with the argument that word-based networks cannot control for the context of words, leading to the identification of spurious relations. Furthermore, we observe that up to 28% of the links identified by networks based on language models include firms from in the same country. This is substantially higher than the share suggested by the word-based network (19.52%) and the SIC network (14.45%). It seems that the language-based networks prioritize domestic links.

Panel B shows the overlap percentages among US firms for the different *BNs*. In the US, the BN_{T5-XXL} contains 69.67% of the relations included in *TNIC*. This is a substantially higher share than the *BOW* network (22.92%). We also find that the highest overlaps might be observed among the networks of large language models (*T5-XXL*, *ADA-002* and *LUMINOUS*) with overlaps up to 85.29%.

Regarding stock return correlation in 2021, we observe substantial differences in the coefficients among US networks. Even though both networks are word-based, the *TNIC* dataset appears to be least correlated with the *BOW* network. We observe a positive coefficient of 49.20%. For context-aware networks, the return correlation is higher with up to 65.79% in the case of the *LUMINOUS* network. Both *T5-XXL* and *ADA-002* yield slightly lower correlations of 65.62% and 64.09%, suggesting that *LLM* networks are more in line with the *TNIC* dataset.

4.3. Performance evaluation

To evaluate the accuracy of the networks, we compare the detected relations to the competitors disclosed by US firms. To achieve this, we adopt the approach outlined in [Eisdorfer et al. \(2021\)](#) and extract competitor names from the Item 1 "competition" subsection of US annual reports (10-K filings). By comparing the detected relations with the disclosed competitors, we should gain insights into the accuracy of our networks.⁹ Note that we restrict our analyses to files published in the most recent year (2022). We further restrict our networks to US firms and peers in all networks to allow for a comparison with TNIC that does not provide information on international peers. We identify 1189 disclosed competitors for 734 unique firms covered in the TNIC dataset.

Furthermore, we follow the approach of [Guo et al. \(2023\)](#) and extract economic peers from merger and acquisition filings submitted to the SEC in 2022. More specifically, we consider the "opinion of the financial advisor" section of M&A filing documents like PREM14A, DEFM14A, and S-4/A, where advising investment banks enclose their opinion on the acquisition. If they conduct a comparable company analysis to evaluate the terms of a deal, they often disclose a list of firms they base their research on. We extract these company names and compare them to the different business networks.¹⁰ Due to variations in the structure of these reports, identifying the relevant sections can be challenging. Nevertheless, we identify 600 disclosed competitors for 108 unique US firms covered by the TNIC dataset.

[Table 5 about here.]

Panel A provides different evaluation metrics of the various networks concerning the identification of disclosed competitors. We find large differences across the networks if we consider how many disclosed competitors may be identified (recall score) based on the most related firms. While 8,96% of the highest-ranked peers in TNIC2021 are disclosed as competitors, the highest-ranked peers in a word-based network cover only 1.56% of the disclosed firms. Networks based on language models perform substantially better and discover up to 4.2%. The four-digit

⁹Instead of using the StanfordNER project, we use the python package "Spacy" and apply a transformer-based entity recognition model to identify company names from text. We match the recognized organizations with the EDGAR database operated by the SEC, the CRSP master file, and the company names of the Refinitiv database.

¹⁰To ensure that we do not detect the names of the advising banks that might be mentioned in these sections, we drop the 1% of the names with the highest number of occurrences (6%) across the various documents.

SIC network is somewhere in between with 2.64%. This pattern seems to persist with an increase in the number of predictions considered but disappears once we consider at least fifty highest-ranked predictions. While the *LUMINOUS* model identifies 41.94%, the TNIC2021 dataset covers only 38.02%.

Considering the entire network, the BN_{BOW} identifies 15.47% of the disclosed competitors. Context-aware networks perform substantially better. The $BN_{LUMINOUS}$ performs best and identifies 53.98% of all disclosed US competitors, which is also higher than *TNIC2021* with 43.98%. The BN_{T5-XXL} follows with 52.3%, and the $BN_{ADA-002}$ network correctly identifies 46.17% of all disclosed US competitors. This performance is impressive, given that our networks rest on substantially less information than TNIC2021. In contrast, if we use the *SIC* industry network, we can only identify 21.19% of the disclosed competitors, suggesting that our networks are more accurate than industry networks.

However, focusing on the *recall* score is insufficient. The reason is that there might be differences in the average number of predictions across the networks. If some networks, on average, predict a higher number of economic links, this could introduce an upward bias in the *recall* score. We therefore also calculate the *precision* score, which is identified as the number of correctly identified competitors divided by the total number of recognized relations. If we compare the *precision* score across the different networks, we find that all context-aware networks achieve values at least as high as the TNIC dataset. This suggests that the outperformance we have seen so far is not attributable to a mechanical effect of having more firms in the business network.

In Panel B, we employ the same evaluation metrics concerning the list of comparable companies disclosed by advising investment banks in the context of M&A deals. Similar to our previous findings, we find that the recall score of the Sentence Transformer model is highest once we consider more than 30 predictions. While a network based on the most recent TNIC information identifies up to 50.06% of the disclosed peers, the *T5-XXL* network identifies slightly more with 55.85% of the relations.

5. First Potential Application: Lead-lag effect

Given the previous evaluation results, we showcase the usability of our *BNs* by investigating the well-documented lead-lag effect (Hou, 2007; Cohen and Frazzini, 2008; Menzly and Ozbas, 2010; Cohen and Lou, 2012; Huang, 2015; Müller, 2019; Ali and Hirshleifer, 2020; Hoberg and

Phillips, 2018) in global stock markets. The lead-lag effect indicates a cross-predictability in returns from one stock to another, which suggests a gradual information diffusion in security markets that is inconsistent with market efficiency. These informational spillover effects have been primarily documented for economically-linked stocks. Following this argumentation, we should be able to detect spillover effects using our *BNs*.¹¹

To do so, we construct equally weighted calendar-time portfolios from 1996 to 2021. At the start of each month, we calculate the average past month’s performance of the economically linked firms on the firm level. We then pursue long (short) investments in the 20% of stocks whose most similar firms performed best (worst) in the previous month. We evaluate these portfolios using the five factors from Fama and French (2015) and momentum and short-term reversal. To ensure that outliers do not drive our results, we drop the lowest and highest 0.5% of the returns from our dataset.

Note that we construct portfolios based on our dataset of business descriptions as of 2022 in section 5.1, before we control for a potential look-ahead bias by using historical descriptions in Section 5.2. The reason why we do not use historical descriptions in the first place is the limited availability of historical descriptions which originate from SDC Platinum and hence restricted to firms being involved in M&A cases in the past.

5.1. Recent business descriptions

5.1.1. US lead-lag effects

[Table 6 about here.]

In Panel A in Table 6, we present the factor exposure of different US long-short portfolios based on the past month’s performance of US peers suggested by the various networks. Overall, our results reveal significantly positive seven-factor alphas regardless of the network used, confirming the existence of a lead-lag effect in the US.

We find that a *BOW* portfolio exhibits the lowest alpha with 63 bps, which aligns with our previous findings. Context-aware networks generate higher alphas up to 103 bps with a t-statistic of 9.41 in the case of the *LUMINOUS* portfolio. On the contrary, we observe that portfolios utilizing the four-digit SIC network yield inferior results (79 bps) yet surpassing the *BOW* portfolio. We also examine a portfolio based on a modified version of the four-digit SIC

¹¹Given that our approach does not allow us to differentiate between suppliers, customers, and competitors, we cannot explain these spillover effects in more detail.

network that excludes those relations identified in languagemodel-based networks (SIC_{NN}). This portfolio generates 26 bps lower alpha than one comprising all firms within the same four-digit SIC code. This compelling evidence indicates that *LLM* networks effectively identify economically linked firms within an industry.

It is important to note that a direct comparison between the portfolio returns of our context-aware *BNs* and the TNIC dataset is not possible. This is because *TNIC* is time-varying and thus controls for changes in business descriptions over time, while our *BNs* do not. As it is unclear whether controlling for these changes results in higher or lower alphas, we construct two portfolios, one that uses the entire TNIC dataset and another that only considers the most recent information (denoted as *TNIC2021*). The unrestricted long-short portfolio generates a monthly alpha of 167 bps. These results are broadly in line with [Hoberg and Phillips \(2018\)](#) who obtain monthly alphas between 190 and 230 bps for a shorter evaluation period (1997-2012).

In contrast, the portfolio restricted to the most recent information generates an alpha of 103 bps only, which coincides with the best performing language-based portfolio. The significantly lower 64 bps alpha for the static version suggests that the performance of our *BNs* are rather downward than upward biased.

To understand whether information from international peers is less efficiently priced, we also construct US portfolios based on past returns of global peers and present the results in Panel B. We observe a substantial increase in alpha across all networks. The *LUMINOUS* portfolio again generates the highest alpha with 149 bps representing a 46 bps increase, followed by 147 bps and 144 bps for the *T5-XXL* and *ADA-002* portfolios. The increases in alpha observed across the networks based on language models are significant at the 5% level.¹² This finding not only suggests that investors overlook value-relevant information from foreign peers, but also highlights the importance of considering international markets rather than focusing on the US only. Moreover, considering that the alpha increase we notice for the *BOW* portfolio lacks statistical significance, it further emphasizes the higher accuracy of *LLMs*.

¹²We test the significance of the differences in alpha by first calculating the difference between the portfolios' returns. Then, we regress these portfolio return differences on the seven pricing factors and test the statistical significance of the constant.

5.1.2. Global lead-lag effects

We proceed to replicate prior calculations on global markets using global factor data¹³, as US factor data is unsuitable for explaining international stock returns.

[Table 7 about here.]

In Table 7, we present the factor exposure of global long-short portfolios based on economic links with US and global firms. Overall, our results reveal significantly positive seven-factor alphas regardless of the network used, confirming the existence of a global lead-lag effect.

In Panel A, we restrict our networks to US peers and observe similar alphas for equally weighted investments into global portfolios compared to US portfolios. Again the *LUMINOUS* portfolio generates the highest alpha with 112 bps, outperforming simple industry networks. The *T5-XXL* portfolio follows with 104 bps and a portfolio based on the model provided by OpenAI (*ADA-002*) generates 96 bps monthly seven-factor alpha.

