

Remote Work across Jobs, Companies, and Space*

Stephen Hansen
University College London

Peter John Lambert
London School of Economics

Nick Bloom
Stanford University

Steven J. Davis
Hoover Institution & Chicago Booth

Raffaella Sadun
Harvard University

Bledi Taska
Lightcast

March 1, 2023

Abstract

The pandemic catalyzed an enduring shift to remote work. To measure and characterize this shift, we examine more than 250 million job vacancy postings across five English-speaking countries. Our measurements rely on a state-of-the-art language-processing framework that we fit, test, and refine using 30,000 human classifications. We achieve 99% accuracy in flagging job postings that advertise hybrid or fully remote work, greatly outperforming dictionary methods and also outperforming other machine-learning methods. From 2019 to early 2023, the share of postings that say new employees can work remotely one or more days per week rose more than three-fold in the U.S and by a factor of five or more in Australia, Canada, New Zealand and the U.K. These developments are highly non-uniform across and within cities, industries, occupations, and companies. Even when zooming in on employers in the same industry competing for talent in the same occupations, we find large differences in the share of job postings that explicitly offer remote work.

Keywords: remote work, hybrid work, work from home, job vacancies, text classifiers, BERT, pandemic impact, labour markets

JEL Codes: E24, O33, R3, M54, C55

*Thanks for outstanding research assistance to Yabra Muvdi, who built and estimated the classification algorithm, and to Miaomiao Zhang and Kelsey Shipman, who supported the data analysis. Hansen gratefully acknowledges financial support from ERC Consolidator Grant 864863, Lambert from the London School of Economics STICERD PhD research grant and the Commonwealth Scholarship Commission, Bloom from the Smith Richardson and John Templeton Foundations, Davis from the Templeton Foundation and the Booth School of Business at the University of Chicago, and Sadun from Harvard Business School. Selected visualisations and data accompanying this paper can be found at www.WFHmap.com.

1 Introduction

The COVID-19 pandemic propelled an enormous uptake in hybrid and fully remote work. Over time, it has become clear that this shift will endure long after the initial forcing event. Looking forward, U.S. survey data say that one-quarter of full workdays will happen at home or other remote location after the pandemic ends, five times the pre-pandemic rate (Barrero et al., 2021). The pandemic also drove large, enduring increases in remote work in dozens of other countries (Criscuolo et al. 2021, Aksoy et al. 2022). There are few, if any, modern precedents for such an abrupt, large-scale shift in working arrangements.

Most previous efforts to quantify and characterize this shift rely on surveys of workers and employers and assessments of remote-work feasibility by occupation. We rely instead on the information contained in job vacancy postings. Specifically, we consider the full text of over 250 million postings in five English-speaking countries. In doing so, we apply state-of-the-art language-processing methods to analyze the text and determine whether the job allows for remote work. We fit, test, and refine our language-processing model using 30,000 classifications generated by human readings. We also identify the city, employer, industry, occupation, and other attributes associated with each job vacancy.

Vacancy postings pertain to the flow of new jobs rather than the stock.¹ In addition, postings that promise remote work two days a week, for example, entail a commitment—or at least a statement of intent—that extends into the future. For both reasons, postings need not show the same pattern of remote work as the currently employed. Indeed, the remote-work share of postings lags far behind the remote-work share of employment in the pandemic’s early stages. And while the incidence of remote work among the employed fell markedly in the two years after spring 2020, we show that the remote-work share of postings rose sharply over the same period.

Our approach to studying the remote-work phenomenon has several noteworthy strengths. First, our data cover almost all vacancies posted online by job boards, employer websites, and vacancy aggregators from 2014 to 2022 in our five countries. Coverage on this scale is infeasible with survey methods. Second, postings typically describe the job and its attributes in considerable detail, as suggested by a median posting length of 347 words. Comparable

¹Vacancy distributions by industry, employer size, and worker turnover rates differ greatly from the corresponding employment distributions, and the differences are highly sensitive to labor market conditions. See Davis et al. (2012).

detail is hard to obtain from other sources, especially at scale.² Third, we apply frontier methods to develop a language-processing model that reads and classifies postings in an automated manner. The model achieves a 99% accuracy rate in flagging jobs that allow for remote work, greatly outperforming dictionary methods. Our model also outperforms a variety of other methods. Fourth, the combination of scale, rich text data, and automation lets us characterize the shift to remote work in a highly granular manner. We track the evolution of remote work at a monthly frequency in hundreds of occupations, thousands of cities, tens-of-thousands of employers, and in city-by-occupation and employer-by-occupation cells. We post many of these statistics at www.WFHmap.com, with frequent updates.

The share of postings that say new employees can work remotely one or more days per week was tiny before the pandemic: 1% or less in Australia, Canada and New Zealand as of 2019, about 3% in the U.K., and about 4% in the U.S. From 2019 to 2022, this remote-work share rose more than three-fold in the U.S and five-fold or more more in the other countries. As of January 2023, the remote-work share exceeds 10% of postings in Australia, Canada, the U.K, and the U.S., and it appears to be on an upward trajectory in all five countries.

Remote-work posting shares vary greatly across occupations, industries, and cities. Looking across occupations, the remote-work share correlates positively with computer use, education, and earnings. Finance, Insurance, Information and Communications have especially high remote-work shares. Chicago, London, New York, San Francisco, Toronto, and other cities that function as business service hubs have high remote-work shares. These differences have widened since the pandemic struck. According to a linear least-squares regression, 63% of the variation across occupations in 2022 remote-work shares is accounted for by their 2019 shares. In contrast, just 28% of city-level variation in 2022 remote-work shares is accounted for by 2019 shares.

We also find that the shift to remote work is highly non-uniform across same-industry employers, even when they are recruiting in the same occupational category. This emergent heterogeneity on the demand side expands opportunities to satisfy preferences over remote work on the supply side.³ Our non-uniformity result also carries another important message: in many occupations, it is misleading to think of remote-work suitability as a purely technological constraint. Remote-work intensity is, instead, an outcome of choices about job design

²Previous research exploits the detail in vacancy postings to study technical change, the cyclical nature of skill requirements, their relationship to wages, how compensation and other job attributes affect applicant flows and, of course, to classify jobs in a fine-grained manner. Examples include Modestino et al. (2016), Deming and Kahn (2018), Hershbein and Kahn (2018), Davis and Samaniego de la Parra (2020), Forsythe et al. (2020), Marinescu and Wolthoff (2020), and Acemoglu et al. (2022).

³On preference heterogeneity in regards to remote work, see Bloom et al. (2015), Mas and Pallais (2017), Wiswall and Zafar (2018), Barrero et al. (2021), and Aksoy et al. (2022).

and how to operate an organization. These choices are influenced by the external environment and subject to shock-induced shifts. In line with this view, Aksoy et al. (2022) find that employers plan higher levels of work from home after the pandemic ends in countries that experienced longer and stricter government-mandated lockdowns during the pandemic.

Several recent papers use job postings to study the remote-work phenomenon. See Draca et al. (2022) for the U.K., Alipour et al. (2020) for Germany, and Bamieh and Ziegler (2022) for Austria. Bai et al. (2021) use pre-pandemic postings in the U.S. to construct firm-level indexes of remote-work feasibility, which they relate to post-pandemic performance as measured by sales, net income, and equity returns. Perhaps the closest forerunner to our paper is Adrjan et al. (2021), who use postings data at the country-sector-month level to study remote work from January 2019 to September 2021. Previous studies use dictionary methods (keyword search criteria) to identify postings that allow for remote work.

In contrast, we apply a large-scale language-processing framework to the task, an approach that is rapidly diffusing in data-science applications but, thus far, is little used in economics. First, we pre-train the DistilBERT framework (Sanh et al., 2020) on one million text chunks drawn from vacancy postings.⁴ This step familiarizes DistilBERT with the (heterogeneous) structure of job ads. Second, we consider 10,000 text chunks drawn from vacancy postings and assign three human readers to label each one. The reader assigns a positive label if the text says the job allows work from home (or other remote location) one or more days per week. Third, we fit the pre-trained DistilBERT framework to the human labels to obtain a model that classifies each posting as follows: the job allows hybrid or fully-remote work arrangements, or it does not. Finally, we apply the resulting “Work from Home Algorithmic Measure” (WHAM) to classify all 250 million job postings.

WHAM greatly outperforms a previously used dictionary in classifying the remote-work status of vacancy postings. The dictionary method yields high classification error rates that vary greatly over time and across occupations. Expressions like “home or office working possible” and “work from home care facilities” and “requires a Home Office work permit” suggest some of the difficulties that arise when applying dictionary methods to job ads. Logistic regressions and GPT-3 offer large improvements over dictionary methods. WHAM offers even larger improvements, substantially outperforming all other methods when applied to vacancy postings, as measured by accuracy rates, precision, and F1 scores.

Few previous works in economics combine a frontier language-processing framework with

⁴DistilBERT is a smaller, faster version of the BERT framework introduced by Devlin et al. (2019), which has 60,000 Google Scholar citations as of February 2023. BERT and DistilBERT exploit machine-learning tools and are pre-trained on the full English-language Wikipedia corpus and the Toronto Book Corpus. For a helpful non-technical overview of BERT, see Luktevich (2022).

human-generated labels to develop an automated classification model and to quantify its performance. In the only example we know, Shapiro et al. (2022) develops a BERT-based model and finds little gain relative to dictionaries in detecting news sentiment.⁵ However, they use fewer than 1,000 human-labeled text examples in fitting their BERT-based model, which may explain why it yields small performance gains.

A prominent line of research classifies occupations as suitable or unsuitable for remote work based on descriptions of work activities and experiences.⁶ Our analysis highlights some limitations of this approach. First, remote-work intensity is a malleable feature of jobs, occupations, and organizations. Second, classifications based on suitability assessments explain little of the variation in remote-work posting shares. For occupations that Dingel and Neiman (2020) classify as unsuitable to be done entirely from home, the remote-work share of U.S. postings in 2022 ranges from 0 to 50% with a mean of 5% and standard deviation of 7%. For occupations they classify as suitable for work from home, the share ranges from 0.3 to 74% with a mean of 18% and standard deviation of 12%.

Another prominent line of research surveys workers and employers to study working arrangements. Barrero et al. (2020), Bartik et al. (2020), Bick et al. (2022) and Brynjolfsson et al. (2020) document and characterize the enormous uptake in work from home in spring 2020. Bartik et al. (2020), Barrero et al. (2021) and Ozimek (2020) use employer plans and other forward-looking survey data to forecast that the big shift to remote work will endure. Relative to our approach, the survey-based approach is more useful for eliciting the perceptions, attitudes, and expectations of workers and employers. Our approach offers several other distinct advantages, as discussed above.

The next section describes our vacancy posting data and develops our classification model. Section 3 assesses the model’s performance in absolute terms, and relative to other approaches. Section 4 sets forth our main findings related to remote-work intensity over time and across countries, cities, occupations, and more. We also compare our remote-working posting shares to survey-based measures of remote work. Section 5 concludes.

⁵Bajari et al. (2021) and Bana (2022) use BERT to predict prices from Amazon product reviews and wages from job posting text, respectively. Each paper achieves high predictive performance. Their applications don’t involve the use of human-generated labels.

⁶Dingel and Neiman (2020) is the most influential example. Other examples include del Rio-Chanona et al. (2020), Mongey et al. (2021), and Adams-Prassl et al. (2022). Like us, Adams-Prassl et al. (2022) concludes that remote-work intensity varies greatly across jobs within occupations.

2 Data and Measurement

To measure remote-work posting shares, we exploit a near-universe of online job postings from January 2014 through January 2023 for our five countries.

We extract 10,000 text sequences from selected postings and ask humans to read them. Each sequence is about 45 words long, and the average posting has about six sequences. Breaking postings into sequences facilitates human and algorithmic classification at scale, as we discuss below. Our human readers answer this question: ‘Does this text explicitly offer an employee the right to remote-work one or more days a week?’, yielding a binary classification. The pairwise agreement rate between readers exceeds 90 percent.

We turn to the ‘DistilBERT’ language-processing framework to build a text-classification model for our purposes.⁷ DistilBERT is pre-trained on thousands of books and the English-language Wikipedia corpus, which helps it interpret the intended meaning of a given document or passage. We further pre-train on roughly one million text sequences drawn from our corpus of vacancy postings. This further pre-training step familiarizes the framework with the nature of the text in vacancy postings.

After pre-training, we use the human labels to train, or fit, a bespoke classification model for predicting these labels. We call this model the ‘Working-(from)-Home Algorithmic Measure’ (*WHAM*). We will show that *WHAM* achieves near-human performance in its classification task, and that it outperforms a variety of other approaches. We describe our approach in some detail, because we think it has useful applications to many other text-analysis tasks in economics and other fields.

2.1 Job Vacancy Data

We examine online vacancy postings collected by Lightcast (formerly Emsi Burning Glass), an employment analytics and labor market information firm. Lightcast scrapes postings from more than fifty thousand online sources that include vacancy aggregators, government job boards, and employer websites. Lightcast claims to cover a “near-universe” of online postings in our five countries during the period covered by our analysis. See Appendix A for a detailed description of our data and pre-processing steps.

⁷DistilBERT is a direct descendant of BERT, which is widely used in industry. BERT stands for Bidirectional Encoder Representations from Transformers. Transformers are a deep-learning method in which every output element is connected to every input element of a text sequence. This allows the meaning of a particular word to depend on the context of surrounding words, which as we show below is crucial in our setting. See Phuong and Hutter (2022) for a formal overview of how Transformers work. Vaswani et al. (2017) is the seminal contribution.

Burke et al. (2020) compare vacancy postings in Lightcast data for the United States to job vacancy data from the U.S. Job Openings and Labor Turnover Survey (JOLTS). The two sources are reasonably well aligned, but the JOLTS data show larger vacancy shares in food services, public administration and construction and smaller shares in finance, insurance, healthcare, social assistance and educational services.

For each online vacancy posting in our dataset, we have access to a plain text document scraped from the job listing. We also observe the posting date, employer name, occupation, location of the employer, industry, and more. We consider postings listed from January 2014 and January 2023, dropping those with an unknown occupation (less than 1%). We use a 5% random sample of postings before January 2019, and the universe of postings thereafter. The resulting dataset covers more than 250 million online vacancy postings in five countries, spanning 5.2 million employers and nearly 40 thousand cities. Table 1 provides more information.

For our baseline results, we re-weight the postings in each country-month cell to match the U.S. occupational distribution of new online vacancy postings in 2019. Appendix B reports selected results for alternative weighting schemes.

2.2 The Measurement Problem

The measurement problem we face is to determine whether each job posting allows a new hire to work remotely, understood here to encompass both fully remote and hybrid positions. We adopt a binary classification approach, and refer to a ‘positive’ posting as one that mentions the ability to work remotely, and a ‘negative’ posting as one that does not. For positions that offer hybrid working arrangements, we use a threshold of at least one day per week for our positive classification⁸ This approach effectively measures an employer’s willingness to commit *ex ante* to offering flexibility in work location. Negative postings may in fact be associated with work-from-home positions, for example because the ability to work from home is assumed by market participants to be feasible in particular jobs, or because the employer prefers to bargain over work arrangements during the hiring process rather than make a prior binding commitment. We return below to discuss the interpretation of our measure, and first focus on developing an accurate and robust classification.

The most precise way of classifying postings is arguably via direct human reading. Given the size of our data, however, this approach is not feasible to scale and some means of automated classification is required. The most standard approach adopted in the text-as-

⁸In principle our measurement approach could be extended to the intensive margin (days per week), but for simplicity we begin with the this binary classification.

data literature in economics is to use a dictionary of keywords whose presence is assumed to indicate a positive classification. As an initial step, we use the keywords in Table C.1 to classify job postings as positive or negative. While we do not claim the dictionary of terms is fully optimized, it is in line with others in the literature for classifying postings as work-from-home (Adrjan et al., 2021).