If we construct portfolios based on global peers, we observe a substantial increase in alpha across all *LLM* networks that is significant at the 1% level. The *ADA-002* portfolio exhibits the highest seven-factor alpha of 235 bps with a t-statistic of 9.49, although the other two language model-based portfolios (*T5-XXL* and *ADA-002*) generate comparably strong alphas of 234 and 225 bps. Compared to investments in US stocks, the increase in alpha induced by also considering international links seems to be larger for global stocks. This could indicate that information from foreign peers is less efficiently priced in international markets.

Conversely, the *BOW* portfolio yields a smaller alpha of 173 bps but surpasses the four-digit *SIC* portfolio. Importantly, when constraining the portfolio to within-industry peers not included in language-model networks (*SIC4_{NN}*), the alpha increases only marginally to 75 bps. It seems that our *LLM* networks thus capture all relevant international firm relations.

Given the substantially higher alphas observed for portfolios that incorporate global peers, we suggest that focusing solely on economic links with US peers is not sufficient, highlighting the contribution of our global business networks once again. Furthermore, researchers and practitioners should also broaden their scope and consider economic relations with firms outside the industry, particularly in global settings. This is supported by the substantial difference in

¹³The international factors were calculated by following the methodology mentioned on the website of Kenneth R. French as closely as possible. For additional information about the construction of global asset pricing factors, we refer to [Huber et al. \(2023\)](#), who analyze the suitability of competing asset pricing models for international stock markets.

alpha between our language model-based portfolios and the SIC4 portfolio.

5.1.3. *The impact of limits to arbitrage*

We rerun the previous experiment using portfolios of smaller and larger stocks to explore the potential influence of limits to arbitrage on our results. Smaller stocks are those below the fifth NYSE size decile, whereas larger stocks belong to the fifth decile and beyond.

[Table 8 about here.]

According to Table 8, a significant fraction of the lead-lag effect may be explained with limits to arbitrage. While we observe highly significant alphas up to 192 bps for investments into US stocks below the fifth NYSE size decile using TNIC, we observe a substantially smaller monthly seven-factor alpha of 88 bps for larger stocks. A similar pattern can be observed in the case of our language-based networks and the four-digit SIC network. Nevertheless, the observed alphas remain highly significant at the 1% level, suggesting that limits to arbitrage are not the only explanatory factor.

If we consider global peers, we observe an overall similar pattern. Here, we observe substantially lower alphas for investments into larger stocks than smaller ones. Limits to arbitrage play an essential role but may only partially explain the existence of the lead-lag effect in the US.

Furthermore, our analysis shows that the increase in alpha for portfolios based on global peers is substantially higher for firms below the fifth NYSE size decile. A portfolio that invests in small US stocks based on past global peer performance generates up to 75 bps higher alpha, as in the case of *T5-XXL*, than a portfolio based on past US peer performance. This difference is highly significant at the 1% level.

In contrast, we do not observe significant differences for similar portfolios of large stocks. Although one might expect larger firms to benefit more from global operations, larger firms might also be more diversified. Information from individual markets thus could be less relevant for large firms that serve various foreign markets.

We observe similar patterns when constructing global portfolios. Portfolios consisting of stocks above the fifth NYSE size decile generate alphas of up to 125 bps, indicating that factors beyond limits to arbitrage may contribute to the lead-lag effect. Additionally, similar to our findings in the US, we note a substantial increase in alpha of up to 155 bps when investing in small stocks based on the past performance of global rather than US peers. This difference

exhibits high significance at the 1% level, with t-statistics reaching up to 6.24. For larger stocks, the increase in alpha is not statistically significant, except for the *BOW* portfolio. While the lead-lag effect is smaller for larger stocks if we construct portfolios based on US links, it remains highly significant with alphas up to 90 bps.

We further find evidence suggesting that the performance gap between *LLM* networks and the four-digit SIC networks is larger for smaller stocks. While a SIC-based network generates an alpha of 150 bps, investments according to *T5-XXL* generate up to 254 bps per month. This indicates that links with firms from different industries provide value-relevant information that investors do not timely price.

5.1.4. Novelty of identified economic links

To understand to what extent our networks capture novel links, we construct various business networks based on alternative approaches suggested in the literature. First, we use the four-digit SIC industry networks introduced earlier (BN_{SIC4}). Second, we follow the arguments of Müller (2019) and identify economic links using similarities in firm characteristics (BN_{SESM}). Third, we construct a network based on past stock price correlations of firms within an industry. This approach combines the work by Gatev et al. (2006) and De Franco et al. (2011). We calculate correlations among firms within the same four-digit SIC code (BN_{CORR}) rather than the entire universe of stocks to reduce the likelihood of spurious correlations.

Finally, we construct a network based on shared analyst coverage, motivated by the paper of Ali and Hirshleifer (2020), who find that any of the lead-lag effects suggested in the literature turn insignificant when controlling for shared analyst coverage. The rationale is that analysts likely cover similar stocks to profit from their domain-specific information advantage. To do so, we closely follow the approach of Ali and Hirshleifer (2020) and identify all analysts who covered a particular stock in the preceding twelve months. Then, we determine all other stocks covered by at least one of these analysts in the same period. We repeat this process for each stock in each month during the investment horizon from 1996 to 2021. To ensure comparability and avoid selection bias, stocks in our *BNs* that are not covered by at least one analyst in at least one month are omitted from the analysis. The shared analyst network is updated every month to account for temporal variations.

[Table 9 about here.]

Table 9 shows the results of the Fama MacBeth regressions for US and global stocks. In line

with our previous findings, we observe a highly significant positive coefficient for our lead-lag measure BN_{T5-XXL} in and outside the US (columns one and four). Stocks whose economically linked firms outperformed in the past month thus achieve a larger return in the subsequent month. The R^2 values are 0.0754 for the US and 0.0622 internationally.

Although there is a decrease in the magnitude of the coefficient, the shared analyst network only partially captures the links identified by BN_{T5-XXL} . Our coefficient of interest remains highly significant if we control for shared analyst coverage (columns two and five). At the same time, the increase in the R^2 with the inclusion of BN_{ANA} suggests that analysts identify relations that are not included in the BN_{T5-XXL} . Once controlling for the other networks mentioned above, we observe a smaller but significantly positive BN_{T5-XXL} coefficient of 0.0766 in the US and 0.1370 globally with a t-statistic of 2.42 and 3.74, respectively (see columns three and six). Put differently, US (global) stocks in the highest quintile achieve a roughly 32 basis point (68 basis point) higher return than those in the lowest quintile, holding constant all other control variables.

In conclusion, the results obtained from the Fama MacBeth regressions presented in Table 9 indicate strong evidence of relevant relations within our BNs that may not be identified via shared analyst coverage. This is evident both in the US and internationally. Additionally, these economic links are not revealed through traditional industry classification, similar firm characteristics, or correlations in the returns of stocks operating in the same industry, highlighting the unique contribution of our business networks in explaining lead-lag effects across the globe.

5.2. Historical business descriptions

Using contemporary descriptions to analyze historical economic relationships can lead to an upward bias due to the inclusion of future information. To address this bias, updating our networks using historical business descriptions would be ideal. While we cannot retrieve historical descriptions for our entire dataset, we obtain descriptions for a large share of firms from SDC Platinum. This platform provides business descriptions for acquiring and targeted companies involved in M&A deals.

We obtain yearly updated business networks on the basis of historical business descriptions as follows. First, we pinpoint all deals involving publicly traded firms before January 1st, 1996. Subsequently, we select those firms that exhibited active trading in 1995. From there, we extract the latest description for each firm from all business descriptions attributed to M&A

deals prior to 1996.¹⁴ We then duplicate the procedure detailed in Section 2.1 to generate the business networks employed for portfolio creation in 1996. This exact process is perpetuated annually for the period spanning from 1997 through to 2021.

To compare the performance of portfolios based on recent and historical business descriptions, we limit the networks to stocks available in both sets of business descriptions.

[Table 10 about here.]

According to Panel A in Table 10, which represents investments from 1996 until 2021, a US portfolio based on US peers identified via recent descriptions generates a 107 bps monthly seven-factor alpha, while a portfolio based on historical descriptions generates a 90 bps alpha. The effect persists if we construct US portfolios using global peers (146 bps vs. 135 bps). However, these differences are not statistically significant.

The same trend is evident internationally. Global portfolios, based on historical descriptions, generate up to 212 bps monthly seven-factor alpha which is 16 bps lower than the alpha obtained from investments according to business descriptions as of 2022. This difference is not statistically significant

To ensure that our results are not biased by the first years in our evaluation periods, where less historical descriptions might be available¹⁵, we also investigate shorter evaluation periods from 2000 (2005, 2010) until 2021. Here, we observe even smaller differences in alpha. This is in line with the idea that the look-ahead bias should be more pronounced for earlier evaluation periods. Furthermore, we find compelling evidence supporting the persistence of the lead-lag effect in recent years. Specifically, we observe highly significant monthly seven-factor alphas, reaching up to 174 bps, throughout the period from 2010 to 2021.

Overall, we may conclude that a potential look-ahead bias introduced by using recent description may only marginally explain the existence of the lead-lag effect.

5.3. *Look-ahead bias of language models?*

Our analysis might contain a subtle forward-looking bias due to the language models we employ. To illustrate this, consider the term *COVID-19*. Until the virus's widespread emergence

¹⁴Similar to our dataset of business descriptions obtained from Refinitiv Workspace, we exclude descriptions consisting of less than ten words.

¹⁵Our dataset contains 11274 (18341, 29457, 47800) business descriptions allocated to M&A deals that occurred prior to 1996 (2000, 2005, 2010).

in late 2019, the term *COVID-19* was obscure and context-less. As a result, language models trained up until 2019 fail to recognize this term as referencing a global pandemic. Conversely, models trained on more current data are equipped with this understanding. For instance, OpenAI’s GPT-3 model is trained on vast amounts of text from recent decades up to September 2021 and thus should be able to interpret this term.

This hypothesis can be examined by comparing the cosine similarity between the embeddings of the terms *COVID-19* and *pandemic* using OpenAI’s embedding model (*text-embedding-ada-002*) and a Sentence Transformer model proposed by Reimers and Gurevych (2019) trained on data prior to the outbreak of the *COVID-19* pandemic (*nli-bert-base*). With the Sentence Transformer model, a relatively low cosine similarity of 0.45 is observed. In contrast, using OpenAI’s embedding model, we obtain a markedly higher cosine similarity of 0.89. We attribute these results to OpenAI’s model recognizing the frequent co-occurrence of *COVID-19* and *pandemic* in the training data, thereby indicating a high semantic similarity. Therefore, using a contemporary language model to backtest the performance of an investment strategy could potentially induce a look-ahead bias, as the model is privy to events that had not yet occurred.