An issue that becomes immediately apparent upon inspecting job postings that are classified by keywords is the presence of notable errors, which Table 2 illustrates. False positives include references to companies' home offices and working in homes dedicated to health-care provision. A second, and perhaps most worrying, source of false positives is that the structure of job ads shifts during COVID-19 in a way correlated with the presence of false positives. This is due to the fact that, after early 2020, many postings feature a new text field indicating whether home work is allowed, and then explicitly state it is not—a naive application of the dictionary method would infer from this text field that the job posting allows working from home.⁹ Table 2 also lists examples of false negatives, which illustrates the many and complex ways that companies can use to describe remote work. Accounting for this linguistic variety with a fixed set of keywords is a major challenge.

2.3 Our Approach to Classification

Our approach to address the classification errors in the dictionary approach has three steps. First, we use at least three human auditors to read and classify 10,000 pieces of text extracted from job ads which produces 30,000 labels. Second, we train a modern machine learning algorithm using these human classifications. Third, we take this predictive model out-of-sample to classify each job ad as either positive or negative. The hope is to scale the accuracy of human reading—which can only be deployed on a small fraction of data—to the entire dataset. While this approach is common in the machine learning literature, it is not often used in economics, even though it appears to hold great promise. We call the final model used in this paper the *Working-from-Home Algorithmic Measure (or WHAM) model*.

The main text contains an overview of our methodology, with further details in Appendix C.

⁹One approach to correcting this problem is to extend the dictionary to incorporate negation (e.g. to treat as a negative classification the phrase ‘this is not a remote work position’). In section 3 we show that this indeed improves measurement accuracy but not by as much as our proposed solution below.

2.3.1 Breaking Up Job-Ad Text into Sequences

While we ultimately wish to classify job postings, we initially label and classify smaller units of text we refer to as *sequences*. The first reason for doing so is that human labeling of entire job postings is prone to a high error rate because of their length and complexity. The second reason is that the typical posting has a great deal of information unrelated to remote work, for example descriptions of the skills required for the job, the tasks involved, etc. Mixing text relevant for work-from-home with a great deal of irrelevant text introduces noise into the classification algorithm.

The procedure for generating sequences has three salient features. First, postings always begin with a job title, e.g. “Software Programmer familiar with R and Python.” We extract these as a single sequence. Second, the beginning of each posting typically has a number of bullet points or other structured fields. In most cases, these also form a single sequence.¹⁰ Finally, the remainder of a posting is typically structured like standard prose with a succession of paragraphs. Each paragraph is taken as a single sequence, unless it passes a length threshold. In this case, we break it into multiple sequences of consecutive sentences.

This procedure produces approximately 1.6 billion sequences out of the over 250 million job postings.

2.3.2 Human Labels for Training and Evaluation

From the sample of sequences, we first chose 10,000 to label manually. One quarter of these sequences was chosen at random from the set of sequences that contained a set of dictionary terms listed in Table C.1. Another quarter was chosen to contain a broad set of terms that might reflect work-from-home language, including the generic terms ‘remote’, ‘home’, ‘work’, ‘location’; any word that begins with ‘tele’; and any two-word sequence that begins with ‘remote’. Another quarter consisted of sequences that might confound a classifier, including ‘home repairs’, ‘nursing home’, ‘remote construction’, etc. The final quarter was a random sample of sequences not satisfying the three aforementioned criteria. Each portion of the label sample is balanced across year-quarter from 2014Q1 through 2021Q3. We also balance the sample evenly across countries¹¹ to account for varying English idioms in different geographic locations.

We used Amazon Mechanical Turk to generate labels. To ensure high-quality workers,

¹⁰The exception is if the number of distinct structured fields is too large, in which case we split them into multiple sequences.

¹¹We draw one quarter of this dataset from each of USA, UK, Canada, and a further one quarter from the pooled Australia and New Zealand data.

we set up an initial screening test that required prospective workers to label 20 sequences that we had previously manually classified. Only workers that made at most one error were allowed to proceed to label the full set. Another quality control strategy was to pay around 25% above typical market rates for labeling tasks. This motivated workers who passed initial screening to continue on the project.¹²

Each of the 10,000 sequences were labeled by three distinct workers. There is a high agreement rate among workers: 66.9% of examples are unanimous negative examples and 25.5% are unanimous positive examples. The remaining 7.6% examples are evenly balanced between one dissenting vote for either positive or negative. Note that, while half of the sample was chosen to contain a word with the potential to denote work-from-home, only 29.2% of the sequences receive a majority of positive votes.

2.3.3 Developing the WHAM model

In the last five years, the field of natural language processing has been revolutionized by models that allow the meaning of word sequences to arise by how they interact. Consider the sentences ‘Some of the deep-sea wells we operate are in remote locations’ and ‘We are pleased to offer opportunities for remote work’. Each includes the word ‘remote’ but only the latter is a positive example of remote work. The important point is that the interaction of ‘remote’ with surrounding context words determines the overall meaning of these sentences. Moreover, not all context words are equally informative: for example, in the first sentence ‘deep-sea’, ‘wells’, and ‘locations’ are more important than ‘some’ and ‘we’ in understanding the meaning of ‘remote’. *Self-attention* (Vaswani et al., 2017) is a mathematical construct that allows vector representations of individual words to interact with each other to form new vectors that encode the meaning of sequences. These interaction weights effectively determine which words should be “paid attention to” in resolving these meanings. Self-attention is the key idea that powers models such as *BERT*, *RoBERTa*, *GPT*, *GPT-3*, *PALM*, and, most famously, *ChatGPT*. For more details on these models, collectively called Transformers, see Ash and Hansen (2023).

The particular Transformer model we adopt in the development of WHAM is *DistilBERT* (Sanh et al., 2020). DistilBERT is based on Google’s BERT model (Bidirectional Encoder Representations from Transformers, Devlin et al., 2019), which when originally released set important new performance benchmarks for common NLP tasks (since eclipsed by larger-

¹²In general workers appeared engaged and focused on the labeling task. We received communication from multiple workers seeking to clarify ambiguous cases, which went above and beyond what AMT required for payment.

scale models). Since 2020, Google has used BERT to process its online search queries. DistilBERT ‘distills’ the information in BERT by training a simplified model to reproduce the same output as BERT. The main advantage is that DistilBERT obtains an expressive language model with far fewer parameters which reduces processing time and estimation costs.

We make two main modifications to the off-the-shelf DistilBERT model to build *WHAM*. First, the initial set of parameters of off-the-shelf DistilBERT is obtained by predicting randomly deleted words in generic English from surrounding context words. We instead update these parameters to predict randomly deleted words in a sample of 900,000 job posting sequences which is balanced across all years and countries. This step creates word representations that are specific to the language of job postings.

Second, we further modify off-the-shelf DistilBERT to predict human labels from vector representations of job posting sequences. We split our labeled sequences into training and test sets of 5,950 and 4,050 sequences, respectively. The prediction problem is conducted at the label rather than sequence level, so there are $3 * 5,950 = 17,850$ total observations in the core training sample of labeled data.¹³ Appendix C details how we specify the prediction model’s hyper-parameters. Table 3 provides an illustration of which words are influential in the classification problem, and compares the WHAM approach to a dictionary approach. Crucially, the weights attached to words are learned by the algorithm rather than imposed *ex ante* by researchers, and the weight on a particular word depends on the surrounding words. In the following section, we compare the test-set accuracy of the estimated model with that of other algorithms in the literature and show its performance is outstanding.

2.3.4 Predicting Remote Work Language at Scale

Finally, we use the estimated prediction model to assign a continuous probability to all sequences in our corpus. The higher the probability, the more confidence the model has that this sequence denotes an offer of remote work arrangements. Figure B1 plots a histogram of the share of sequences that fall in different probability intervals. The distribution is bimodal at the lowest and highest probability bins, with the former dominating the distribution. As expected, most sequences do not contain work-from-home language because, as we show below, most job postings do not explicitly mention the possibility to work from home and, among those that do, the majority of sequences discuss other features of the position. The

¹³During an initial exploratory phase, we labeled a sample of around 10,000 additional sentences (rather than sequences) using a combination of Mechanical Turk, hired research assistants, and ourselves. Since these are also potentially informative, we include them in the training set. In most cases, these sentences only received a single label and so in total generate 11,574 additional labels in the training set.

bimodality of the distribution shows that the classification algorithm typically produces a clear prediction, in line with human labelers’ high agreement rates. We use an 0.5 threshold for assigning a sequence a positive classification according to WHAM’s predicted probability, but the properties of the predicted probability distribution imply that our results are not sensitive to this particular cutoff.

2.3.5 Aggregating Measurement back to Job Postings

We have conducted all the analysis so far at the sequence level, but are ultimately interested in a job-posting-level classification. For this, we use a simple ‘max’ rule and positively classify a job posting if it contains one or more positive sequences. Table B.1 shows the number of positively classified sequences in each job ad. We can see that among the positive job ads (those with one or more positive sequence), the majority have just a single positive sequence. This reduces concern that the algorithm produces correlated false positive hits at the posting level.¹⁴ This posting-level classification constitutes the final output from our WHAM model, which we use to study the adoption of remote work.

2.3.6 Public Access to WHAM

To allow researchers to interact with and study the properties of our model, we make available a simple online tool that allows one to input arbitrary text and receive a predicted probability as output. The URL is <https://huggingface.co/spaces/yabramuvdi/wfh-app-v2>, which will reproduce the same probabilities as in the paper.¹⁵ One can verify that the examples in Table Table 2 that confound dictionary approaches are correctly classified by WHAM.

2.3.7 Computational Performance of WHAM

One constraint on implementing large-scale NLP models is computational. To provide some performance guidelines, Table Table B.2 tabulates the hardware we use for each step of development, the time taken, and the cost involved. All estimation is done on the Google Cloud Platform. In neither time nor money terms is the implementation of WHAM particularly

¹⁴We have manually read a number of randomly drawn postings with more than five positive sequences, and found no instance of the algorithm failing. In some cases, the scraping procedure that gathers data from online job portals appears to have identified as a single job ad a succession of postings by recruitment agencies. In other words, the measurement error arises from the data itself rather than the classification approach.

¹⁵Many thanks to Yabra Muvdi for estimating the model and making it accessible. The model is subject to revision, at which point the predicted probabilities for a given text may change. Users who find systematic biases in the predictions are welcome to contact the authors with their findings, which can be incorporated into future work.

costly from a computational perspective: the total run time for all steps is 51 hours, and the total cost is approximately \$1,500. Our view is that researchers should therefore not view computational costs as a major impediment to adopting large language models.

On request, we make available all code to efficiently train WHAM and apply it out-of-sample. Interested researchers should register interest at WFHmap.com.

3 Assessing the Performance of WHAM

Above we highlight instances in which the presence or absence of keywords is insufficient to correctly classify a selection of job posting texts due to the complexity of surrounding context. In order to quantify the gains from adopting our approach, we now undertake a systematic comparison of the ability of different algorithms to correctly classify unseen texts. To do so, we adopt a standard approach in the machine learning literature and randomly split the 10,000 human-labeled sequences into training and test sets (of sizes 5,950 and 4,050, respectively). We then train WHAM just on the training data and use the fitted model to assign a predicted value to each test-set observation. By way of comparison, we also use the following alternative methods for classifying test-set observations (full details of each approach are in Appendix C):

1. *All Zero*. Each test-set observation is assigned a 0 to match the modal outcome.
2. *Dictionary*. We use a term set similar to that from Adrjan et al. (2021),¹⁶ and count an observation as positive if it contains a term from this set.
3. *Dictionary with Negation*. Shapiro et al. (2022) shows that accounting for negation can improve the performance of dictionaries. We adopt a similar method and only count the presence of a dictionary term as indicating remote work when a negation term does not appear in the surrounding context.
4. *Logistic Regression*. Adams-Prassl et al. (2020) uses Lightcast data from the UK to measure the prevalence of flexible work schedules, i.e. the times at which work must be completed, from job posting text. The paper uses humans to manually annotate 7,000 texts, and fits a (penalized) logistic regression model for classification. The features of the logistic regression are the word frequencies in a given document. We implement a similar logistic model on our training data and use it to classify test data.

¹⁶The terms are reported in Table Table C.1. The remote work measures in Adrjan et al. (2021) are based on data from Indeed which potentially has a different structure from the Lightcast data.

5. *Logistic Regression with Negation.* We expand the feature set of the logistic regression to incorporate negation and re-estimate it on the training data.
6. *GPT-3.* Brown et al. (2020) introduced GPT-3, a large language model capable of performing a variety of natural language tasks with limited or no training examples to learn from. We query the model with the prompt “Identify if the text offers the possibility of remote work at least one day per week” and convert the answer into a 0/1 classification.¹⁷
7. *WHAM with Generic English.* Rather than pre-train DistilBERT using job posting text, we use its off-the-shelf word embeddings trained on general English.

Table 4 reports the test-set performance for all methods. A straightforward metric is the error rate, i.e. the fraction of mis-classified texts. On this measure, WHAM outperforms all other methods with an error rate of 0.02.¹⁸ GPT-3 has an error rate three times that of the baseline model, while the dictionary method’s error rate is eight times higher. On the other hand, the pre-training of the model with job posting text generates only a modest improvement over generic English.

A more standard performance metric in the machine learning literature is the F_1 score which accounts for both a classifier’s ability to recover the true positives (*recall*) as well as the share of predicted positives that are true positives (*precision*).¹⁹ The F_1 score varies between 0 and 1, where higher values indicate better performance. Again, we observe that WHAM substantially outperforms all other measures.²⁰

One concern is that the distribution of positive and negative postings in the test data does not correspond to that of the full population of job postings: the data extracted for labeling is specifically designed to over-represent positive cases. To obtain a sense of classification

¹⁷Deploying GPT-3 on the full Lightcast dataset would be prohibitively expensive at current costs, but we still report its test-set performance for benchmarking purposes. More recently, ChatGPT, a successor model to GPT-3, has generated a great deal of public interest. ChatGPT is largely built on an underlying model that OpenAI calls *text-davinci-03* whereas GPT-3 is built on *text-davinci-02*. In our experiments, the former outperforms the latter, so we only report results for GPT-3.

¹⁸This error rate is consistent across countries and years. When broken down by country, the test set error rate is 0.02 in each case. When broken down by year, the set error rate is 0.02 in each year except for 2015 (error rate 0.03) and 2014 (error rate 0.01).

¹⁹Among a set of classified observations, let TP, FP, TN, and FN be the number of true positives, false positives, true negatives, and false negatives, respectively. Precision is $TP / (TP + FP)$, and recall is $TP / (TP + FN)$.

²⁰An alternative dictionary for measuring remote work adoption is proposed in Draca et al. (2022) which uses our same UK Lightcast sample. The overall error rate of this dictionary in the full test data set is 0.19 and for the test data set arising from the UK is 0.17. Interestingly, the F_1 score we obtain for logistic regression (0.81) is similar to that reported by Adams-Prassl et al. (2020) for classifying flexible work scheduling (0.83, see Table 3 of that paper).

accuracy on the full population, we create a simulated dataset of $1000 * 4,050 = 4,050,000$ observations, 3% (97%) of which are sampled with replacement from the set of positive (negative) test set examples. Table 4 reports the same metrics as Table 4 but computed on this more unbalanced dataset. Again, we find that WHAM outperforms all other methods, but in this case the difference in F_1 scores is even starker. Our baseline WHAM achieves a 0.85 F_1 score, while the F_1 score of GPT-3 falls to 0.52 and other methods have even worse performance. Moreover, pre-training becomes more important as the F_1 score for WHAM with generic embeddings drops to 0.78. These results arise because, as Table 4 shows, WHAM has a particularly low false positive (FP) rate compared to other methods. When negative examples dominate the evaluation sample, correctly classifying them becomes important for overall performance and WHAM is strong in this dimension. Since this sample’s label composition is more in line with the expected composition of the universe of job postings, our findings highlight the potential gains in accuracy of using our approach.

We view these results as methodologically important because they are among the first, to our knowledge, to quantify the gains of adopting modern NLP methods for text classification in economic environments. There are very few papers in the economics literature that systematically compare different classification approaches, and those that do have not found that large language models outperform simpler approaches. For example, Shapiro et al. (2022) does not report large gains from using BERT over simpler models for classifying news sentiment. One reason that we, in contrast, do find large gains is the size of our training data. Shapiro et al. (2022) trains BERT on 800 labeled articles whereas we have an order of magnitude more training data, which provides more information for estimating the complex ways in which word sequences map into outcomes. We conjecture that other prediction problems using text in economics might similarly benefit from a large training sample combined with sequence embedding models.

A separate question is how WHAM compares to alternative methods on the full data sample. Rather than consider all alternatives, we focus on how WHAM compares to the Dictionary method, which is most common in the literature measuring remote work adoption from job posting text. Figure 1 plots monthly time series of the share of remote work postings in the US sample from 2019 through early 2023.²¹ The patterns present in both series differ markedly. According to the Dictionary method, the remote work share surged at the onset of the COVID-19 pandemic, peaked in early 2021, and fell markedly throughout 2021 before stabilizing in 2022. In contrast, the WHAM method suggests a more modest immediate

²¹These time series are computed using the approach we adopt for the baseline results discussed in the next section, and are not the simple raw positive share.

reaction to the pandemic followed by a steady growth rate thereafter. Two features of the Dictionary series are of note: First, the initial COVID-19 shock drove a large number of both real and negated mentions of remote work arrangements, so this series increases much more dramatically than the WHAM series. Second, towards the end of 2022 a handful of very large job boards altered their structure to partially address this issue of negation. Importantly, this second event appears not to have induced a discontinuity in our WHAM measure, likely because it is robust to changes in structure so long as the intended meaning remains consistent. Clearly, then, the choice of measurement approach can have important quantitative implications even in aggregate.²²

Of course, aggregate comparisons between methods can mask underlying differences at more granular levels. To illustrate this, we compute the growth rate in remote work adoption according to the Dictionary method and WHAM from 2019 to 2022 for individual SOC2 occupations, pooling all 2019 postings and 2022 postings together. In these two years, the Dictionary method appears similar to WHAM but with an upward shift of around five percentage points. However, as Figure 2 reveals, there are large differences in the specific occupations that each method associated with growth in remote work adoption. According to the Dictionary method, the ‘Food Preparation and Serving’ occupation has experienced highest growth in adoption, while for WHAM the highest-growth occupation is ‘Computer and Mathematical’. Moreover, according to WHAM all occupations experienced positive growth in adoption, whereas adoption rates fall for the ‘Farming, Fishing, and Forestry’ occupation according to the Dictionary method. The higher accuracy of WHAM in the sample of human labels suggests its ranking of occupations is more reliable. In the next section we provide a more in-depth analysis of occupation-level heterogeneity according to WHAM.

In sum, the WHAM model displays a very high classification accuracy—relative to human labels—and differs markedly from the most popular alternative approach in the literature based on keyword search. This difference is especially pronounced since 2020, even at the aggregate level. We believe our approach to measurement provides a highly accurate classification of remote work offers in the text of job postings, and base the remainder of the paper on analyzing its output.

²²The patterns in the Dictionary series need not match those from Adrjan et al. (2021) even though we use a similar set of keywords, as the structure of the Lightcast data could differ in important ways from that of the Indeed data that Adrjan et al. (2021) use.

4 Results

In this section we document how the *percent of new remote work vacancies*—the fraction of all new vacancies which explicitly offer the right to work remotely one or more days per week—has changed over time. We document this across countries, occupations, cities, and employers. This covers both hybrid and fully remote work.

This section is organised as follows: First, we look at the percent of remote work vacancies across each of our five countries. We plot this as a monthly time series, spanning January 2014 to January 2023. Second, we compare the percent of new remote work vacancies across broad and narrowly defined occupations, contrasting our measurements in 2019 2022. We show that the substantial rise since the onset of COVID is highly uneven across occupations, and find that occupations with the highest 2019 percentage of remote work were the most likely to top the list in 2022. We also compare occupation-level classifications used in the literature to our measurement. Third, we compare the percentage of new vacancies offering remote work arrangements across cities. We show that cities with higher remote work percentages in 2019 do not strongly predict higher percentages by 2022 (unlike occupations). This suggests that additional confounding city-level characteristics have played an important role in the adoption of remote work. We also compare a monthly time series across a selection of US cities. Fifth, we compare our measures to survey information from the American Communities Survey (ACS). We show that MSA’s which have a high remote work share of vacancies in our data also have high fractions of the population who selected “Worked from home” when asked about their commuting methods. Fifth, we show that the percentage of remote work vacancies posted by employers who operate in the same industry, and search for the same talent, can vary widely.

4.1 Remote Work across Countries

How did the share of advertised hybrid and fully remote work differ across countries prior to, during and after the pandemic? In Figure 3 we plot the monthly time series of the share of advertised remote work for the US, UK, Canada, Australia and New Zealand. For each country and in each month, this figure reports the weighted-mean of the percent of remote work vacancies across nearly 800 narrow occupation groups. We weight each group based on the share of vacancies in this group in the USA during 2019. Our baseline results utilise this method to reduce the impact of compositional differences, both across time and across countries. Three high-level facts emerge:

1. **Unprecedented and sharp increase of advertised remote work at the onset**

of COVID-19

In March-April 2020, the share of new job vacancies which advertised remote work saw a sharp rise across all countries. On average, the increase from February 2020 to April 2020 was 200%. While this immediate increase occurred across all our countries, the level-change was most pronounced in countries with a more severe initial COVID outbreak (USA, UK and Canada).

2. Sustained growth thereafter

Since the large spike in March-April 2020, there has been sustained growth in the percent of advertised remote work. In level-terms, this growth has been most pronounced in the UK (where COVID lockdowns most lingered and were most severe relative to the other countries in the sample). We also see evidence of higher growth rates in Australia and NZ, as their pandemic experience worsened during 2021. In all countries, the growth in advertised remote work has continued long after the forcing event of the pandemic subsided. An additional reason for this high and persistent growth is that our measure of new job vacancies lags the stock of employees actually working from home, possibly because employers were slow to accept this as a permanent practice.

3. Substantial heterogeneity across countries, even before the pandemic

The USA had nearly 4% advertised remote work share in 2019, the highest of any country. The UK was only marginally lower, where as Australia, Canada and New Zealand had respectively half, a third, and a tenth the share of the US in 2019. By mid 2022 the spread in levels is much greater, but proportional differences have reduced.

In our robustness exercises, we also look at the raw shares of remote work, i.e. without the re-weighting applied to our baseline Figure 3. Comparing the unweighted Figure B.2 to Figure 3 tells us the direction and magnitude of the impact that occupation composition plays in our results. For example, in mid 2022 the difference between the UK and USA is 8 percentage points using the raw data and 4 percentage points after re-weighting. This suggests that roughly half of the difference in advertised remote work shares between the US and UK is accounted for by differences in the types of jobs being advertised, which is unsurprising as the UK is on the whole more skewed towards white-collar jobs with a higher propensity to be worked from home.

4.2 Remote Work across Jobs

We first show the share of advertised remote work by grouping job ads into broad occupation groups (based on two-digit SOC 2010 classifications), which Figure 4 reports. For this, we look only at data from the United States. The differences across broad occupation groups varies greatly. In 2019, we see that just one-in-twenty positions of all job ads in ‘Computer and Mathematical’ occupations explicitly offered remote work arrangements in their postings, whereas in 2022 this share raises to a more one-third of new ads offering remote work. As one might expect, the share of advertised remote work correlates positively with computer use, education, and earnings and is lower in occupation groups which require specialised equipment or customer interactions. Lastly, Figure 4 provides some evidence that the 2019 shares of remote work correlates with post-pandemic shares.

To investigate the relationship between 2019 and 2022 shares further we next turn to an analysis at the detailed ONET occupation-level. We group our US job vacancies into granular occupations (using O*NET definitions), and plot both the 2019 and 2022 percent of advertised remote work (on a log-scale), presented in Figure 5. After dropping a handful of data points with fewer than 250 postings in 2019, we retain 875 O*NET occupations.²³ Figure 5 also shows the feasibility classification according to Dingel and Neiman (2020). A black circle represents jobs which these authors classify as ‘not suitable for full-time telework’, and an orange triangle denotes the opposite²⁴. An unweighted ordinary-least-squares trend line is also depicted in blue. Our main takeaways from Figure 5 are:

- The bivariate unweighted-OLS fit using a log-log specification yields an R^2 of 0.63, which shows that—for a given ONET occupation—the share of vacancies which advertised remote work pre-pandemic was strongly predictive of the share post-pandemic.
- The slope coefficient from the bivariate unweighted-OLS model shows that the elasticity of 2019 percent to 2022 percent was 0.76%.²⁵
- Across all ONET occupations depicted, the mean share of new postings which advertised remote work was 4% in 2019 and 10% in 2022.
- There is substantial variation in the share of advertised remote work across occupations, which grows over time. Across all ONET occupations, the standard deviation in the

²³In total, there are 1,016 O*NET occupations. Our sample of O*NET codes which have greater than 250 vacancy postings in 2019 is 875. This attrition is expected, for example a number of military occupations are not present in our data.

²⁴These are taken from the authors replication data, accessed April 2022, which can be found here.

²⁵Our ordinary least-squares estimates impose a power-law coefficient, given the log-log specification.

shares of advertised remote work was 5% in 2019 and 11% in 2022.

- Dingel and Neiman (2020)’s classification can account for a small part of the variation in the 2022 levels of advertised remote work. For occupations that they classify as unsuitable to be done entirely remotely, the share of advertised remote work in 2022 ranges from 0 to 51% with a mean of 5% and standard deviation of 7%. For occupations they classify as suitable for telework, the share ranges from 0.3 to 74% with a mean of 18% and standard deviation of 12%.²⁶

We view three key points of difference between our measurement approach and those measures which assess telework feasibility for each occupation. First, since our measurement works at the job vacancy level and not the occupation level, our measure offers more variation and a signal heterogeneity in remote work feasibility within occupations and across firms. Second, whereas the feasibility measures treat each job as a collection of tasks, our measure combines both task-feasibility as well as employer and employee preferences, labour market forces, past experience with remote arrangements, and so on.²⁷ The third reason for the discrepancy is that our measurement exercise will likely have some amount of under-reporting, as employers may not explicitly advertise remote work in their vacancies but none-the-less allow such arrangements.

4.3 Remote Work across Cities

Next we compare the percent of new vacancy postings which advertised remote work across cities. Job posting are matched to a city based on specific locations listed on the website from which it was scraped, or else mentioned in-text.²⁸

²⁶In a few cases, the D&N machine classification appears very inaccurate. For example, travel agents have been classified as ‘not teleworkable’, although both before and after the pandemic roughly 1-in-3 jobs advertised remote work. This is likewise the case for ‘Advertising Sales Agents’ and ‘Interpreters & Translators’. Some of these outliers appear to be resolved by the hand coded measure, but these data are only available at a higher level of occupational-aggregation.

²⁷A clear example of the differences between our measurement approach and Dingle and Neiman (2020) is for teaching jobs. For example, while D&N correctly classify jobs for “physical education teaching” as being *feasible* for full time home working (i.e. via a virtual class room), we know anecdotally that this arrangement was very taxing on staff and avoided as soon as normal schooling resumed. We find that teaching jobs in general (and “physical education teachers” in particular) have some of the lowest shares of advertised remote work of any job, highlighting that feasibility and actual behaviour can vary markedly.

²⁸Since the predominant remote work arrangements are hybrid, the location of the work site remains a key feature of most jobs. However, in the case of a ‘fully remote’ position this analysis becomes less precise. We plan to refine our measurement approach in future work to distinctly classify ‘hybrid’ vs ‘fully remote’ work arrangements, but have thus far concluded that the majority of remote work jobs offer hybrid arrangements

Figure 6 shows the percent of advertised remote work across a selection of large international cities, both for 2019 and 2022. We see that the percentages vary widely. For example, in 2022, 1-in-4 new job postings in Washington (DC) advertised remote work arrangements, compared to 1-in-14 in Perth, Australia. The substantial increases as well as the large heterogeneity in these shifts can be seen both across- and within-countries.

Further evidence of the large shift in both levels and spread of remote work ads is shown in Figure 7, which plots all cities in our data with more than 250 new vacancies in 2019. The mean (standard deviation) increased from 4% (2%) in 2019 to 10% (5%) in 2022. An unweighted OLS regression line fitted to these city-level data show a much lower coefficient of determination (R^2) of 0.28, compared to the value of 0.63 when running the same exercise across occupation-groups. This highlights that the 2019 shares are far less informative predictor of post-pandemic shares at the city level.

This sizable increase in the levels and spread of remote work across cities, as well as the weak relationship between 2019 and 2022 shares, poses an interesting question: What are the city-level determinants of remote work adoption? We hypothesize that a mix of institutional features, infrastructure quality, pandemic severity (both in disease and policy) and the composition of jobs and firms in each city are all important factors. We leave a more formal tests of these predictions to future work.

We next turn to more granular monthly time series for selected US cities, shown in Figure 8. As well as illustrating the granularity of our data, a number of interesting features emerge from these time series:

- Cities from the North-East and West regions (e.g. San Francisco, Boston, New York) all experience similar increases at the outset of the pandemic, but have very different growth levels subsequently. By 2023, these differential growth rates result in very dispersed levels.
- We also see substantial fluctuations over time in these North-East and Western cities. These fluctuations appear to be correlated across series, for example the July 2021 dip occurs in SF, Boston, Colorado, and to a lesser extent NYC.
- By contrast, cities from the South show far less growth since COVID and also less volatility. Savannah and Miami Beach appear to have partially reverted back to pre-pandemic shares of advertised remote work.
- Note that in this exercise, we do not re-weight the data, such that much of the variation across cities is likely to be driven by differences in occupation and industry composition.

We leave as future work a mapping from our time-series measures and forcing events, such as shelter-in-place orders.

4.4 Comparing Job Advert Measurement to Survey Responses across MSAs

Our measurement of remote working utilises new job postings, which is conceptually a very different empirical object to measures of the share of employees / work days conducted in peoples homes. To understand how these different measurement approaches might relate (if at all) to one another, we utilise recent survey evidence form the American Communities Survey (ACS).²⁹ Specifically, we use the (survey weighted) share of 2021 employees across Metropolitan Statistical Area’s (MSAs) who, when asked about their commuting method, respond that they ‘Worked from home’³⁰. We compare this to the fraction of new job ads from each MSA which advertised remote work in 2022.

Figure 9 compares the our measure of remote work from job ads to the ACS’s survey measure. A least-squares regression line (shown in blue) has a slope coefficient of 0.36. Strictly interpreted, this suggests that *ceteris paribus* an MSA with 10% more employees who respond ‘Working at home’ to the ACS commute question would accompany a 3.5% increase in the percent of new job ads offering hybrid or fully remote arrangements. The overall fit of this least squares regression line is rather high, with a coefficient of determination (R^2) of 0.55. Taken in tandem, this evidence suggests our measurement approach relates to the stock of remote workers. This, along with the many advantages of working with job postings (large scale, very high granularity, long historical time series) support the use-case for our data.

4.5 Remote Work across Companies

Ultimately, the decision to advertise remote work arrangements is made by each employer who is searching for talent. By-and-large, workers value the flexibility to work some days remotely, with survey evidence estimating that a typical worker would sacrifice 6% of their

²⁹The American Community Survey (ACS) is a demographics survey program conducted by the U.S. Census Bureau. The ACS regularly gathers information previously contained only in the long form of the decennial census, such as ancestry, citizenship, educational attainment, income, language proficiency, migration, disability, employment, and housing characteristics.