We expect the relevance of this bias to differ across different research questions. While Lopez-Lira and Tang (2023) suspect a forward-looking bias of language models in the context of sentiment prediction and therefore restrict their sample period to October 2021 until December 2022, we believe this bias might be less severe in our setting. The reason is that the business descriptions under study typically do not reference specific events, and we only compute similarities instead of predicting future developments.

Nevertheless, we control for a potential look-ahead bias by leveraging differences in the training data of different models. An ideal setting would involve using a language model trained exclusively on data available before our evaluation period. To the best of our knowledge, there exists no *LLM* that is solely trained on data prior to our sample period from 1996 to 2021. At the same time, given the considerable volume of data and computational resources needed to train *LLMs*, we lack the resources to train such a model ourselves.

Despite these circumstances, we can leverage the discrepancies in the knowledge foundations of the above-mentioned language models to test if our networks suffer from a forward-looking bias. If this is the case, we would expect significant differences in the networks when analyzing firms heavily affected by major events in 2020, in contrast to others. We suspect that the *COVID-19* pandemic had a substantial impact on the *Pharmaceutical and Biotechnology* sector,

mainly due to the quest for a *COVID-19* vaccine. Meanwhile, *Gas and Oil producers* and firms operating in the *Mining* sector should have been less affected. Thus, we expect a lower correlation between the cosine similarities of the two different models among stocks within the *Pharmaceutical and Biotechnology* sector than in most other sectors.

[Table 11 about here.]

According to Table 11, we find that the correlation between the cosine similarities of the Sentence Transformer and the OpenAI network is not significantly smaller for stocks of the *Pharmaceuticals and Biotechnology* sector compared to other sectors. With a correlation of 0.59, the *Pharmaceuticals and Biotechnology* sector is among the top five sectors with the highest within-sector correlations of the two language models. While the correlation is indeed higher for the *Mining* sector, a vast number of sectors that should be less affected by *COVID-19*, for instance, the sectors *Industrial Engineering* and *Media*, show a lower correlation. These findings align with our initial hypothesis that a potential look-ahead bias introduced by a language model is negligible in our setting.

6. Second Potential Application: Identifying M&A targets using Business Networks

Our business networks provide a potential opportunity to shed light on the criteria firms consider when selecting takeover targets. We therefore gather data on public firms who acquired publicly traded companies from 2000 until 2022 from Refinitiv SDC Platinum. We then calculate how many target firms may be identified with the different networks.¹⁶

[Figure 3 about here.]

According to Figure 3, up to 9% of the acquired firms had the most similar business description according to the *Global T5-XXL* network, in contrast to less than 2% for global word-based networks. This difference is significant at the 5% level.

This disparity remains consistent even when considering a larger set of relations. For instance, while approximately 42% of target companies rank among the top 100 peers as per *Global*

¹⁶Note that we use the networks based on the most recent business descriptions. Using the descriptions obtained from SDC Platinum instead might introduce a selection bias, given that we only have access to descriptions of firms who engage in M&A activities.

T5-XXL network, fewer than 8% show up in global word-based networks (*Global BOW*). We also conduct a bootstrapping exercise, wherein we randomly select 100 stocks from the same four-digit SIC code. Our findings indicate that, on average, only about 10% of the target companies are identifiable through this method, suggesting that our networks are more accurate than narrow industry classifications.

Finally, we also calculate how many of the target firms may be identified when restricting our business networks to domestic peers (*Domestic T5-XXL*). We find that roughly 20% of the target firms are included in the 100 highest ranked domestic peers. This 20 percentage points lower recall score highlights the importance of considering networks that span the globe rather than domestic ones.

As a next step, we establish a logistic regression framework to systematically account for the possibility that companies are more inclined to acquire firms in the same industry, country, or with lower market capitalization. The dependent variable is a dummy variable "acquired," set at one if a firm pair represents a M&A deal and 0 otherwise. We compare the average cosine similarities between firm pairs in both groups. Suppose we discover that the cosine similarity of two business descriptions remains significant after controlling for other factors. In that case, we can conclude that business descriptions aid in predicting which firms are most likely purchased.

[Table 12 about here.]

Table 12 confirms that firms operating in the same industry and country are more likely to be targeted by acquiring firms and that smaller firms are more likely to be acquired. We also find evidence suggesting that acquiring firms tend to purchase firms with higher debt ratios, greater profitability, and more cash on hand. The odds ratios are highly significant, with t-statistics up to 35.88.

However, even after controlling for these factors, we find that the similarity between the business descriptions of the two firms is higher if a firm pair represents a real M&A deal. This effect is statistically significant, with a t-statistic of around 11, and is robust to different sample sizes of non-deal firm pairs. If we collect 100 randomly drawn firms for every deal, we observe a substantial increase in the Pseudo-R², from about 0.420 to 0.511. This increase is only slightly lower for larger deal/non-deal ratios. We thus argue that our networks could provide new insights concerning the global M&A literature.

7. Network Modification for further research

This paper introduces business networks as a tool to uncover a company’s economic relationships. Unlike other networks, ours can reveal potential competitors, suppliers, and customers. Depending on the research question, scholars may focus on specific relationships where only certain types of relationships are important. For example, researchers are not interested in supplier or customer links when assessing firm-specific competition intensity. To remove these links from our networks, researchers can first identify suppliers by searching for the term ”supplier” in the business descriptions of all linked firms, as firms that act as suppliers tend to include this keyword in their descriptions. Of course, this strategy only applies to firms that are not suppliers themselves, as in this case, the identified firms are potential competitors. However, this filtering strategy should produce solid results for all other firms.

To filter out firm-specific links with customers, researchers can identify the industries where customers typically operate. For example, a car dealer would typically be considered a customer of *Ford* and classify as a ”General Retailer”. Researchers can therefore define all related firms from this industry as customers of *Ford*. The remaining links can then be interpreted as competition links. However, to narrow down the list of competitors, researchers can further condition on industry membership (e.g., same four, three, or two-digit SIC) or add a size-difference threshold. For instance, researchers could only treat firms with a similar market capitalization as close competitors.

If we apply these filters on the economically linked firms of *Ford* (using the BN_{T5-XXL}), we identify six competitors, two suppliers, and nine customers. The name of the firms and the corresponding short business descriptions may be found in Table 13.

[Table 13 about here.]

Overall, we find evidence suggesting that the proposed filters can help to differentiate between competitor, supplier, and customer links. While all firms classified as competitors are car manufacturers, the two identified suppliers indeed operate as such. We further observe that most identified customers are car dealers which sounds reasonable.

It is worth noting that it may be possible to differentiate between competitor, supplier, and customer relations more effectively by implementing more advanced filtering methods or manual oversight. Future research may explore the development of a machine learning classifier to improve the accuracy of identifying and distinguishing these types of links.

8. Conclusion

This study introduces a novel approach to identifying economic links through textual data. Our method applies advanced context-aware natural language processing techniques to over 79,000 business descriptions of publicly traded stocks. By doing so, we identify business networks that model global economic links.

Instead of extracting business information from international annual reports, which is error-prone, we construct our networks based on business descriptions obtained from Refinitiv Workspace. Rather than using traditional word-based methods to identify similarities in business descriptions, we vote in favor of the application of large language models to deal with the limited number of business-related words present.

We evaluate our business networks (*BNs*) in various dimensions. We find that *LLM* business networks contain a substantially higher share of disclosed competitors than word-based networks. Furthermore, we showcase the usability of the *BNs* by investigating two potential applications.

First, we construct calendar time portfolios to capture lead-lag effects in the US and globally. We find that *LLM* networks outperform industry networks. Furthermore, we find evidence suggesting that investors should consider international in addition to domestic economic links. Globally, we obtain seven-factor alphas of up to 125 bps for investments into stocks above the fifth NYSE size decile. We further compare our business networks with networks based on shared analyst coverage, similar firm characteristics, traditional industry classifications, and stock return correlation within industries. Our results suggest that our *BNs* contain relations that these methods may not discover. We also run tests to control for potential look-ahead biases caused by using recent descriptions and applying contemporary language models. Our findings reveal no indications of a significant look-ahead bias within our analyses.

Second, we study M&A deals and find that target firms show up disproportionately often in our business networks. Industry membership alone may not explain this finding. We, therefore, run a logistic regression and find that firms are more likely to acquire firms with a higher business description similarity, controlling for factors like industry, country, size, and other fundamentals.

Given their global scale and broad coverage, our networks can reveal information on global economic links and provide the foundation for more accurate, firm-specific controls for competitor, supplier, and customer performance.

- Aleph-Alpha, 2023. Luminous performance benchmarks.
- Ali, U., Hirshleifer, D., 2020. Shared analyst coverage: Unifying momentum spillover effects. *Journal of Financial Economics* 136 (3), 649–675.
- Breitung, C., Müller, S., 2022. When firms open up: Identifying value relevant textual disclosure using simbert. Available at SSRN.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33, 1877–1901.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., Nori, H., Palangi, H., Ribeiro, M. T., Zhang, Y., 2023. Sparks of artificial general intelligence: Early experiments with gpt-4.
- Cohen, L., Frazzini, A., 2008. Economic links and predictable returns. *The Journal of Finance* 63 (4), 1977–2011.
- Cohen, L., Lou, D., 2012. Complicated firms. *Journal of financial economics* 104 (2), 383–400.
- Cohen, L., Malloy, C., Nguyen, Q., 2020. Lazy prices. *The Journal of Finance* 75 (3), 1371–1415.
- Conneau, A., Kiela, D., 2018. Senteval: An evaluation toolkit for universal sentence representations. arXiv preprint arXiv:1803.05449.
- De Franco, G., Kothari, S. P., Verdi, R. S., 2011. The benefits of financial statement comparability. *Journal of Accounting research* 49 (4), 895–931.
- Devlin, J., Chang, M.-W., Lee, K., Toutanova, K., 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
- Eisdorfer, A., Froot, K., Ozik, G., Sadka, R., 12 2021. Competition Links and Stock Returns. *The Review of Financial Studies*Hhab133.
- Fama, E. F., French, K. R., 2015. A five-factor asset pricing model. *Journal of financial economics* 116 (1), 1–22.
- Fama, E. F., MacBeth, J. D., 1973. Risk, return, and equilibrium: Empirical tests. *Journal of Political Economy* 81, 607–636.