³⁰ACS respondents are instructed to “Mark (X) ONE box for the method of transportation used for most of the distance” which suggests that only those who work in a fully remote capacity should select this box, since persons with 1+ days of commute per week have more mileage from that commute mode.

salary to receive this amenity (Barrero et al., 2021). Thus, one important reason why employers have increasingly chosen to offer remote work arrangements even after the pandemic is to attract workers. Similarly, remote work arrangements can also lessen the burden of distance and allow firms to recruit for talent in wider geographic areas. Again, this deepens the labour market and may facilitate matching with better candidates. Another reason why we see that employers are offering remote work arrangements in their vacancy listings might be due to learning. Most CEO’s comment that mass remote-work of staff would have been unthinkable prior to 2020, yet the forced experimentation during COVID-19 has left many with at least an indifference to such practices and at most tangible evidence of the productivity benefits these bring. Finally, firms—especially those who are expanding quickly—may see remote work arrangements as a way to reducing office space and energy consumption. On the other hand, the need to adjust internal processes to a fully or partially remote workforce may also inhibit firms from explicitly committing to this work arrangement.

Our analysis of employers is by no means exhaustive, and we leave for future work a more in-depth match to firm-level covariates. The first piece of analysis illustrates that the prevalence of employers who explicitly offer remote work arrangements in their vacancy postings varies greatly, even among same-industry firms recruiting in the same occupational category. Figure 10 takes selected employers, and finds:

- Figure 10: Panel A shows the share of remote work vacancies posted by four large aerospace manufacturing firms (NAICS code 3364). We consider only management occupations in this panel, and find that both Boeing and Lockheed Martin explicitly offer remote arrangements in half of their postings in 2022.³¹ We further see that Northrop Grumman makes explicit offers of remote work in less than one-in-four management job vacancy postings. In contrast, SpaceX made no explicit offers for such arrangements in any of its new job listings in 2022. All of these firms explicitly offered minimal amounts of remote work in 2019.
- Figure 10: Panel B shows selected insurance firms who advertise vacancies for workers in the mathematical science occupations. We chose this occupation because we know it has a high national share of remote work vacancy postings. We see that United Health had a sizable fraction (52%) of vacancies which explicitly offered remote work,

³¹Without further analysis, we cannot say if a 50% share of remote work vacancies in 2022 results from some switch in behaviour (e.g. if firms post uniformly in time, and switch to fully remote halfway through the year, we would calculate a 50% share) or else if this is driven by some more granular cross-sectional difference between jobs with and without remote arrangements on offer. This will be addressed in future revisions.

even pre-pandemic in 2019, which grew to 80% by 2022. Mutual of Omaha had a more modest pre-pandemic remote work share, but by 2022 mentions such practices in nearly all vacancies targeting mathematicians. Humana saw more than a doubling in its remote work share, but remains substantially lower than its peers.

- Figure 10: Panel C conducts the same exercise for selected auto manufacturing firms who hire engineers. Almost no explicit offers of remote work were made in 2019. During 2022, Honda explicitly offers 1-in-2 new engineering hires the right to work remotely at least one day per week. GM offers less than half this number, and Ford less than a 6th. Tesla job postings make almost no offers of remote work in either 2019 or 2022.

5 Conclusion

This paper’s first contribution is to develop a methodology for classifying job postings as offering fully remote or hybrid work arrangements. We take an off-the-shelf, large-scale Transformer model and adjust it to both account for the specific language structure of postings and, more importantly, to predict tens of thousands of human-labeled sequences. The resulting WHAM algorithm substantially outperforms existing methods in terms of out-of-sample classification accuracy, including the language models that underlie GPT-3 and ChatGPT. To the best of our knowledge, this is the first attempt in the literature to assess the gains from adopting large language models for economic measurement, and our results suggest the promise of the method more broadly.

We then apply WHAM to the full universe of Lightcast data across five English-speaking countries (USA, UK, Canada, Australia, and New Zealand). This generates a dataset of remote and hybrid work adoption whose scale, granularity, and high frequency extend well beyond what is possible to achieve with surveys. We use this to zoom in on countries, occupations, industries, cities, and firms and, in each case, document a high degree of heterogeneity in remote work adoption since the pandemic. Moreover, this heterogeneity is not simply a function of pre-pandemic conditions. For example, the incidence of remote and hybrid work across cities in 2019 explains relatively little of the cross-city increase in adoption since. We conjecture that the heterogeneity we document has its roots in myriad forces, including worker and firm preferences, competitive pressures in the labour market, and local norms. An important topic for future research, which our measures can help advance, will be to quantify the relative importance of these factors.

Many of the data series in this paper are available through a companion website WFHmap.com which we will continue to update regularly going forward.

References

- Acemoglu, D., Autor, D., Hazell, J., and Restrepo, P. (2022). Artificial Intelligence and Jobs: Evidence from Online Vacancies. *Journal of Labor Economics*, 40(S1):S293–S340.
- Adams-Prassl, A., Balgova, M., and Qian, M. (2020). Flexible Work Arrangements in Low Wage Jobs: Evidence from Job Vacancy Data. *SSRN Electronic Journal*.
- Adams-Prassl, A., Boneva, T., Golin, M., and Rauh, C. (2022). Work that can be done from home: Evidence on variation within and across occupations and industries. *Labour Economics*, 74:102083.
- Adrjan, P., Ciminelli, G., Judes, A., Koelle, M., Schwellnus, C., and Sinclair, T. (2021). Will it stay or will it go? Analysing developments in telework during COVID-19 using online job postings data.
- Aksoy, C. G., Barrero, J. M., Bloom, N., Davis, S. J., Dolls, M., and Zarate, P. (2022). Working From Home Around the World.
- Alipour, J. V., Falck, O., and Schüller, S. (2020). Germany’s Capacities to Work from Home.
- Ash, E. and Hansen, S. (2023). Text Algorithms in Economics. Unpublished Manuscript.
- Bai, J. J., Brynjolfsson, E., Jin, W., Steffen, S., and Wan, C. (2021). Digital Resilience: How Work-From-Home Feasibility Affects Firm Performance.
- Bajari, P., Cen, Z., Chernozhukov, V., Manukonda, M., Wang, J., Huerta, R., Li, J., Leng, L., Monokroussos, G., Vijaykumar, S., and Wan, S. (2021). Hedonic prices and quality adjusted price indices powered by AI. Working Paper CWP04/21, Cemmap.
- Bamieh, O. and Ziegler, L. (2022). Are remote work options the new standard? Evidence from vacancy postings during the COVID-19 crisis. *Labour Economics*, 76:102179.
- Bana, S. H. (2022). Work2vec: Using Language Models to Understand Wage Premia. Unpublished Manuscript.
- Barrero, J. M., Bloom, N., and Davis, S. J. (2020). COVID-19 Is Also a Reallocation Shock.
- Barrero, J. M., Bloom, N., and Davis, S. J. (2021). Why Working from Home Will Stick.
- Bartik, A. W., Cullen, Z. B., Glaeser, E. L., Luca, M., and Stanton, C. T. (2020). What Jobs are Being Done at Home During the Covid-19 Crisis? Evidence from Firm-Level Surveys.
- Bick, A., Blandin, A., and Mertens, K. (2022). Work from Home Before and After the COVID-19 Outbreak.
- Bloom, N., Liang, J., Roberts, J., and Ying, Z. J. (2015). Does Working from Home Work? Evidence from a Chinese Experiment. *The Quarterly Journal of Economics*, 130(1):165–

- Brown, T., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Brynjolfsson, E., Horton, J. J., Ozimek, A., Rock, D., Sharma, G., and TuYe, H.-Y. (2020). COVID-19 and Remote Work: An Early Look at US Data.
- Burke, M., Sasser Modestino, A., Sadighi, S., Sederberg, R., and Taska, B. (2020). No Longer Qualified? Changes in the Supply and Demand for Skills within Occupations. Federal Reserve Bank of Boston Research Department Working Papers, Federal Reserve Bank of Boston.
- Criscuolo, C., Gal, P., Leidecker, T., Losma, F., and Nicoletti, G. (2021). The role of telework for productivity during and post-COVID-19. (31).
- Davis, S. J., Faberman, R. J., and Haltiwanger, J. C. (2012). Recruiting Intensity during and after the Great Recession: National and Industry Evidence. *American Economic Review*, 102(3):584–588.
- Davis, S. J. and Samaniego de la Parra, B. (2020). Application Flows.
- del Rio-Chanona, R. M., Mealy, P., Pichler, A., Lafond, F., and Farmer, J. D. (2020). Supply and demand shocks in the COVID-19 pandemic: An industry and occupation perspective. *Oxford Review of Economic Policy*, 36(Supplement_1):S94–S137.
- Deming, D. and Kahn, L. B. (2018). Skill Requirements across Firms and Labor Markets: Evidence from Job Postings for Professionals. *Journal of Labor Economics*, 36(S1):S337–S369.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dingel, J. I. and Neiman, B. (2020). How many jobs can be done at home? *Journal of Public Economics*, 189:104235.
- Draca, M., Duchini, E., Rathelot, R., Turrell, A., and Vattuone, G. (2022). Revolution in Progress? The Rise of Remote Work in the UK.
- Forsythe, E., Kahn, L. B., Lange, F., and Wiczer, D. (2020). Labor demand in the time of COVID-19: Evidence from vacancy postings and UI claims. *Journal of Public Economics*,

189:104238.

- Hershbein, B. and Kahn, L. B. (2018). Do Recessions Accelerate Routine-Biased Technological Change? Evidence from Vacancy Postings. *American Economic Review*, 108(7):1737–1772.
- Luktevich, B. (2022). BERT Language Model. <https://www.techtarget.com/searchenterpriseai/definition/BERT-language-model>.
- Marinescu, I. and Wolthoff, R. (2020). Opening the Black Box of the Matching Function: The Power of Words. *Journal of Labor Economics*, 38(2):535–568.
- Mas, A. and Pallais, A. (2017). Valuing Alternative Work Arrangements. *American Economic Review*, 107(12):3722–3759.
- Modestino, A. S., Shoag, D., and Ballance, J. (2016). Downskilling: Changes in employer skill requirements over the business cycle. *Labour Economics*, 41:333–347.
- Mongey, S., Pilossoph, L., and Weinberg, A. (2021). Which workers bear the burden of social distancing? *The Journal of Economic Inequality*, 19(3):509–526.
- Ozimek, A. (2020). The Future of Remote Work.
- Phuong, M. and Hutter, M. (2022). Formal Algorithms for Transformers.
- Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2020). DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter. *arXiv:1910.01108 [cs]*.
- Shapiro, A. H., Sudhof, M., and Wilson, D. J. (2022). Measuring news sentiment. *Journal of Econometrics*, 228(2):221–243.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is All you Need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Wiswall, M. and Zafar, B. (2018). Preference for the Workplace, Investment in Human Capital, and Gender*. *The Quarterly Journal of Economics*, 133(1):457–507.

TABLES AND FIGURES

Table 1: Counts of Vacancy Postings, Employers, and Cities, January 2014 to January 2023

(1)	(2)	(3)	(4)
Country	Vacancies	Employers	Cities
New Zealand	1,700,523	36,201	67
Australia	8,607,160	197,870	59
Canada	11,711,357	712,577	3,691
United Kingdom	74,576,747	876,103	2,268
United States	161,872,915	3,485,630	31,635
Total	258,468,702	5,308,381	37,720

Note: Reported counts pertain to the universe of online postings from January 2019 onwards and a 5% random sample from 2014 to 2018, after we drop about 6% of the postings in the data-cleaning steps described in Appendix A. We rely on Lightcast’s proprietary algorithm to identify employers and cities.

Table 2: Examples of Classification Errors in Dictionary Methods

False-Positive Examples:	False-Negative Examples:
<p>We are looking for a Deputy Home Manager with domiciliary care experience to join our company. You will work from home care facilities with a strong track record of quality service.</p>	<p>We encourage our people to explore new ways of working - including part-time, job-share or working from different kinds of locations, including their home. Everyone can ask about it.</p>
<p>Schedule: * 10 Hour Shift * 8 Hour Shift Work remotely: * No</p>	<p>With a hybrid mix of time at home as well as our corporate office, this role will suit an analytical, process orientated and people focused payroll professional who thrives in a fast-paced environment.</p>
<p>Applicants must also have: * Ability to work as part of a team, in a fast paced environment * Experience in a 4 or 5 star hotel * Previous experience working in remote locations</p>	<p>We see the value in work-life balance, so whether you like to get a surf in before work, like to head home in time to pick up the kids or you just like working from the comfort of your own home now and then, we want to support you.</p>
<p>You may work on renovation projects, store reorganizations, new store openings, and store closings. May respond to managerial or Home Office requests for special reports, information, or for help on special projects.</p>	<p>The interviews for this role are likely to be conducted remotely using Microsoft Teams or Zoom. It is also expected that relevant work within these roles may be done remotely, within the UK.</p>

Note: The left column provides examples of how a dictionary method falsely classifies a vacancy posting as saying the job allows remote work. The right column shows examples of how it falsely classifies a vacancy posting as not saying that the job allows remote work. Bold font designates dictionary keywords, and yellow shading highlights text that helps determine a correct classification. These examples are based on actual vacancy postings in our dataset and the dictionary used in Adrjan et al. (2021).

Table 3: WHAM Attention Weights Compared to Dictionary Keywords

WHAM View:

Schedule:

** 10 Hour Shift

** 8 Hour Shift

Work remotely:

** No

We are looking for a Deputy Home Manager with domiciliary care experience to join our company. You will work from home care facilities with a strong track record of quality service.

We encourage our people to explore new ways of working - including part-time, job-share or working from different kinds of locations, including their home. Everyone can ask about it.

Dictionary View:

Schedule:

** 10 Hour Shift

** 8 Hour Shift

Work remotely:

** No

We are looking for a Deputy Home Manager with domiciliary care experience to join our company. You will work from home care facilities with a strong track record of quality service.

We encourage our people to explore new ways of working - including part-time, job-share or working from different kinds of locations, including their home. Everyone can ask about it.

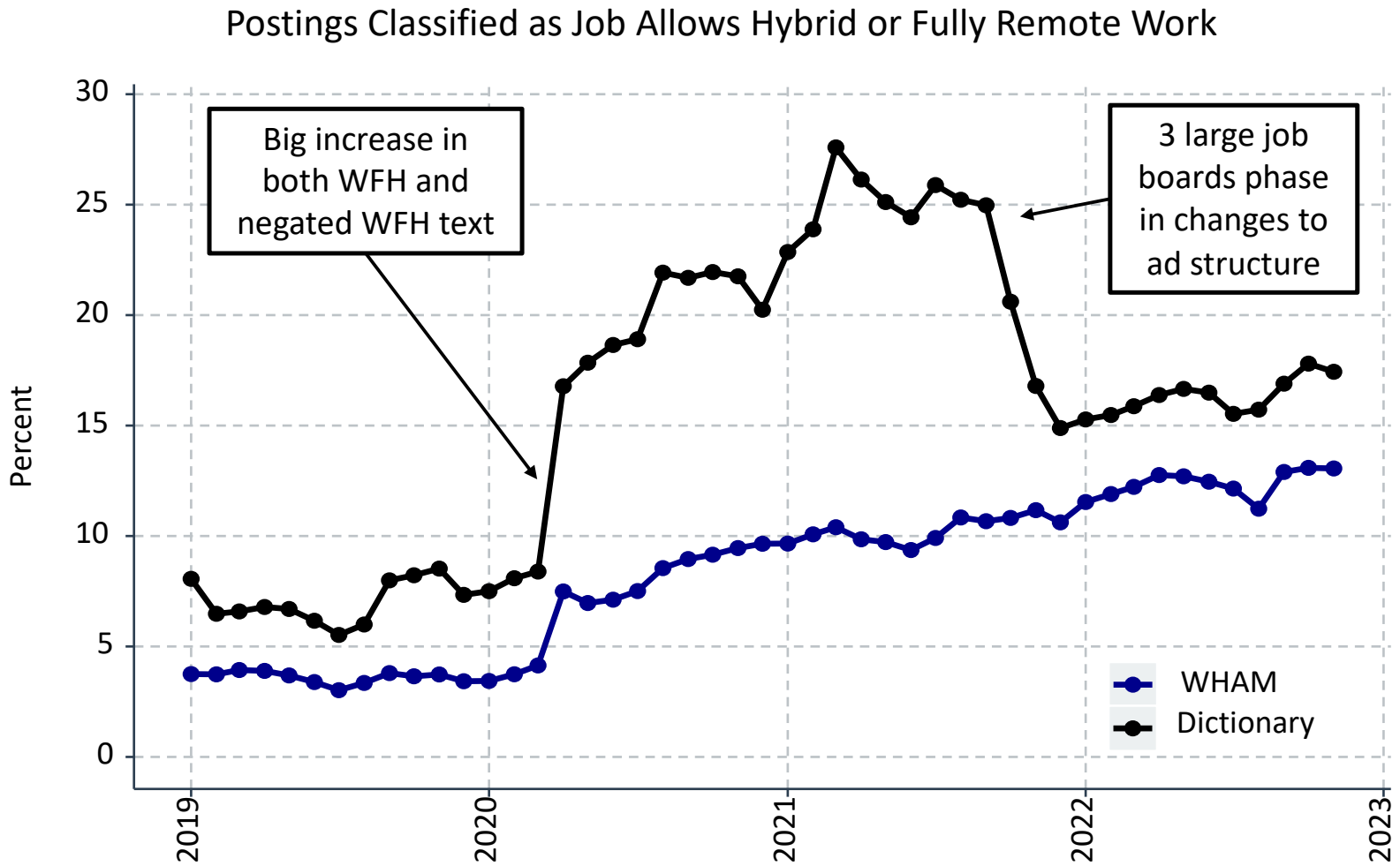
Note: The left column illustrates the role of attention weights in WHAM classifications of vacancy postings, where darker shadings pertain to higher weights. The right column illustrates the application of dictionary methods to the same text passages, where highlight text pertains to keywords.