- Gatev, E., Goetzmann, W. N., Rouwenhorst, K. G., 2006. Pairs trading: Performance of a relative-value arbitrage rule. *The Review of Financial Studies* 19 (3), 797–827.
- Grüner, A., Finke, C., et al., 2018. Lead-lag relationships in international stock markets revisited: Are they exploitable? *International Journal of Financial Research* 9 (1), 8–30.
- Guo, F., Liu, T., Tu, D., 2023. Neglected peers in merger valuations. *The Review of Financial Studies*, hhad004.
- Hoberg, G., Phillips, G., 2010. Product market synergies and competition in mergers and acquisitions: A text-based analysis. *The Review of Financial Studies* 23 (10), 3773–3811.
- Hoberg, G., Phillips, G., 2016. Text-based network industries and endogenous product differentiation. *Journal of Political Economy* 124 (5), 1423–1465.
- Hoberg, G., Phillips, G. M., 2018. Text-based industry momentum. *Journal of Financial and Quantitative Analysis* 53 (6), 2355–2388.
- Hou, K., 2007. Industry information diffusion and the lead-lag effect in stock returns 20, 1113–1138.
- Huang, X., 2015. Thinking outside the borders: Investors’ underreaction to foreign operations information 28, 3109–3152.
- Huber, D., Jacobs, H., Müller, S., Preissler, F., 2023. International factor models. *Journal of Banking and Finance* forthcoming, 1–16.
- Ibriyamova, F., Kogan, S., Salganik-Shoshan, G., Stolin, D., 2019. Predicting stock return correlations with brief company descriptions. *Applied Economics* 51 (1), 88–102.
- Lee, C. M., Sun, S. T., Wang, R., Zhang, R., 2019. Technological links and predictable returns. *Journal of Financial Economics* 132 (3), 76–96.
- Li, F., Lundholm, R., Minnis, M., 2013. A measure of competition based on 10-k filings. *Journal of Accounting Research* 51 (2), 399–436.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V., 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.

- Lopez-Lira, A., Tang, Y., 2023. Can chatgpt forecast stock price movements? return predictability and large language models. Available at SSRN 4412788.
- Menzly, L., Ozbas, O., 2010. Market segmentation and cross-predictability of returns. *The Journal of Finance* 65 (4), 1555–1580.
- Moskowitz, T. J., Grinblatt, M., 1999. Do industries explain momentum? *The Journal of finance* 54 (4), 1249–1290.
- Müller, S., 2019. Economic links and cross-predictability of stock returns: Evidence from characteristic-based styles. *Review of Finance* 23 (2), 363–395.
- Ni, J., Ábrego, G. H., Constant, N., Ma, J., Hall, K. B., Cer, D., Yang, Y., 2021. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. arXiv preprint arXiv:2108.08877.
- OpenAI, 2023. Gpt-4 technical report.
- Parsons, C. A., Sabbatucci, R., Titman, S., 2020. Geographic lead-lag effects. *The Review of Financial Studies* 33 (10), 4721–4770.
- Reimers, N., Gurevych, I., 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need. *Advances in neural information processing systems* 30.

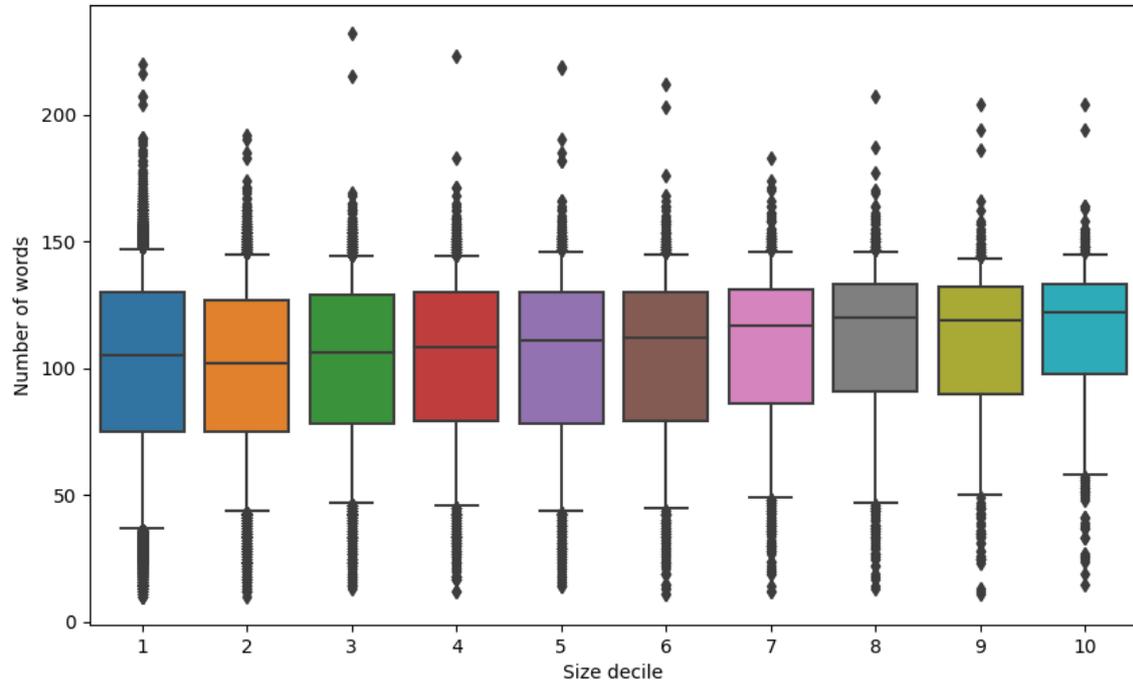


Figure 1: Distribution of the number of words in a business description across different NYSE size deciles.

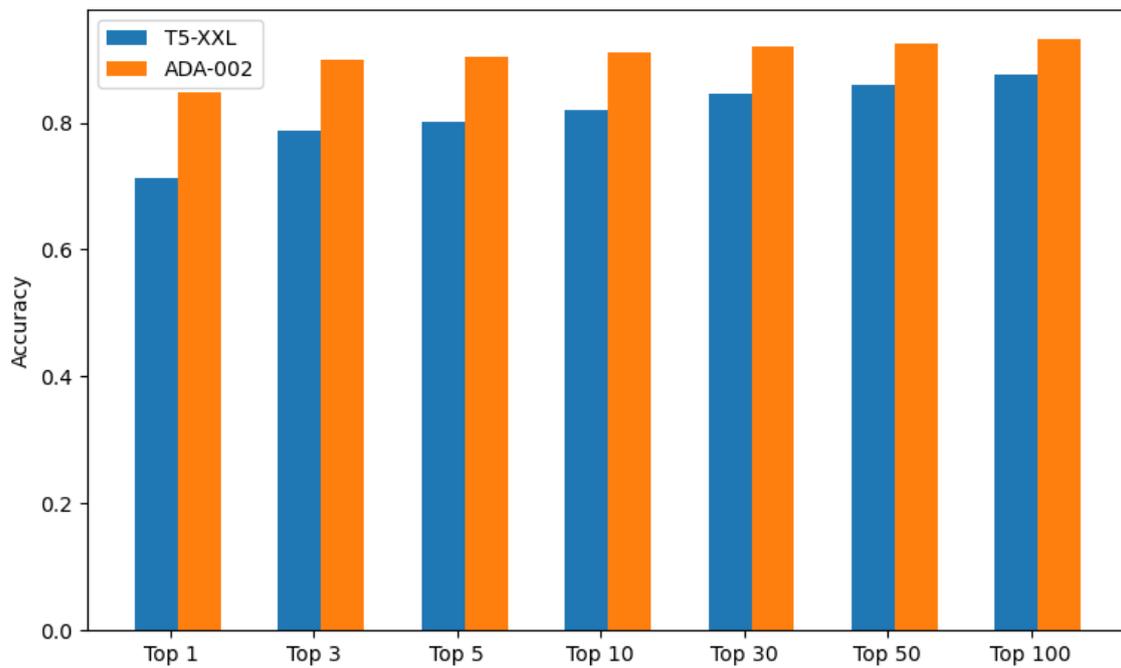


Figure 2: Proportion of US firms whose business description is ranked among the most similar descriptions from Refinitiv to the Item 1 section of the same firm. *Top 1* indicates how often the correct firm is ranked as most similar. *Top 3* shows how often the correct firm is ranked among the most similar firms. The same logic applies to the other groups.

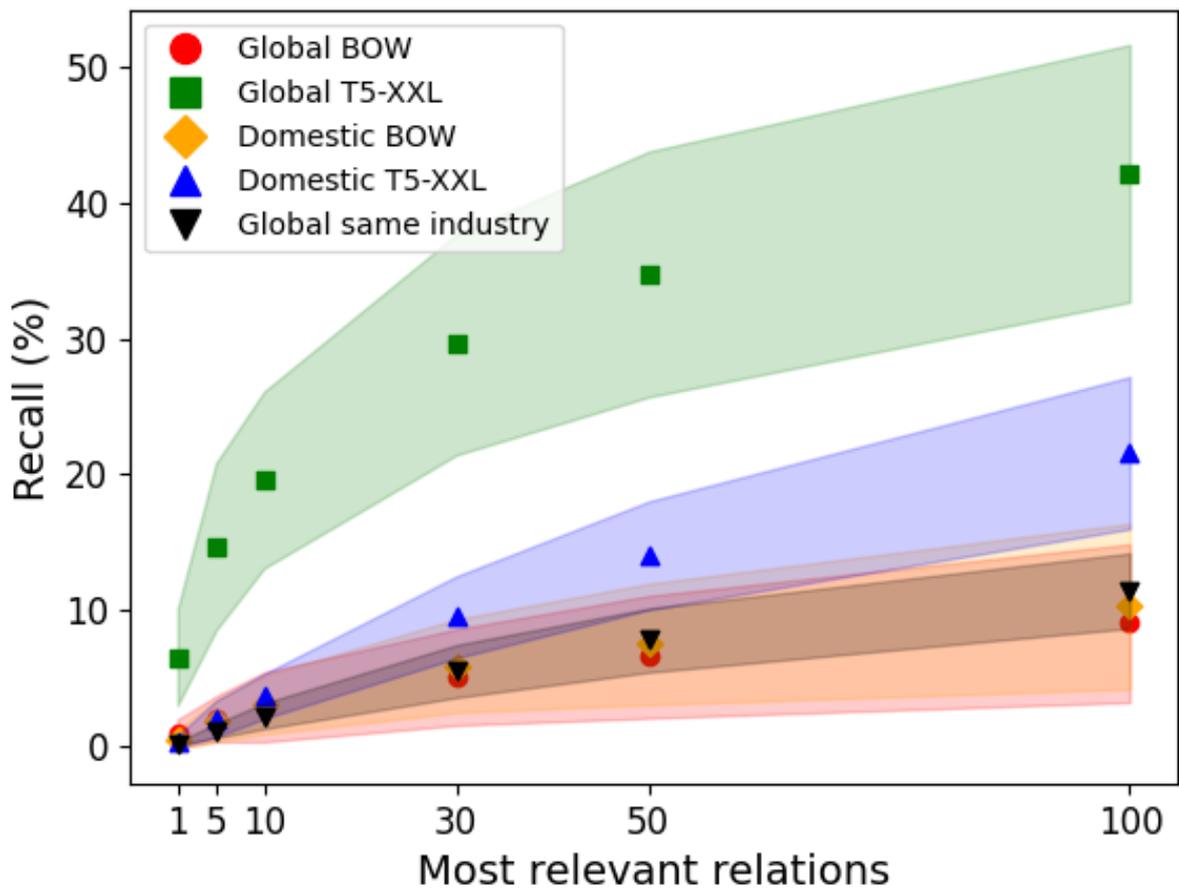


Figure 3: This figure examines the proportion of M&A target firms that are members of various business networks. Additionally, we randomly select stocks from the same industry as a control group (Random same industry). The shadow in our visualizations represents the 95% confidence intervals of the data.

Table 1: **Number of business descriptions by country**

| Country | #Firms | Country | #Firms | Country | #Firms |
|----------------|--------|-------------|--------|----------------|--------|
| Argentina | 125 | India | 3897 | Portugal | 141 |
| Australia | 3225 | Indonesia | 826 | Qatar | 52 |
| Austria | 187 | Ireland | 156 | Romania | 163 |
| Bahrain | 47 | Italy | 722 | Russia | 604 |
| Bangladesh | 147 | Japan | 5767 | Serbia | 117 |
| Belgium | 264 | Jordan | 254 | Singapore | 1106 |
| Brazil | 371 | Kazakhstan | 65 | Slovenia | 37 |
| Bulgaria | 307 | Kenya | 60 | South Africa | 862 |
| Canada | 5246 | Korea | 3107 | Spain | 417 |
| Chile | 275 | Kuwait | 221 | Sri Lanka | 263 |
| China | 4637 | Lithuania | 51 | Sweden | 1247 |
| Colombia | 84 | Malaysia | 1408 | Switzerland | 427 |
| Croatia | 131 | Mauritius | 83 | Taiwan | 2475 |
| Czech Republic | 93 | Mexico | 242 | Thailand | 1006 |
| Denmark | 370 | Morocco | 88 | Tunisia | 77 |
| Egypt | 239 | Netherland | 322 | Turkey | 471 |
| Estonia | 23 | New Zealand | 266 | USA | 20835 |
| Finland | 269 | Nigeria | 161 | Ukraine | 123 |
| France | 1782 | Norway | 674 | United Kingdom | 4698 |
| Germany | 1701 | Oman | 135 | Vietnam | 1307 |
| Greece | 397 | Pakistan | 407 | | |
| Hong Kong | 2521 | Peru | 155 | | |
| Hungary | 78 | Philippines | 313 | | |

This table presents the number of available business descriptions that contain at least 10 words on the country-level.

Table 2: **Recent and historical business descriptions of Ford and General Motors**

| Ford | |
|----------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| August 2022 (Workspace) | Ford Motor Company is an automobile company that designs, manufactures, markets, and services a full line of Ford trucks, utility vehicles, cars as well as Lincoln luxury vehicles. The Company operates in three segments: Automotive, Mobility and Ford Credit. The Automotive segment is engaged in developing, manufacturing, distributing, and servicing the vehicles, parts and accessories of Ford and Lincoln vehicles. The Mobility segment primarily includes the development of Ford’s autonomous vehicles and related businesses. The Company also holds ownership is Argo AI, which is a developer of autonomous driving systems, and Spin, which is a micro-mobility service provider. The Ford Credit segment is comprised of the Ford Credit business on a consolidated basis, which is primarily vehicle-related financing and leasing activities. Ford Credit offers a wide variety of automotive financing products to and through automotive dealers throughout the world. |
| 1983-07-25 (SDC) | Ford Motor Co, located in Dearborn, Michigan, manufactures and wholesales automobiles, trucks, automobile parts, industrial trucks and tractors. The company also provides auto-financing services. It was founded in 1903. |
| General Motors | |
| August 2022 (Workspace) | General Motors Company designs, builds and sells trucks, crossovers, cars and automobile parts and provides software-enabled services and subscriptions worldwide. The Company provides automotive financing services through its General Motors Financial Company, Inc. (GM Financial) segment. GM North America (GMNA) and GM International (GMI) develops, manufactures and/or markets vehicles under the Buick, Cadillac, Chevrolet and GMC brands. The Company’s segments include GMNA, GMI, Cruise and GM Financial. Its Cruise segment is engaged in the development and commercialization of autonomous vehicle technology. It offers OnStar and connected services to approximately 22 million connected vehicles globally through subscription-based and complimentary services. It is also developing hydrogen fuel cell applications across transportation and industries, including mobile power generation, class seven/eight truck, locomotive, aerospace and marine applications. |
| 2012-02-28 (SDC) | General Motors Co, located in Detroit, Michigan, manufactures and wholesales trucks, crossovers, cars and automobile parts. It also provides automotive financing services through General Motors Financial Co Inc (GM Financial). GM North America (GMNA) and GM International (GMI) are its automotive segments. GMNA and GMI are meeting the demands of customers with vehicles developed, manufactured and/or marketed under the Buick, Cadillac, Chevrolet and GMC and Holden brands. Its brands offer luxury cars, crossovers, sport utility vehicles (SUVs) and sedans. Its Car-and Ride-Sharing Maven is a shared vehicle marketplace. Through its subsidiary, OnStar LLC (OnStar), it provides connected safety, security and mobility solutions for retail and fleet customers. GM Cruise is its global segment engaged in the development and commercialization of autonomous vehicle technology. It is also a holding company. The Company was founded on September 16, 1908. |
| 2009-09-16 (SDC) | General Motors Co, headquartered in Detroit, Michigan, manufactures and wholesales cars and trucks. The company and its strategic partners manufacture cars and trucks in 31 countries and sell through its brands namely, Buick, Cadillac, Chevrolet, GMC, Daewoo, Holden, Jiefang, Opel, Vauxhall and Wuling. It operates in 157 countries including the United States, China, Brazil, Germany, the United Kingdom, Canada, and Italy. The Company was founded in 1908. |
| 1984-04-15 (SDC) | General Motors Corp, located in Detroit, Michigan, manufactures motor vehicles, related parts, defense and space products, business information and telecommunication systems, locomotives, satellites, test equipment and marine engines. The company also provides financing services include consumer vehicle financing, full-service leasing and fleet leasing, dealer financing and car and truck extended service contracts, residential and commercial mortgage services, commercial and vehicle insurance and asset-based lending. In addition the company offers insurance services including automobile and homeowners insurance, automobile mechanical protection, reinsurance and commercial insurance. The company operates in the US, Canada, Mexico, Europe, Asia Pacific and Latin America and was founded in September 16, 1908. |

This table contains unique historical business descriptions of *Ford* and *General Motors* obtained from Refinitiv Workspace and SDC Platinum. We further denote the announcement date of the deal where the description was first mentioned.

Table 3: Competitors of Ford according to different networks in 2021

| | Name | Country | Sector | NYSE |
|---------------------|----------------|----------------|------------------------|------|
| BOW: 34 Peers | | | | |
| 1 | BMW | Germany | Automobiles and Parts | 10.0 |
| 2 | TOYOTA MOTOR | Japan | Automobiles and Parts | 10.0 |
| 3 | HYUNDAI MOTOR | Korea | Automobiles and Parts | 10.0 |
| 4 | VOLKSWAGEN | Germany | Automobiles and Parts | 10.0 |
| 5 | BP | United Kingdom | Oil and Gas Producers | 10.0 |
| T5-XXL: 26 Peers | | | | |
| 1 | GENERAL MOTORS | USA | Automobiles and Parts | 10.0 |
| 2 | TOYOTA MOTOR | Japan | Automobiles and Parts | 10.0 |
| 3 | HYUNDAI MOTOR | Korea | Automobiles and Parts | 10.0 |
| 4 | DAIMLER | Germany | Automobiles and Parts | 10.0 |
| 5 | TOYOTA INDS. | Japan | Automobiles and Parts | 9.0 |
| ADA-002: 70 Peers | | | | |
| 1 | GENERAL MOTORS | USA | Automobiles and Parts | 10.0 |
| 2 | TESLA | USA | Automobiles and Parts | 10.0 |
| 3 | TOYOTA MOTOR | Japan | Automobiles and Parts | 10.0 |
| 4 | HYUNDAI MOTOR | Korea | Automobiles and Parts | 10.0 |
| 5 | DAIMLER | Germany | Automobiles and Parts | 10.0 |
| LUMINOUS: 278 Peers | | | | |
| 1 | GENERAL MOTORS | USA | Automobiles and Parts | 10.0 |
| 2 | TESLA | USA | Automobiles and Parts | 10.0 |
| 3 | VOLKSWAGEN | Germany | Automobiles and Parts | 10.0 |
| 4 | TOYOTA MOTOR | Japan | Automobiles and Parts | 10.0 |
| 5 | DAIMLER | Germany | Automobiles and Parts | 10.0 |
| TNIC: 10 Peers | | | | |
| 1 | GENERAL MOTORS | USA | Automobiles and Parts | 10.0 |
| 2 | TESLA | USA | Automobiles and Parts | 10.0 |
| 3 | LEAR | USA | Automobiles and Parts | 7.0 |
| 4 | LKQ | USA | Automobiles and Parts | 8.0 |
| 5 | PACCAR | USA | Industrial Engineering | 9.0 |

This table presents the most relevant competitors of the car manufacturer *Ford*, as determined by various networks. Firms allocated to the same NYSE size decile are ranked highest, assuming that the closest competitors should be of similar firm size. We further prioritize domestic peers among those, as competition intensity tends to be most severe domestically. In total, we evaluate a word-based network *BOW*, a Sentence Transformer network (*T5-XXL*), a network based on the OpenAI model (*ADA-002*), and a network based on a model by Aleph Alpha (*LUMINOUS*). We also provide the most similar firms according to the TNIC dataset ([Hoberg and Phillips, 2010, 2016](#)) as of 2021.