Table 4: WHAM Outperforms Other Classification Methods

	Audit Sample			Approximate Random Sample		
	(1)	(2)	(3)	(4)	(5)	(6)
	Error Rate	Precision	F1 Score	Error Rate	Precision	F1 Score
All Zero	.28	.00	.00	.03	.00	.00
Dictionary	.16	.68	.74	.14	.15	.25
Dictionary w/ Negation	.12	.82	.78	.07	.28	.40
Logistic Regression	.11	.81	.81	.07	.26	.40
Logistic Regression w/ Negation	.08	.87	.85	.05	.36	.50
GPT-3	.06	.87	.89	.05	.36	.52
WHAM (Generic English)	.03	.95	.95	.02	.66	.78
WHAM (Baseline)	.02	.97	.97	.01	.75	.85

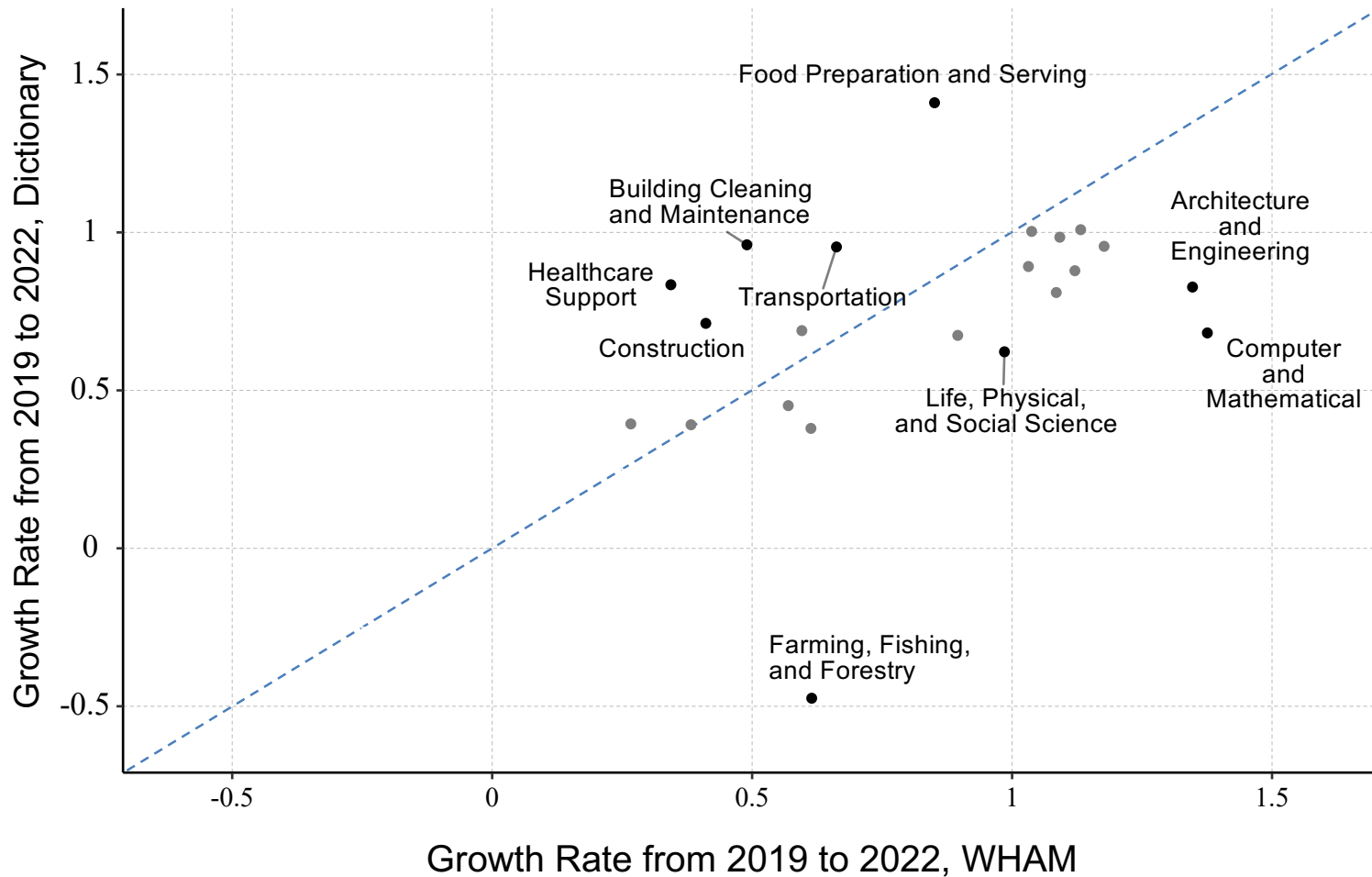
Note: This table reports classification performance metrics, which we calculate using a hold-out sample of human-classified text sequences. “Error Rate” is the overall rate of misclassifications (relative to humans). “Precision” is the ratio of true-positive classifications to the sum of true positives and false positives. “F1 score” is the harmonic mean of Precision and “Recall”, where Recall is the fraction of true positives divided by the sum of true positives and false negatives – i.e., the denominator is the true number of positives, according to human classifications. Columns (1)-(3) uses a 40% random subset of our audit sample, and Columns (4)-(6) uses a sample that approximates a random sample of our full universe of postings. See Appendix B for details, including a description of each algorithm.

Figure 1: WHAM and Dictionary Methods Applied to U.S. Vacancy Postings



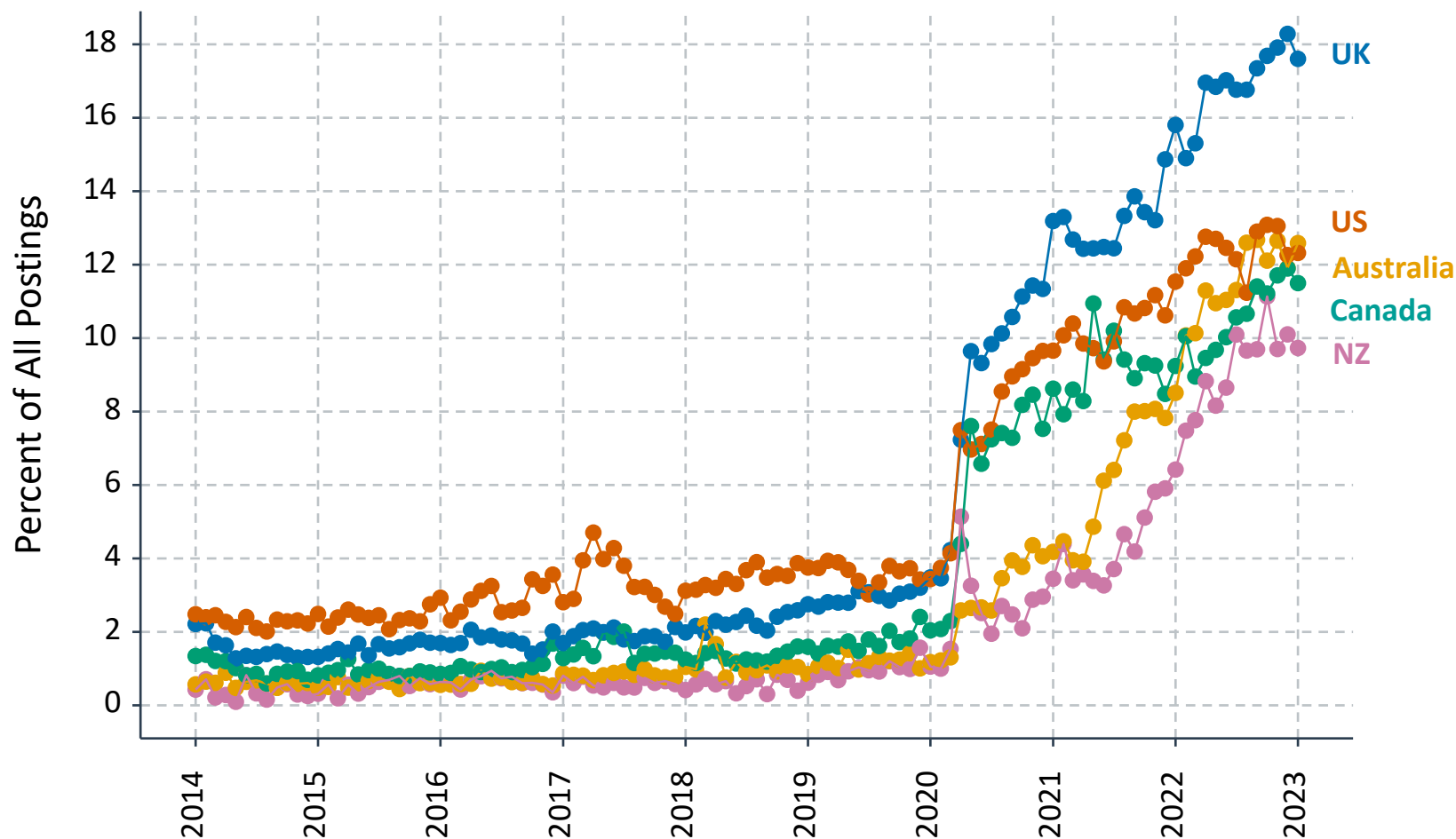
Note: This figure shows the percent of postings that say the job allows one or more remote workdays per week, as classified by WHAM (blue) and a dictionary-based approach (black) using the keywords in Adrjan et al. (2021). For both methods, we reweight the data to match the U.S. occupational distribution of vacancies in 2019 at the six-digit SOC level.

Figure 2: Share of U.S. Postings that Allow Some Remote Work, Growth Rate by Two-Digit Occupations, WHAM Compared to Dictionary Method



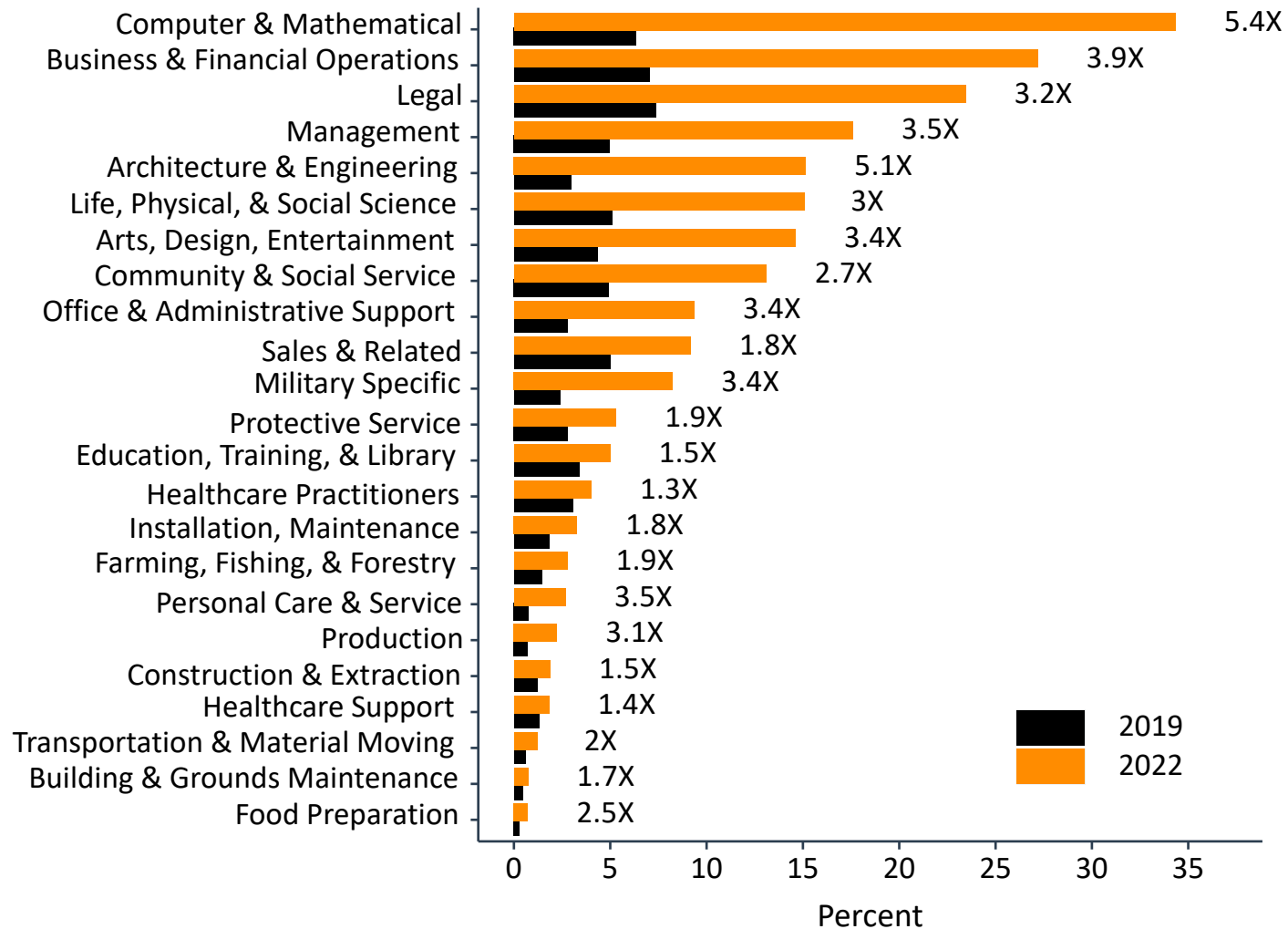
Note: We sort postings into Standard Occupational Classifications (SOC) at the two-digit level and calculate the share of postings that say the job allows for one or more days per week of remote work in 2019 and 2022. We then calculate the DHS growth rate from 2019 to 2022 as $(X_{2022} - X_{2019}) / 0.5 * (X_{2019} + X_{2022})$. For the dictionary method, we use the keywords in Adrjan et al. (2021). The blue-dashed line shows a 45 degree line.