Table 4: **Analysis of US 2021 Business Networks: Relations, Overlaps, and Correlations**

| | BOW | T5-XXL | ADA-002 | LUMIN | TNIC | SIC4 |
|---------------------------------|--------|--------|---------|--------|--------|--------|
| Panel A: Summary | | | | | | |
| Mean Pred. | 179 | 351 | 222 | 339 | 78 | 253 |
| Median. Pred. | 109 | 255 | 145 | 204 | 18 | 141 |
| Same SIC4 (%) | 10.47 | 18.80 | 18.72 | 17.46 | 25.92 | 100 |
| Same SIC-3 (%) | 20.82 | 35.94 | 35.35 | 38.15 | 46.19 | 100 |
| Same SIC-2 (%) | 30.35 | 49.59 | 48.91 | 53.47 | 61.10 | 100 |
| Same Country (%) | 19.52 | 28.10 | 27.28 | 23.16 | 100 | 14.45 |
| Same NYSE decile (%) | 24.03 | 26.73 | 26.87 | 24.50 | 23.50 | 27.53 |
| Panel B: Overlap (%) | | | | | | |
| BOW | 100.00 | 22.34 | 26.68 | 26.49 | 22.92 | 20.51 |
| T5-XXL | 63.47 | 100.00 | 85.29 | 79.10 | 69.67 | 59.77 |
| ADA-002 | 46.59 | 52.42 | 100.00 | 58.14 | 46.89 | 40.59 |
| LUMIN | 60.55 | 63.62 | 76.09 | 100.00 | 57.36 | 49.29 |
| TNIC2021 | 50.84 | 54.38 | 59.55 | 55.67 | 100.00 | 52.43 |
| SIC4 | 22.82 | 23.41 | 25.86 | 23.99 | 26.30 | 100.00 |
| Panel C: Return Correlation (%) | | | | | | |
| BOW | 100.00 | 58.41 | 57.44 | 60.03 | 49.20 | 47.20 |
| T5-XXL | | 100.00 | 81.03 | 82.81 | 65.62 | 63.41 |
| ADA-002 | | | 100.00 | 81.52 | 64.09 | 60.40 |
| LUMIN | | | | 100.00 | 65.79 | 61.96 |
| TNIC2021 | | | | | 100.00 | 53.43 |
| SIC4 | | | | | | 100.00 |

This table provides various analyses of the similarity of the different networks restricted to US firms and peers. Panel A provides insights on the average number of relations per firm and the share of ties within the same industry, country, and size. Panel B shows the pairwise overlaps across the networks. Panel C displays a correlation matrix of average returns of related US firms, as identified by different business networks. Our analysis includes the TNIC dataset restricted to the fiscal year 2021 (TNIC2021), a word-based network (BOW), a Sentence Transformer model (T5-XXL), a model provided by OpenAI (ADA-002) and a model supplied by Aleph Alpha (LUMIN). We further construct and evaluate networks based on four-digit SIC industry classifications.

Table 5: **Detection rate of disclosed US competitors**

| | BOW | T5-XXL | ADA-002 | LUMIN | TNIC2021 | SIC4 |
|------------------------------------------------|-------|--------|---------|-------|----------|-------|
| Panel A: Item 1 (%) | | | | | | |
| Recall 1 | 1.56 | 3.52 | 3.48 | 4.20 | 8.96 | 2.64 |
| Recall 5 | 4.80 | 11.08 | 12.59 | 12.87 | 22.07 | 8.44 |
| Recall 10 | 7.32 | 17.03 | 19.59 | 19.35 | 28.19 | 11.92 |
| Recall 30 | 12.08 | 30.51 | 32.63 | 35.15 | 36.31 | 16.31 |
| Recall 50 | 13.31 | 36.39 | 36.39 | 41.94 | 38.02 | 17.79 |
| Recall 100 | 15.07 | 43.66 | 40.34 | 49.82 | 40.46 | 20.27 |
| Recall Total | 15.47 | 52.30 | 43.98 | 53.98 | 43.98 | 21.19 |
| Precision Total | 0.42 | 0.53 | 0.72 | 0.68 | 0.57 | 0.52 |
| Panel B: Comparable Company Analysis (M&A) (%) | | | | | | |
| Recall 1 | 1.25 | 1.70 | 1.36 | 1.82 | 2.95 | 1.36 |
| Recall 5 | 3.75 | 5.79 | 6.24 | 6.70 | 9.31 | 5.79 |
| Recall 10 | 6.58 | 10.56 | 9.76 | 12.03 | 12.71 | 8.40 |
| Recall 30 | 12.37 | 19.18 | 16.46 | 20.66 | 20.54 | 15.10 |
| Recall 50 | 16.69 | 24.63 | 21.34 | 26.33 | 25.43 | 17.37 |
| Recall 100 | 19.64 | 35.53 | 28.38 | 34.39 | 33.03 | 21.79 |
| Recall Total | 22.47 | 55.85 | 41.09 | 52.89 | 50.06 | 23.61 |
| Precision Total | 0.61 | 0.57 | 0.68 | 0.66 | 0.65 | 0.58 |

This table evaluates the accuracy of the networks. Panel A provides an overview of how many US competitors disclosed in the competition subsection of Item 1 in financial statements may be identified by considering the top 1 (5, 10, 30) or the entire set of sufficiently related firms (recall score) according to the different networks. The recall score is calculated by dividing the number of firms in a network by the total number of disclosed competitors. The precision score is calculated by relating the number of correctly identified competitors to the total number of identified peers. In Panel B, we use the same metrics to identify competitors discovered by investment banks during the valuation of a US target firm, a process called comparable company analysis (CCA). We use several models, including the TNIC dataset restricted to the fiscal year 2021 (*TNIC2021*), a word-based network (*BOW*), a Sentence Transformer model (*T5-XXL*), a model provided by OpenAI (*ADA-002*) and a model supplied by Aleph Alpha (*LUMIN*). We further construct and evaluate networks based on four-digit SIC industry classifications.

Table 6: US Business Network spillover effect

| | TNIC | TNIC2021 | BOW | T5-XXL | ADA-002 | LUMIN | SIC4 | $SIC4_{NN}$ |
|---------------------------------------------|-----------------------|----------------------|----------------------|-----------------------|----------------------|-----------------------|----------------------|----------------------|
| Panel A: US portfolio based on US peers | | | | | | | | |
| MKTRF | 7.76 (1.44) | 0.06 (0.01) | 5.12** (2.17) | 6.53** (2.15) | 5.07* (1.79) | 4.98 (1.58) | 5.24* (1.95) | 4.58** (2.16) |
| SMB | 11.66 (1.54) | 13.51** (2.19) | 7.11** (2.2) | 8.29* (1.94) | 8.12** (2.0) | 7.92* (1.8) | 4.48 (1.26) | 1.71 (0.63) |
| HML | -3.04 (-0.3) | -1.63 (-0.21) | -1.47 (-0.33) | -2.85 (-0.51) | -1.4 (-0.27) | -2.4 (-0.4) | -2.8 (-0.59) | 0.28 (0.07) |
| WML | -6.37 (-1.05) | -2.89 (-0.63) | -0.62 (-0.2) | -0.65 (-0.17) | -0.44 (-0.12) | -1.05 (-0.27) | -1.21 (-0.4) | 0.28 (0.14) |
| RMW | -0.84 (-0.07) | -4.72 (-0.44) | -1.01 (-0.19) | -1.78 (-0.25) | -3.29 (-0.51) | -2.39 (-0.33) | -1.34 (-0.22) | -4.87 (-1.02) |
| CMA | 20.56 (1.2) | 10.89 (0.75) | 10.85 (1.61) | 12.16 (1.35) | 12.07 (1.4) | 11.5 (1.24) | 8.68 (1.22) | 4.91 (0.9) |
| ST Reversal | -97.42*** (-13.64) | -79.35*** (-14.4) | -31.21*** (-7.54) | -42.93*** (-8.49) | -39.62*** (-8.12) | -43.28*** (-8.24) | -36.46*** (-8.85) | -26.71*** (-7.89) |
| Alpha | 1.67*** (8.81) | 1.03*** (6.34) | 0.63*** (7.81) | 0.95*** (8.85) | 0.89*** (8.68) | 1.03*** (9.41) | 0.79*** (8.63) | 0.53*** (7.31) |
| Panel B: US portfolio based on global peers | | | | | | | | |
| MKTRF | - | - | 9.55* (1.94) | 8.96 (1.64) | 5.95 (1.01) | 5.92 (1.09) | 5.62 (1.01) | 4.96 (1.11) |
| SMB | - | - | 11.1 (1.63) | 14.7** (2.05) | 16.42** (2.19) | 13.78** (1.97) | 11.82* (1.71) | 11.82** (2.06) |
| HML | - | - | -1.5 (-0.17) | -3.07 (-0.32) | -2.15 (-0.21) | -2.47 (-0.27) | -0.56 (-0.06) | 6.6 (0.94) |
| WML | - | - | 4.16 (0.7) | 3.03 (0.48) | 2.12 (0.32) | 2.2 (0.38) | 1.1 (0.18) | 0.11 (0.02) |
| RMW | - | - | 4.66 (0.38) | 2.92 (0.23) | 2.86 (0.21) | 1.76 (0.14) | 0.37 (0.03) | -4.36 (-0.46) |
| CMA | - | - | 23.32* (1.74) | 31.1* (1.88) | 29.31* (1.77) | 26.77* (1.79) | 21.55 (1.55) | 7.44 (0.7) |
| ST Reversal | - | - | -66.54*** (-7.92) | -87.54*** (-10.75) | -86.78*** (-9.84) | -87.03*** (-11.08) | -74.86*** (-9.06) | -59.84*** (-8.1) |
| Alpha | - | - | 0.96*** (6.21) | 1.47*** (8.13) | 1.44*** (7.83) | 1.49*** (8.53) | 1.25*** (7.3) | 0.84*** (6.06) |