Figure 3: Vacancy Postings that Explicitly Offer Hybrid or Fully Remote Work Rose Sharply in All Five Countries from 2020



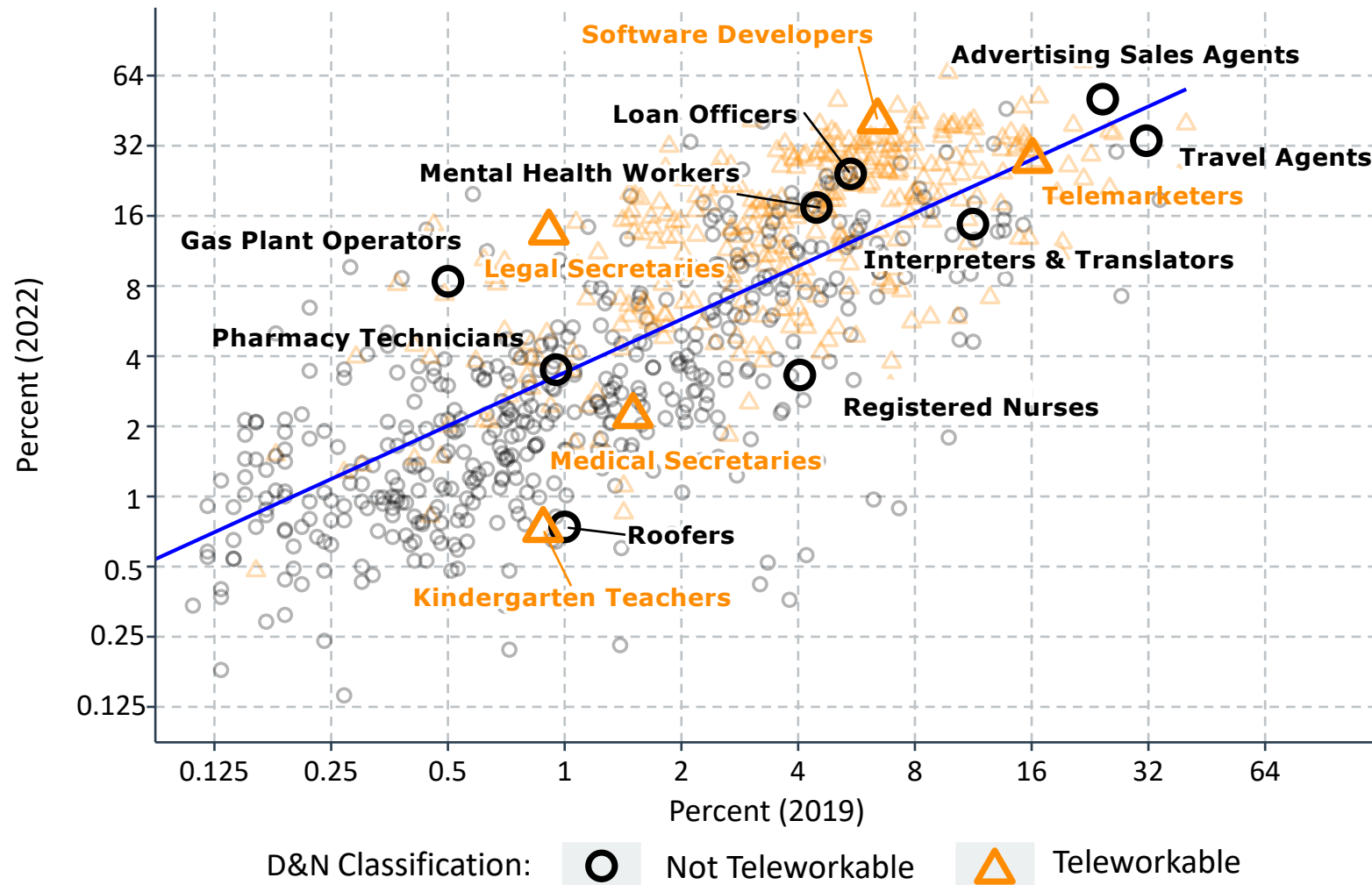
Note: This figure shows the percent of vacancy postings that say the job allows one or more remote workdays per week, encompassing both hybrid and fully-remote working arrangements). We compute these monthly, country-level shares as the weighted mean of the own-country occupation-level shares, with weights given by the U.S. vacancy distribution in 2019. Our occupation-level granularity is roughly equivalent to six-digit SOC codes. See Appendix B for the corresponding raw series and series based on alternative weighting schemes.

Figure 4: Professional, Scientific and Computer-Related Occupations Have the Highest Shares of Postings that Offer Hybrid or Fully-Remote Work, U.S. Data



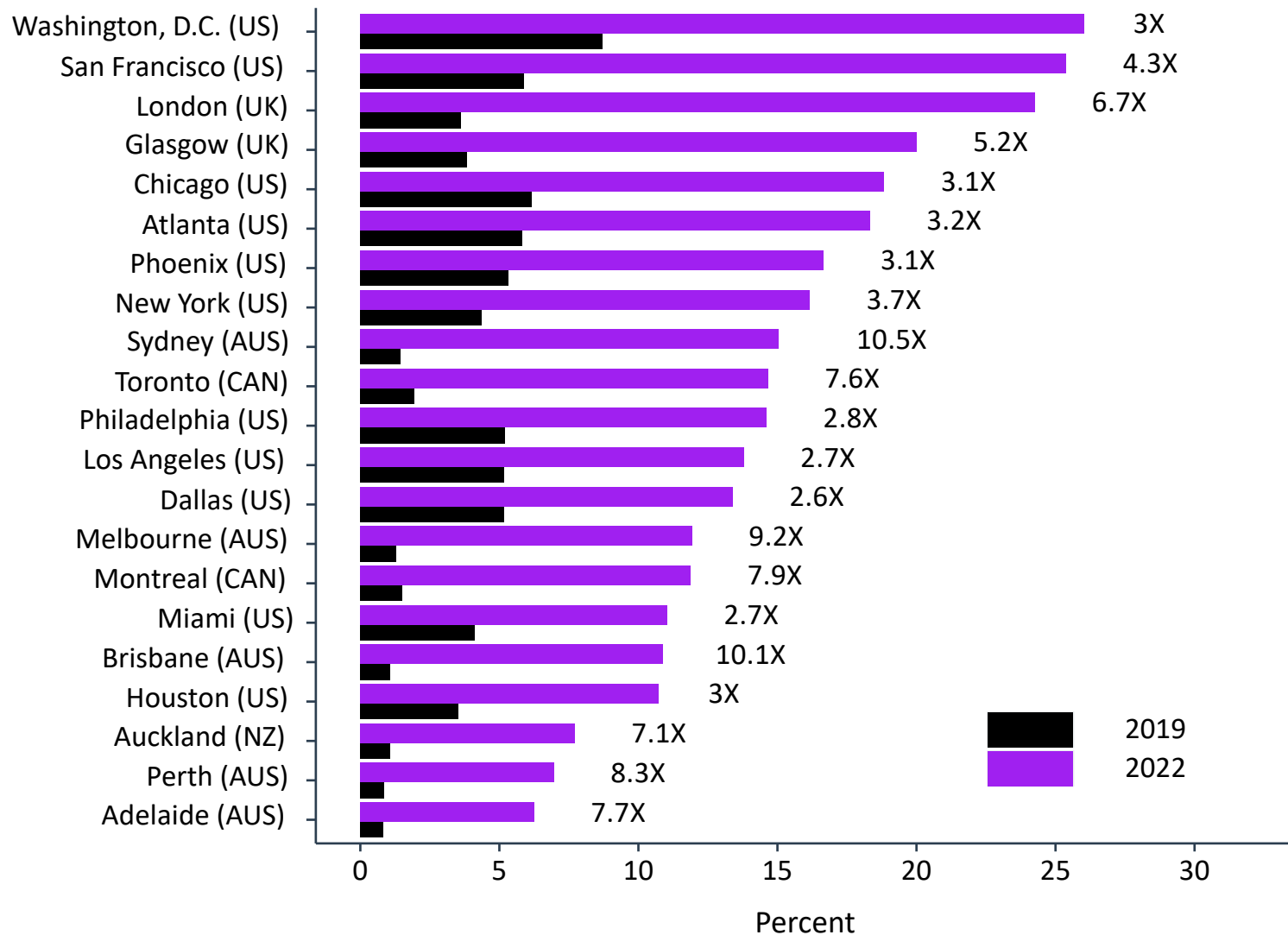
Note: Each bar reports the percent of vacancy postings that say the job allows one or more remote workdays per week in the indicated period and occupation group (two-digit SOC).

Figure 5: The Share of Vacancy Postings that Explicitly Offer Hybrid or Fully Remote Work Rose in Almost Every Occupation, U.S. Data



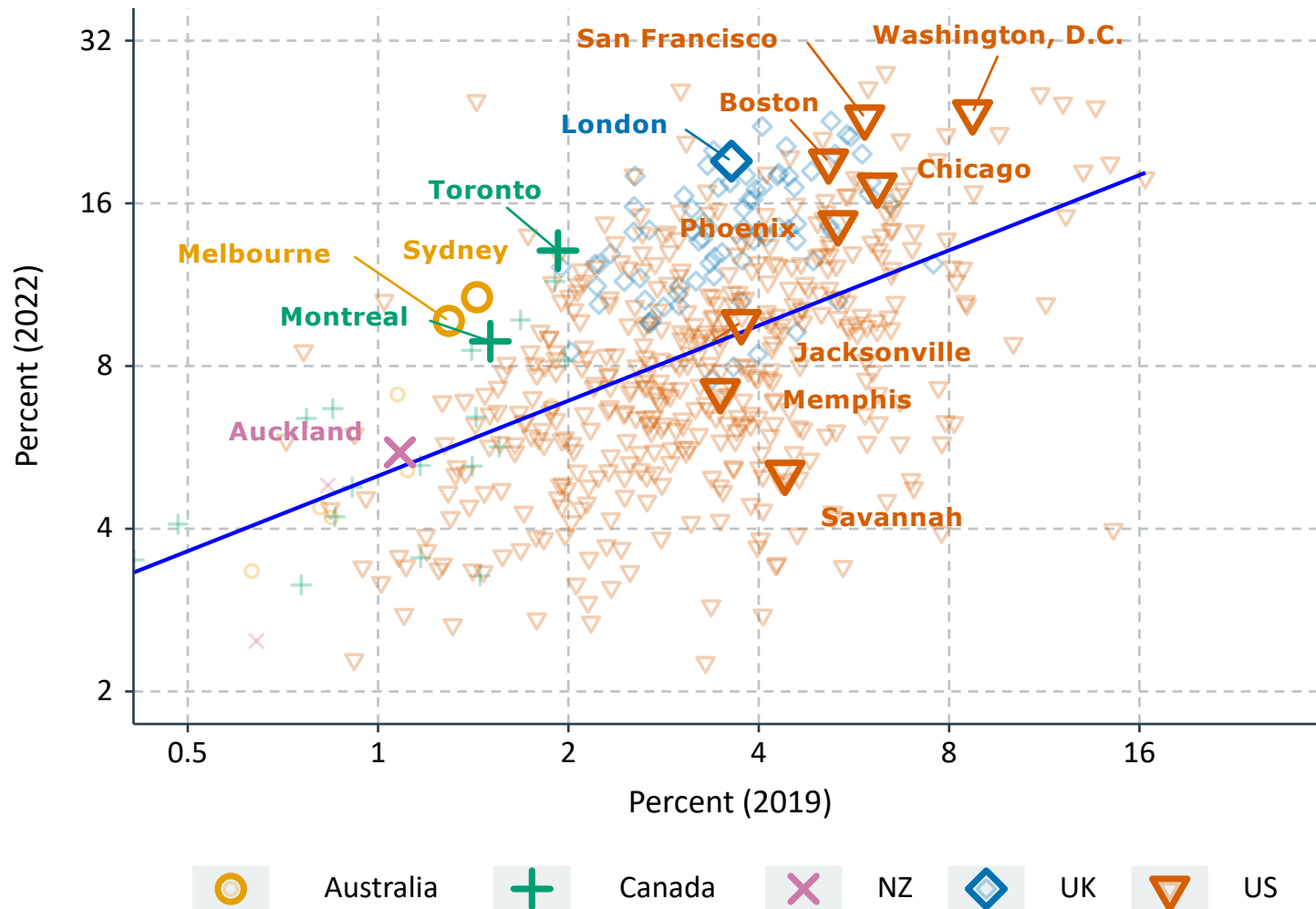
Note: This figure plots the percent of postings that say the job allows one or more remote workdays per week for 875 occupations in 2019 and 2022. We define occupations by ONET codes, dropping those with fewer than 250 posting in 2019. The line shows the unweighted OLS fit: $\log(y) = 1.22 + 0.76 \log(x)$, which has an R^2 value of 0.63. The color and shape denote whether Dingle & Neiman (2020) classify the occupation as feasible for fully remote working.

Figure 6: The Share of Vacancy Postings that Explicitly Offer Hybrid or Fully Remote Work Varies Widely across Major Cities



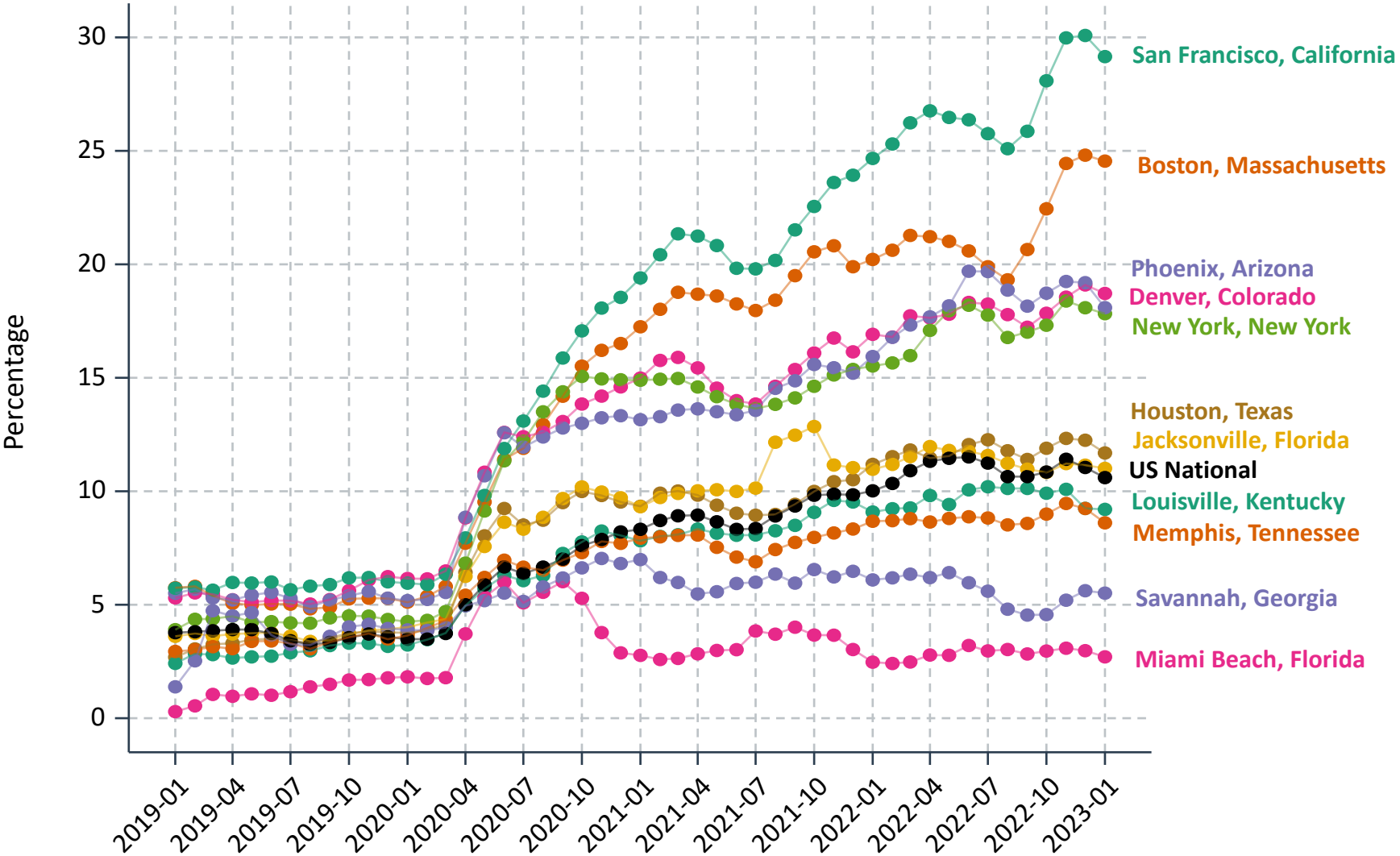
Note: Each bar reports the percent of vacancy postings that say the job allows one or more remote workdays per week in the indicated period and city. City refers to the location of the establishment or firm that is hiring.