We study the lead-lag relationship among US firms by constructing calendar-time portfolios that are rebalanced every month. We go long (short) in the 20% stocks whose most similar firms showed the best (worst) performance in the previous month. We restrict the networks to US peers in Panel A, whereas in Panel B, we consider the entire universe of global peers. We use several models, including the full TNIC dataset and a version that is restricted to the fiscal year 2021 (*TNIC2021*), a word-based network (*BOW*), a Sentence Transformer model (*T5-XXL*), a model provided by OpenAI (*ADA-002*) and a model provided by Aleph Alpha (*LUMIN*). We further construct and evaluate networks based on four-digit SIC industry classifications. Finally, we evaluate the portfolio constructed based on firms within the same four-digit SIC that are not included in any language model-based network (*SIC_{NN}*). We report seven-factor alphas (five-factor model plus momentum and short-term reversal). We denote the t-statistics of the coefficients in parentheses. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 7: Global Business network spillover effects

| | BOW | T5-XXL | ADA-002 | LUMIN | SIC4 | $SIC4_{NN}$ |
|-------------------------------------------------|----------------------|----------------------|----------------------|----------------------|-----------------------|---------------------|
| Panel A: Global portfolio based on US peers | | | | | | |
| MKTRF | 3.06 (1.34) | 4.6 (1.46) | 3.04 (0.98) | 2.55 (0.79) | 2.47 (0.97) | 2.03 (1.04) |
| SMB | 1.49 (0.26) | -1.3 (-0.17) | -0.25 (-0.03) | -0.72 (-0.09) | -2.33 (-0.36) | -4.48 (-0.95) |
| HML | 3.29 (0.43) | -0.78 (-0.08) | -0.99 (-0.11) | -2.59 (-0.27) | -0.15 (-0.02) | 4.25 (0.67) |
| WML | 0.69 (0.21) | 1.54 (0.38) | 1.07 (0.27) | 0.76 (0.18) | 0.28 (0.09) | 1.65 (0.68) |
| RMW | -6.36 (-1.2) | -7.72 (-1.13) | -7.53 (-1.19) | -7.65 (-1.14) | -4.9 (-0.84) | -8.16* (-1.93) |
| CMA | 7.31 (0.75) | 9.92 (0.8) | 11.0 (0.9) | 10.93 (0.83) | 0.22 (0.02) | -3.95 (-0.51) |
| ST Reversal | -31.97*** (-8.81) | -43.11*** (-9.54) | -39.37*** (-8.56) | -43.47*** (-9.49) | -37.05*** (-10.29) | -27.8*** (-9.26) |
| Alpha | 0.69*** (7.43) | 1.04*** (8.23) | 0.96*** (7.95) | 1.12*** (8.73) | 0.86*** (8.1) | 0.59*** (7.09) |
| Panel B: Global portfolio based on global peers | | | | | | |
| MKTRF | -1.31 (-0.3) | 0.75 (0.14) | -0.41 (-0.08) | 0.11 (0.02) | 4.05 (1.28) | 4.03 (1.54) |
| SMB | -12.15 (-1.01) | -15.22 (-1.0) | -18.26 (-1.15) | -15.26 (-1.08) | 1.59 (0.18) | 5.43 (0.84) |
| HML | 11.69 (0.73) | 16.14 (0.75) | 18.24 (0.85) | 14.29 (0.72) | 8.53 (0.8) | 7.18 (0.95) |
| WML | 6.88 (1.28) | 9.78 (1.51) | 10.28 (1.5) | 8.12 (1.3) | 2.92 (0.66) | 1.14 (0.34) |
| RMW | -18.78* (-1.78) | -24.26* (-1.72) | -26.03* (-1.78) | -22.67* (-1.72) | -7.67 (-1.05) | -4.79 (-0.94) |
| CMA | -24.52 (-1.16) | -28.72 (-1.02) | -31.2 (-1.07) | -22.04 (-0.84) | 1.49 (0.11) | 5.35 (0.64) |
| ST Reversal | -47.54*** (-8.14) | -59.57*** (-8.39) | -62.51*** (-8.59) | -59.05*** (-8.75) | -49.0*** (-10.12) | -38.09*** (-9.1) |
| Alpha | 1.73*** (9.11) | 2.34*** (9.6) | 2.35*** (9.49) | 2.25*** (9.98) | 1.42*** (10.79) | 0.75*** (7.43) |

We study the lead-lag relationship among global firms by constructing calendar-time portfolios that are rebalanced every month. We go long (short) in the 20% stocks whose most similar firms showed the best (worst) performance in the previous month. We restrict the networks to US peers in Panel A, whereas in Panel B, we consider the entire universe of global peers. We use several models, including the full TNIC dataset and a version that is restricted to the fiscal year 2021 (*TNIC2021*), a word-based network (*BOW*), a Sentence Transformer model (*T5-XXL*), a model provided by OpenAI (*ADA-002*) and a model provided by Aleph Alpha (*LUMIN*). We further construct and evaluate networks based on four-digit SIC industry classifications. Finally, we evaluate the portfolio constructed based on firms within the same four-digit SIC that are not included in any language model-based network (SIC_{NN}). We report seven-factor alphas (five-factor model plus momentum and short-term reversal). We denote the t-statistics of the coefficients in parentheses. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 8: **Abnormal return of Business Network portfolios**

| | TNIC | TNIC2021 | BOW | T5-XXL | ADA-002 | LUMIN | SIC4 | $SIC4_{NN}$ |
|-------------------------------------------------|-------------------|-------------------|-------------------|--------------------|--------------------|--------------------|-------------------|-------------------|
| Panel A: US portfolio based on US peers | | | | | | | | |
| < 5 th NYSE decile | 1.92*** (9.67) | 1.34*** (6.98) | 0.68*** (8.23) | 0.98*** (9.07) | 0.92*** (9.04) | 1.07*** (9.65) | 0.8*** (8.69) | 0.54*** (7.3) |
| ≥ 5 th NYSE decile | 0.88*** (4.39) | 0.3* (1.77) | 0.42*** (4.54) | 0.77*** (6.02) | 0.69*** (5.26) | 0.78*** (5.89) | 0.66*** (5.95) | 0.44*** (4.32) |
| Panel B: US portfolio based on global peers | | | | | | | | |
| < 5 th NYSE decile | - | - | 1.12*** (6.51) | 1.75*** (8.93) | 1.64*** (8.25) | 1.72*** (9.04) | 1.36*** (7.41) | 0.96*** (6.43) |
| ≥ 5 th NYSE decile | - | - | 0.42*** (2.92) | 0.71*** (3.79) | 0.78*** (4.19) | 0.78*** (4.29) | 0.72*** (3.97) | 0.46*** (3.12) |
| Panel C: Global portfolio based on US peers | | | | | | | | |
| < 5 th NYSE decile | - | - | 0.74*** (7.93) | 1.06*** (8.53) | 0.98*** (8.34) | 1.15*** (9.09) | 0.87*** (8.34) | 0.59*** (7.15) |
| ≥ 5 th NYSE decile | - | - | 0.48*** (4.51) | 0.87*** (5.92) | 0.82*** (5.46) | 0.9*** (5.9) | 0.71*** (5.54) | 0.51*** (4.17) |
| Panel D: Global portfolio based on global peers | | | | | | | | |
| < 5 th NYSE decile | - | - | 1.89*** (9.61) | 2.54*** (10.19) | 2.53*** (10.02) | 2.43*** (10.63) | 1.5*** (11.79) | 0.78*** (7.81) |
| ≥ 5 th NYSE decile | - | - | 0.88*** (4.41) | 1.2*** (4.47) | 1.24*** (4.66) | 1.25*** (4.92) | 0.95*** (5.55) | 0.56*** (4.36) |

We study the lead-lag relationship among smaller and larger stocks in the US and internationally by constructing calendar-time portfolios that are rebalanced every month. Smaller stocks are those below the fifth NYSE size decile, and larger ones in the fifth NYSE size decile and beyond. We go long (short) in the 20% stocks whose most similar firms showed the best (worst) performance in the previous month. We use several models, including the full TNIC dataset and a version that is restricted to the fiscal year 2021 (*TNIC2021*), a word-based network (*BOW*), a Sentence Transformer model (*T5-XXL*), a model provided by OpenAI (*ADA-002*) and a model provided by Aleph Alpha (*LUMIN*). We further construct and evaluate networks based on four-digit SIC industry classifications. Finally, we evaluate the portfolio constructed based on firms within the same four-digit SIC that are not included in any language model-based network (SIC_{NN}). We report seven-factor alphas (five-factor model plus momentum and short-term reversal). We denote the t-statistics of the coefficients in parentheses. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 9: **Fama MacBeth: Assessing the novelty of BN_{T5-XXL}**

| | (1) | (2) | (3) | (4) | (5) | (6) |
|---------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| | US | US | US | Global | Global | Global |
| BN_{T5-XXL} | 0.158*** (3.68) | 0.0970*** (2.98) | 0.0766** (2.42) | 0.277*** (5.75) | 0.187*** (5.03) | 0.137*** (3.74) |
| BN_{ANA} | | 0.146*** (3.80) | 0.125*** (3.83) | | 0.151*** (4.19) | 0.0828** (2.20) |
| BN_{corr} | | | 0.00672 (0.26) | | | 0.0332 (1.23) |
| BN_{SESM} | | | 0.0861** (2.42) | | | 0.149*** (5.15) |
| BN_{SIC} | | | 0.0129 (0.61) | | | 0.123*** (5.27) |
| <i>Return</i> | -0.219*** (-5.87) | -0.246*** (-7.10) | -0.259*** (-7.60) | -0.140*** (-4.36) | -0.193*** (-6.69) | -0.261*** (-9.19) |
| <i>SIZE</i> | -0.0445 (-1.06) | -0.0460 (-1.10) | -0.0512 (-1.22) | 0.00101 (0.04) | -0.0277 (-1.10) | -0.0350 (-1.22) |
| <i>Beta</i> | 0.0686 (1.11) | 0.0670 (1.11) | 0.0756 (1.31) | -0.0280 (-0.51) | -0.0351 (-0.67) | -0.0305 (-0.56) |
| <i>Booktomarket</i> | 0.0263 (0.58) | 0.0257 (0.58) | 0.0280 (0.66) | 0.0829** (2.33) | 0.0525 (1.47) | 0.0590 (1.44) |
| <i>Momentum</i> | 0.0209 (0.42) | 0.0209 (0.43) | 0.00982 (0.20) | 0.144*** (2.83) | 0.134*** (2.63) | 0.0701 (1.29) |
| R2 | 0.0754 | 0.0806 | 0.0949 | 0.0622 | 0.0652 | 0.0866 |
| N | 490296 | 478519 | 401924 | 1938326 | 1721934 | 1233446 |