Figure 7: The Share of Vacancy Postings that Explicitly Offer Hybrid or Fully Remote Work Grew at Different Rates across Cities since the Pandemic



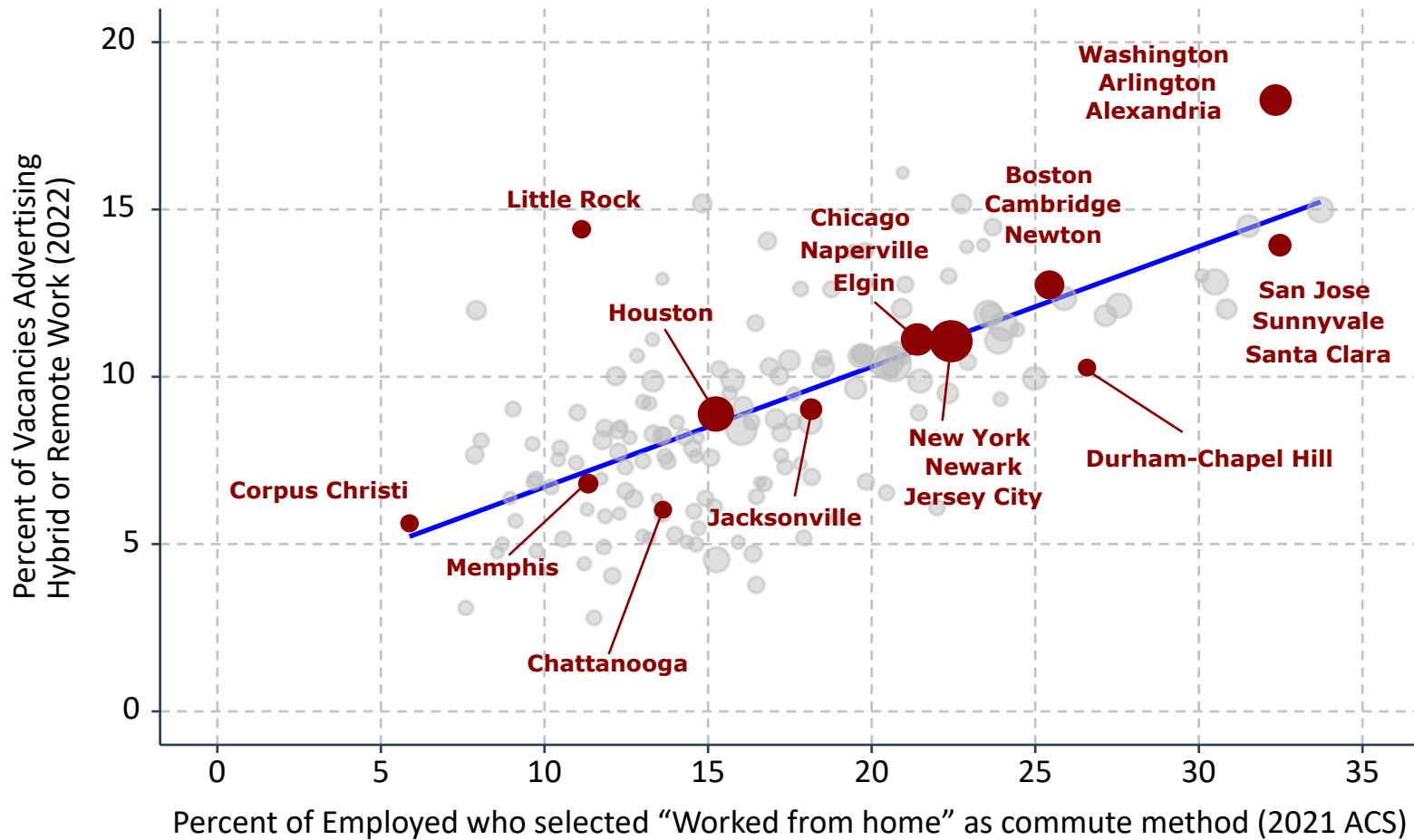
Note: This figure plots the city-level percent of postings that say the job allows one or more remote workdays per week in 2019 and 2022. “City” refers to the location of the establishment or firm that is hiring. The line shows the unweighted OLS fit: $\log(y) = 1.61 + 0.46 \log(x)$, which has an R^2 value of 0.28.

Figure 8: Share of Postings Offering Hybrid or Fully Remote Work vary across US cities



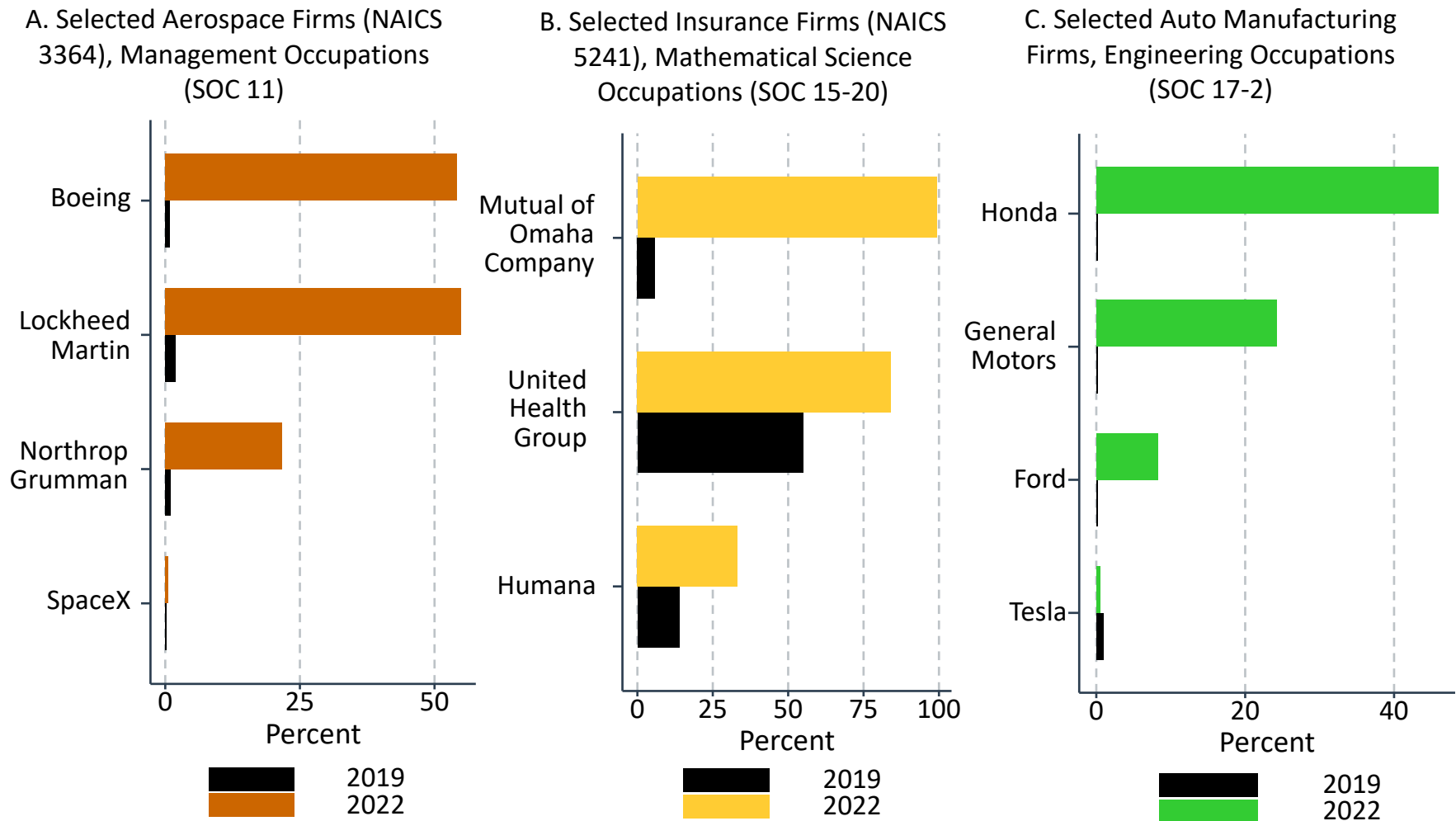
Note: We calculate the monthly share of all new job vacancy postings which explicitly advertise remote working arrangements (i.e. both hybrid and fully-remote), by selected cities. Prior to aggregation at the monthly level, we employ a jackknife filter to remove a small number of outlier days (see Appendix A: Data for further details). This figure shows the 3-month moving average. Cities chosen above are selected examples to illustrate the wide cross-city spread.

Figure 9: Share of Vacancy Postings Offering Hybrid or Fully Remote Work Compared to Share of Employed that Designate “Worked from home” as commute method, U.S. Metropolitan Statistical Areas



Note: The vertical scale is the percent of postings in 2022 that say the job allows one or more remote workdays per week (i.e. both hybrid and fully-remote). The horizontal scale is the percent of employees who select “Worked from home” as their commute method in 2021 in the American Communities Survey (ACS). ACS respondents are instructed to “Mark (X) ONE box for the method of transportation used for most of the distance,” which suggests that only those who work in a fully-remote capacity should select this box. (Persons with 1+ days of commute per week have more mileage from that commute mode.) The line shows the unweighted OLS fit: $\log(y) = 3.12 + 0.36 \log(x)$, which has an R^2 value of 0.55. The regression includes one observation that is outside the plotted axes.

Figure 10: The Prevalence of Postings that Allow Hybrid or Fully-Remote Work Varies Greatly, even among Same-Industry Firms Recruiting in the Same Occupational Category



Note: For each firm, year and indicated occupation, we report the percent of U.S. postings that say the job allows one or more remote workdays per week.

ONLINE APPENDIX

A Data Appendix

In this Appendix we provide further commentary on the corpus of online job vacancy postings.

A.1 Data Provider

Our corpus of online job vacancy postings is provided by the labour market and analytics company ‘Lightcast’ (formerly Emsi Burning Glass). Lightcast has been scraping online job vacancy postings in the USA since 2007, and has continued to expand to other countries.

A.2 Web Sources

Each job vacancy posting is scraped by Lightcast from the internet. Specifically, the company scrapes over 50,000 web sources. These sources include private online job vacancy aggregators (e.g. Indeed.com, Monster.com), public online job boards (e.g. New York City Department of Labour’s ‘JobZone’), and employers’ own recruitment web pages (e.g. careers.microsoft.com, usajobs.gov). Lightcast actively audits their list of web sources to ensure data from new websites is on-boarded in a timely manor.³² One of the main competitive advantages of Lightcast’s data product is the breadth of their sources. These data are often referred to in the literature as the ‘near universe’ of online job vacancy postings.

A.3 What’s in the job vacancy posting data?

Once an online job vacancy posting is scraped, Lightcast processes this data to produce three categories of information: (i) plain text, (ii) meta data, and (iii) structured data. A description of each of these categories follows presently:

A.3.1 Plain Text

The plain text of each job ad contains the full textual description of the job, as written by employers. To construct this, Lightcast takes the HTML file scraped from a given website and does two further processing steps. First, it parses out portions of the HTML file which do not contain information about the vacancy (e.g. removing website headers, footers, and side-menu bars). Second, Lightcast takes this portion of HTML which (ideally) contains only information about the job vacancy, and turns it HTML into plain-text.

A.3.2 Meta Data

Each vacancy posting also contains a number of meta-data items. These are immutable properties of each web scraped vacancy. The most important of these is the date the page was scraped. Another important piece of meta-data is the URL from which the posting was scraped.

³²One reason we eschew analysis of the count of postings and instead focus on shares is that the underlying donor pool of online sources is constantly changing.

A.3.3 Structured Data

The most commonly used data product that Lightcast creates is the set of structured data. This dataset contains one row for each job vacancy posting, and a large number of additional information such as the job title, occupation, salary, educational requirements, location, and employer name. These variables are extracted using Lightcast’s own proprietary algorithms. These fields differ from meta data because they may contain missing values and/or measurement error due to imperfect algorithmic extraction.

A.4 Errors and Missing Information

Overall, the data product is a highly informative and accurate product. We view the incidence of errors as very minute, but acknowledge that any dataset with hundreds of millions of observations scraped from over 50,000 sources will never be perfect. Both the structured data and the plain text data require a number of pre-processing steps and the use of algorithmic feature extraction, which in a very small number of cases produce errors (e.g. misclassification of occupations, truncation of plain text, presence of erroneous text). In this subsection we highlight some of the errors we have encountered, and discuss the strategies we employed to ensure our results remain robust to such issues.

A.4.1 Missing Values

A specific value (e.g. the educational requirement for a job) might be missing for at least two reasons: (i) the employer does not mention this explicitly in the text of the job ad, and (ii) the algorithm used to extract this feature from the text failed. The former issue is especially problematic in the context of educational requirements (e.g. we see that very few vacancies for Medical Doctors explicitly mention a requirement to have gone to medical school). This is because certain features of the job will likely be taken as given (for example, specialized degrees for medical doctors). We also see that a large share of vacancy postings do not list the salary (this is almost entirely due to lack of information, and not poor feature extraction). One could employ imputation methods to address these missing values (see Bana (2022), who predict the salary with a very high degree of accuracy from the text). The main strategy employed in this paper was to only utilise covariates which contain fewer missing values, such as occupation classifications and location information.

A.4.2 Erroneous Plain Text

In a very small number of cases we observe that the plain text includes some parts of the website other than the job description. For example, the plain text from one job board in New Zealand included a number of vacancy posting text from ads that were being cross-promoted to the browser, essentially turning each document into a compilation of six job ads.

A.4.3 Truncated Plain Text

In a small number of other cases, the plain text is truncated. For example, we found one employer who listed each jobs location using an interactive link which must be clicked to appear. Since the web scraper only parses static information, this portion of the job ad was missing from the plain text. We conducted extensive tests, and stress that in the vast majority of cases the plain text provides an accurate representation of the job vacancy posting.

A.5 Checking for Correlated Measurement Error

As discussed above, since our measurement of remote working relies on the underlying plain text, some measurement error is inevitable. One concern we took seriously is the possibility that this noise may be correlated within online sources. We discuss our approach to addressing this below.

A.5.1 Validation of Large Job Boards

In some instances the pre-processing of job posting websites includes additional erroneous text. When this occurs, it is very likely to be true for all job postings scraped from the same web source. To ensure our results are not overly sensitive to such issues, we first identify the twenty largest web sources for each country. We then create twenty versions of country-level time series of monthly remote work vacancy shares, leaving one job board out at a time. This process revealed one problematic source from each of Canada, USA, NZ and UK. We found two problematic job boards in Canada. Table A.1 reports the fraction of total job ads that were removed after dropping postings from these sources.