We run Fama MacBeth regressions to test whether our BN_{T5-XXL} contains economic links that might not be discovered by any approaches documented in the literature. We control for analyst coverage (BN_{ANA}), similar firm characteristics (BN_{SESM}), industry membership (BN_{SIC}) and within-industry correlation (BN_{corr}) effects. We consider the time horizon from 1996 until 2021. For ease of interpretability, we measure all independent variables in quintiles. We denote the t-statistics of the coefficients in parentheses. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 10: **Static vs. Historic Business Network spillover effect**

| | US-US | US-Global | Global-US | Global-Global |
|--------------------|-------------------|-------------------|-------------------|-------------------|
| Panel A: 1996-2021 | | | | |
| Recent | 1.07*** (8.85) | 1.46*** (8.56) | 1.16*** (7.96) | 2.28*** (9.39) |
| Historical | 0.9*** (8.16) | 1.35*** (8.1) | 0.96*** (7.23) | 2.12*** (8.88) |
| Panel B: 2000-2021 | | | | |
| Recent | 0.98*** (7.4) | 1.32*** (7.13) | 1.04*** (6.39) | 2.07*** (7.79) |
| Historical | 0.83*** (7.09) | 1.21*** (6.89) | 0.86*** (6.05) | 1.96*** (7.41) |
| Panel C: 2005-2021 | | | | |
| Recent | 0.76*** (6.93) | 0.89*** (5.87) | 0.76*** (6.0) | 1.93*** (6.63) |
| Historical | 0.61*** (7.02) | 0.81*** (5.53) | 0.6*** (6.19) | 1.83*** (6.23) |
| Panel D: 2010-2021 | | | | |
| Recent | 0.81*** (7.3) | 0.67*** (3.95) | 0.82*** (5.75) | 1.74*** (4.58) |
| Historical | 0.65*** (7.55) | 0.69*** (3.96) | 0.7*** (6.43) | 1.67*** (4.18) |

We study the lead-lag relationship among US and global firms by constructing calendar-time portfolios that are rebalanced every month. We go long (short) in the 20% stocks whose most similar firms showed the best (worst) performance in the previous month. We compare the performance of a Sentence Transformer model (*T5-XXL*) based on recent and historical descriptions for different evaluation periods. We report seven-factor alphas (five-factor model plus momentum and short-term reversal) and denote the t-statistics of the coefficients in parentheses. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 11: **Logistic regression analysis of merger likelihood**

| Sector | correlation | #firms |
|------------------------------------------|---------------|-------------|
| Tobacco | 0.7038 | 118 |
| Mining | 0.6977 | 3979 |
| Industrial Metals and Mining | 0.6318 | 1719 |
| Pharmaceuticals and Biotechnology | 0.5878 | 3528 |
| Beverages | 0.5870 | 690 |
| Aerospace and Defense | 0.5853 | 420 |
| Alternative Energy | 0.5669 | 492 |
| Gas, Water and Multiutilities | 0.5650 | 541 |
| Oil and Gas Producers | 0.5544 | 2023 |
| Forestry and Paper | 0.5543 | 543 |
| Food and Drug Retailers | 0.5391 | 647 |
| Automobiles and Parts | 0.5360 | 1382 |
| Chemicals | 0.5274 | 2056 |
| Industrial Transportation | 0.5246 | 1636 |
| General Industrials | 0.5121 | 1034 |
| Personal Goods | 0.5071 | 1937 |
| Travel and Leisure | 0.5006 | 2701 |
| Food Producers | 0.4982 | 2792 |
| Health Care Equipment and Services | 0.4926 | 2135 |
| Technology Hardware and Equipment | 0.4854 | 3094 |
| Electricity | 0.4801 | 1003 |
| Financial Services (Sector) | 0.4742 | 5378 |
| Oil Equipment and Services | 0.4593 | 777 |
| Real Estate Investment Trusts | 0.4570 | 972 |
| Electronic and Electrical Equipment | 0.4520 | 2069 |
| Leisure Goods | 0.4486 | 873 |
| Life Insurance | 0.4420 | 246 |
| Construction and Materials | 0.4412 | 2764 |
| Household Goods and Home Construction | 0.4336 | 1304 |
| General Retailers | 0.4313 | 2805 |
| Real Estate Investment and Services | 0.4280 | 2916 |
| Banks | 0.4251 | 3450 |
| Fixed Line Telecommunications | 0.4031 | 930 |
| Industrial Engineering | 0.4008 | 2283 |
| Software and Computer Services | 0.3959 | 5515 |
| Nonlife Insurance | 0.3901 | 847 |
| Support Services | 0.3776 | 2323 |
| Media | 0.3555 | 2031 |

This table shows the correlation of the business description similarities of stocks operating in the same sector obtained from a Sentence Transformer model and an embedding model provided by OpenAI. While the Sentence Transformer model was released in 2019, the embedding model of OpenAI is derived from GPT-3, which was trained on data up to September 2021.

Table 12: Logistic regression analysis of merger likelihood

| | (1) | (2) | (3) | (4) | (5) | (6) |
|------------------------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| | 1:100 | 1:100 | 1:500 | 1:500 | 1:1000 | 1:1000 |
| <i>Constant</i> | -8.016*** (-47.24) | -13.60*** (-22.13) | -9.564*** (-53.85) | -15.29*** (-23.50) | -10.21*** (-56.09) | -15.99*** (-24.16) |
| <i>SameSIC4</i> | 4.024*** (32.59) | 2.919*** (28.79) | 3.796*** (29.49) | 2.901*** (26.25) | 3.713*** (28.83) | 2.873*** (25.57) |
| <i>SameCountry</i> | 3.388*** (35.88) | 2.635*** (27.71) | 3.344*** (32.15) | 2.670*** (25.82) | 3.329*** (31.54) | 2.670*** (25.32) |
| <i>SizeDiff</i> | 0.625*** (16.43) | 0.635*** (18.20) | 0.605*** (15.84) | 0.612*** (17.06) | 0.600*** (15.52) | 0.610*** (16.55) |
| <i>Debratio_{Target}</i> | 0.0897*** (4.77) | 0.0646*** (3.31) | 0.114*** (6.41) | 0.0823*** (4.36) | 0.116*** (6.43) | 0.0844*** (4.47) |
| <i>ROE_{Target}</i> | 0.0909*** (4.86) | 0.0747*** (3.65) | 0.108*** (5.91) | 0.0952*** (4.71) | 0.114*** (6.29) | 0.100*** (5.02) |
| <i>Cash_{Target}</i> | 0.0908*** (3.02) | 0.0716** (2.45) | 0.0727*** (2.58) | 0.0485* (1.71) | 0.0617** (2.20) | 0.0367 (1.29) |
| <i>Similarity_{T5-XXL}</i> | | 1.443*** (11.23) | | 1.457*** (11.05) | | 1.463*** (11.04) |
| Pseudo R2 | 0.420 | 0.511 | 0.348 | 0.420 | 0.320 | 0.386 |
| N | 680004 | 680004 | 3315204 | 3315204 | 6609204 | 6609204 |

This table presents the results of a logistic regression that examines the relationship between business description similarity and the likelihood of a merger. We use data on mergers and acquisitions from SDC and randomly select 100 (500, 1000) times as many non-merger firm pairs. The cosine similarity of the business descriptions is measured using a Sentence Transformer model (*T5-XXL*). We also control for relevant factors, such as whether the acquiring firm and potential target share the same four-digit SIC code and country, and we calculate the market capitalization difference of the firm pair. Additionally, we consider fundamental information of the (non-) target firms, such as their profitability, cash amount, and debt share. All non-categorical variables are grouped into quintiles for easier result interpretation. We cluster standard errors at the firm level. * indicates significance at the 10% level, ** indicates significance at the 5% level and *** indicates significance at the 1% level.

Table 13: Competitors, suppliers and customers of Ford: T5-XXL

| Name | Shortened business description |
|------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Panel A: Competitor | |
| General Motors | General Motors Company designs, builds and sells trucks, crossovers, cars and automobile parts and provides software-enabled services and subscriptions worldwide. |
| Toyota Motor | Toyota Motor Corp is a Japan-based company engaged in the automobile business, finance business and other businesses. |
| Hyundai Motor | Hyundai Motor Co is a Korea-based company principally engaged in the manufacture and distribution of automobiles. |
| Daimler | Daimler AG (Daimler) is a Germany-based automotive engineering company. |
| Toyota Industries | TOYOTA INDUSTRIES CORPORATION is primarily engaged in the manufacture and sale of automobiles, industrial vehicles and textile machinery. |
| Maruti Suzuki India | Maruti Suzuki India Limited is engaged in the manufacturing, purchasing and sale of motor vehicles, components and spare parts. |
| Panel B: Supplier | |
| Visteon | Visteon Corporation is an automotive supplier that designs, engineers, and manufactures automotive electronics and connected car solutions for the vehicle manufacturers including Ford, Mazda, Volkswagen, General Motors, Renault/Nissan, BMW, Jaguar/Land Rover, Daimler, and Stellantis. |
| Meritor | Meritor, Inc. is a supplier of a range of integrated systems, modules and components to original equipment manufacturers (OEMs) and the aftermarket for the commercial vehicle, transportation and industrial sectors. |
| Panel C: Customer | |
| Autonation | AutoNation, Inc. is an automotive retailer in the United States. |
| Lithia Motors | Lithia Motors, Inc. is a provider of personal transportation solutions. |
| Motus Holdings Ltd | Motus Holdings Limited is a South Africa-based automotive company. |
| Autocanada | AutoCanada Inc. (AutoCanada) is a Canada-based multi-location automobile dealership company. |
| Diesel & Motor Engineering. | Diesel & Motor Engineering PLC is engaged in import, sale and repair of passenger vehicles, commercial vehicles, car parking systems, lamps, batteries, import and sale of vehicle spares, components, accessories, providing lighting solutions and storage systems. |
| Colonial Motor | The Colonial Motor Company Limited is engaged in operating franchised motor vehicle dealerships. |
| Fujian Zhangzhou Development | Fujian Zhangzhou Development Co., LTD. is a China-based company principally engaged in the automobile trading business. |
| City Auto | City Auto Corporation, formerly Tan Thanh Do City Ford Joint Stock Company, is a Vietnam-based company primarily engaged in automobile trading sector. |
| Mercantile investments and finance | Mercantile Investments and Finance PLC is a finance company. The Loans and advances segment includes vehicle loans. |

This table contains the names and a short description of the competitors, suppliers, and customers if we apply the following filter to the economically linked firms (according to BN_{T5-XXL}) of Ford: **Competitor:** Stocks in the same industry (four-digit SIC code) and of similar size (difference in NYSE size decile smaller than two). **Supplier:** Firms with the word "supplier" in their business description. **Supplier:** Firms that operate in the "General Retailer" industry.