A.5.2 Outlier Detection and the Jack-Knife Filter

When we present monthly time series data, we apply an algorithm which filters outlier days whose contribution to the over-all monthly share of vacancy postings offering remote work is at odds with other days in a given month. This filter has a very minimal impact on the results (e.g. we drop less than a quarter of one percent of job postings from the US based on this filter). The few outlier days we do filter out occur when a large number of vacancies get posted on a single day which are concentrated by employer/occupation/web source. Our extensive audits of the data reveal that outlier days are due to compositional discontinuities at the daily frequency, and not caused by measurement error in our algorithm. Our filter is based on the Jack-Knife resampling procedure, and works as follows:

- For a given calendar month M denote S_M as the share of vacancy postings which offer remote work
- For each day $t \in M$, compute the share of remote work postings *excluding* all postings on this focal day t from the calculation. Define this share as $S_{M \setminus \{t\}}$

- If the absolute level deviation between S_M and $S_{M \setminus \{t\}}$ is greater than 2 percentage points, or else if the absolute ratio of their natural logarithms is greater than 0.1, then we classify focal day t as an outlier
- Recalculate the share of vacancy postings for month M excluding all postings on outlier days

This filter alters the data minimally. For example, in the United States, it removes 0.2% of the total number of vacancy postings. The number of postings which are filtered is shown in the below table:

A.6 Representativeness of Online Job Vacancy Postings

Lightcast frequently reviews the representativeness of the job vacancy postings it scrapes, to ensure the information renders an accurate picture of the overall labour market. Both our analysis and that of our data provider, as well as many other papers in the literature who utilise these data, all find a high degree of fidelity between the share of job vacancies across occupations and industries, and other official Government data products which measure similar phenomena.

In our baseline results, we also re-weight the data to reduce sensitivity to shifts in the overall composition of the labour market. The next section discusses this further, but we note that this provides additional robustness to concerns of representativeness

B Supplementary Results

See figures below.

C Supplementary Information on Measurement

C.1 Estimation details for WHAM

WHAM builds from DistilBERT (Sanh et al., 2020), which has a Transformer architecture with six layers and 66 million parameters. It was originally estimated to predict randomly deleted words in a corpus of unpublished books and all English Wikipedia. We use the uncased version of the model.

The first estimation step in WHAM is to pre-train off-the-shelf DistilBERT (Sanh et al., 2020) to predict randomly deleted words in a random sample of 900,000 job posting sequences which is balanced across all years and countries. The total fraction of deleted words is 15%. We use guidelines from the original BERT paper (Devlin et al., 2019) to select the hyperparameters for estimation: a batch size of 8, three training epochs, and a learning rate of $5e-5$.³³

The second estimation step in WHAM is to fine-tune the model to predict human labels. To select the estimation hyperparameters, we use three-fold cross validation and the training data used for the benchmark exercises reported in section 3. We perform an exhaustive search over learning rates $\{2 * 10^{-5}, 3 * 10^{-5}, 5 * 10^{-5}\}$, epochs $\{2, 3, 5\}$ and batch sizes $\{16, 32\}$, and select the set of hyperparameters with the highest average F1 score across training data splits. The resulting choices are $5e-5$, 2, and 16, respectively. The model estimated with these choices solely on the training data is used to determine the test-set performance reported in section 3. To produce output on the entire dataset, we re-estimate the model using all human labels (from both training and test sets) using the same hyperparameters and use this model to predicted remote work on all sequences in the Lightcast data.

C.2 Details for other classification approaches

Section 3 compares various alternatives to WHAM for classifying remote work, and here we provide additional details on these.

C.2.1 Dictionary

We implement the dictionary approach with the following steps:

1. **Preprocessing:** We lowercase all text, remove punctuation symbols (except for hyphens and apostrophes), remove numbers, and replace all sequences of white spaces with a single white space.
2. **Tagging:** We search for the appearance of any of the keyword phrases from Table C.1. For phrases containing multiple words (e.g. ‘work from home’) we allow for any arbitrary combination of white spaces and hyphens separating the words that compose the dictionary keyword (e.g. ‘work-from-home’, ‘work- from- home’).

³³The batch size determines how many text sequences are processed at each step in estimation. The number of epochs determines the total number of times the entire data set is passed through in estimation. The learning rate determines the speed at which the model parameters update in gradient descent.

3. **Binary classification:** Any job posting that contains a match to any of the dictionary keywords is classified as positive.

C.2.2 Negation adjustment

Our strategy for negation adjustment follows that proposed by Shapiro et al. (2022) to capture negation in the context of sentiment analysis. For every keyword match from the dictionary within a job posting, we consider it to be negated if any of the following is true:

1. There is a negation term in any of the three words before the keyword match. The set of negation terms is displayed in Table C.2, and comes from the VADER Sentiment Analysis toolkit.
2. “no” or “not” appear in the two words after the keyword match
3. A word that contains “n’t” is the immediate word after the keyword match

If a job posting is negated we then change its binary label from positive to negative.

C.2.3 Logistic regression

Our approach to logistic regression follows the approach in Adams-Prassl et al. (2020). We start by applying the same pre-processing steps used for the dictionary approach to the job postings: i) lowercase text, ii) remove punctuation (except for the hyphen), iii) remove numbers, and iv) clean white spaces. Next we split the text into individual tokens and build the document-term frequency matrix by using the 5,000 most common tokens. For each keyword in our dictionary (the phrases in Table Table C.1) that is not part of the 5,000 most common tokens, we add a column in the document-term matrix with its counts. Finally, we transform the matrix into its binary form; every entry above one is replaced with a one. This matrix then becomes the set of covariates used to predict human labels via logistic regression with L_1 regularization (LASSO). To determine the LASSO penalty, we use five-fold cross-validation on the training data, and select the regularization parameter that achieves the highest average $F1$ score across the five splits.

C.2.4 Logistic regression with negation

We follow an identical procedure to the one described for logistic regression but we further extend the document-term matrix with one extra column per keyword in the dictionary that indicates that the keyword was negated (according to our negation adjustment described before).

C.2.5 GPT-3

We use OpenAI’s GPT-3 model to generate predictions on the presence of remote work in our job postings. To do this, we craft a simple prompt that instructs the model to *“Identify if the text offers the possibility of remote work at least one day per week”* and ask the model to

generate an answer. Figure C.1 illustrates a particular example using OpenAI’s Playground.

In most cases, the text generated by GPT-3 is a *Yes/No* answer. Sometimes, however, the model generates longer answers (e.g. “temporarily due to covid”). In order to transform these answers into a binary prediction we do the following: i) lowercase the answer of GPT-3 and clean any additional white spaces and ii) if the answer contains “no” as part of its three first characters we assign a zero (no remote work) to the sequence, else we give it a one (remote work).

We test both the ‘text-davinci-02’ model and the ‘text-davinci-03’ model using the same prompt and report performance of the former given its lower error rate with respect to our human labels.

Figure A.1: Fraction of job vacancy postings we drop from sample due to problematic job boards

(1)	(2)	(3)
Source	Country	Percent of Raw Data Dropped
A	USA	6.7
B	UK	3.6
C	NZ	28.9
D	Canada	3.9
E	Canada	3.5

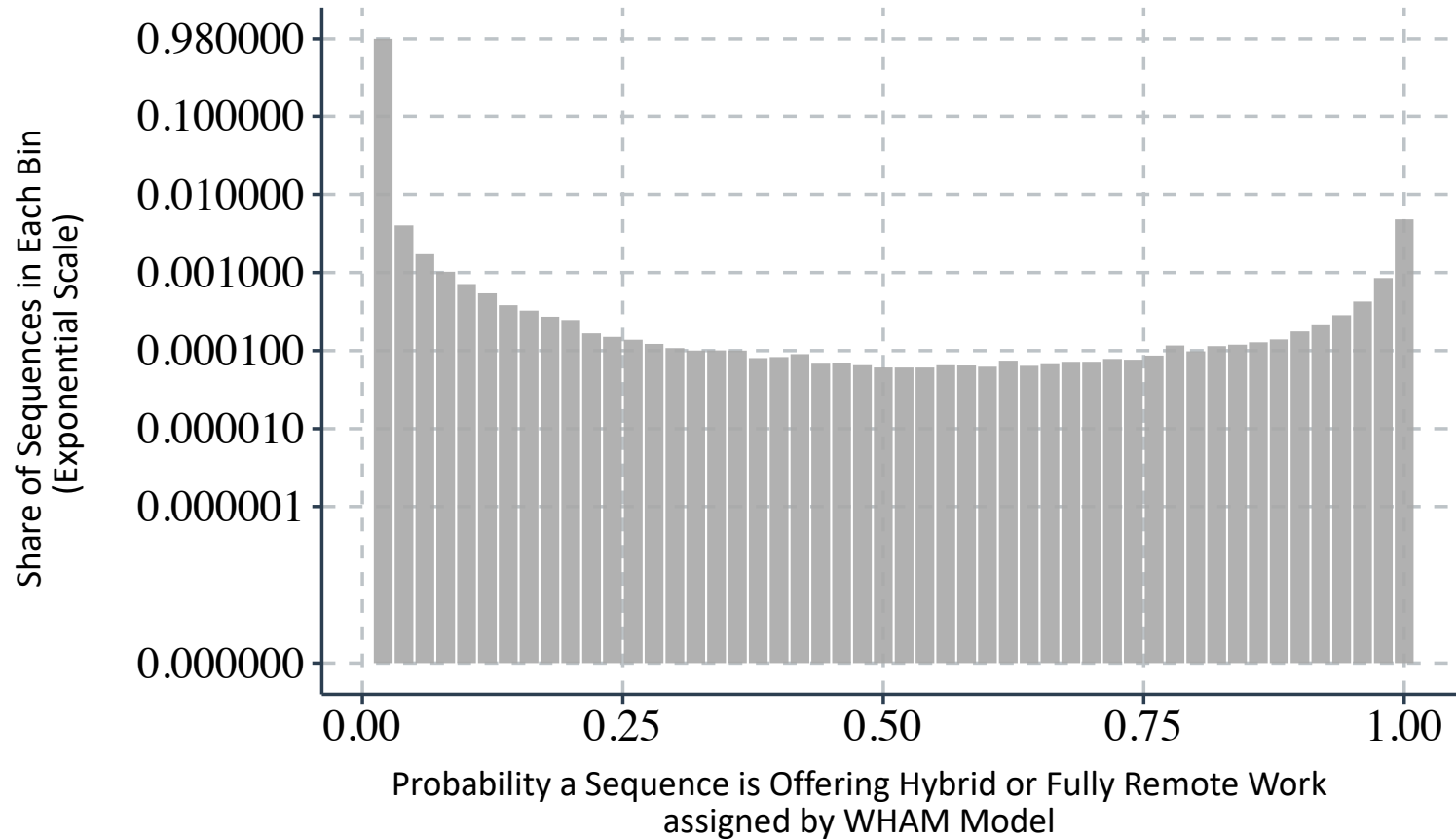
Note: We do not identify these job boards, and instead refer to each by a letter e.g. “A”. This is to avoid any potential conflict with the commercial interests of these websites, though it should be noted that this is not a failure of the job board but rather an issue with our web scraped data. Researchers should contact us if they would like to know the names of each source which we drop. Column (3) reports the fraction of all postings from Jan 2014 to Jan 2023, within the relevant country, that is removed from the raw data after we dropped the corresponding source.

Figure A.2: Jack-Knife Time Series Filter Removes a very small amount of job postings from monthly time series results

(1)	(2)
Country	Percent of Raw Data Dropped
NZ	7.74
Australia	0.78
Canada	1.43
United Kingdom	0.06
United States	0.21

Note: This table records the total fraction of all online job vacancy postings which are dropped from our sample after applying our jackknife filter. This is only imposed when we present monthly frequency time series plots, and works by identifying outliers at a daily frequency prior to monthly aggregation.

Figure B.1: Most Sequences are Assigned a Predicted Probability by WHAM at Extreme Values



Note: WHAM assigns a predicted probability to each sequence in the full job posting dataset using our trained neural network model. This figure presents a histogram of the share of sequences that fall in different bins according to these predictions.

Table B.1: Most Job Postings Either Have Zero or One Sequence that gets Classified as Offering Hybrid or Fully Remote Work Arrangements

(1)	(2)	(3)
Remote Work Sequences	Number of Vacancy Postings	Share Of Total (%)
0	40,006,052	90.4
1	2,682,844	6.1
2	989,084	2.2
3	365,970	0.8
More than 3	201,523	0.5

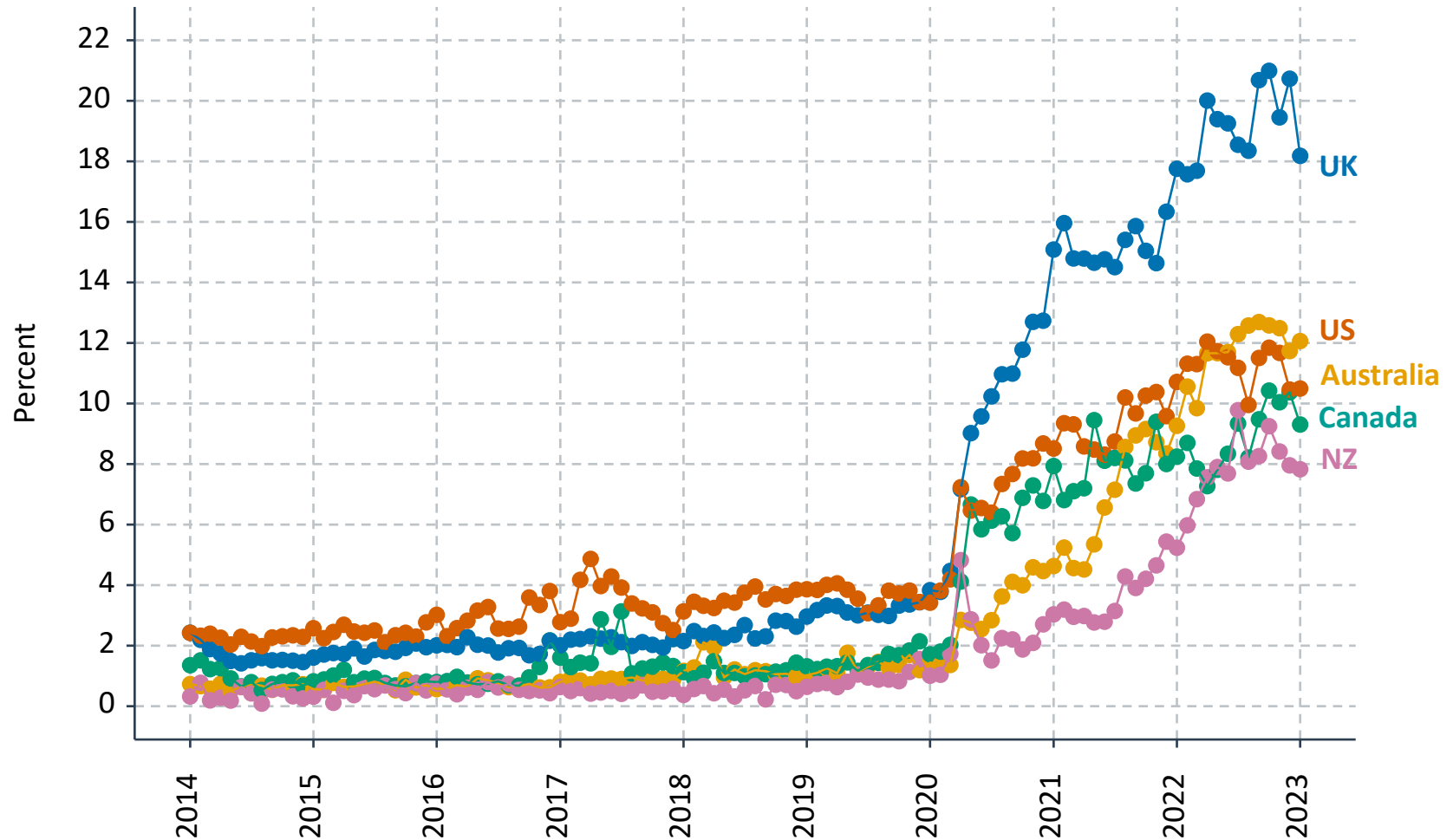
Note: This table tabulates how many text sequences in each US job posting from 2021 are classified as offering remote work (either hybrid or fully remote) according to WHAM. A typical job ad is split into six sequences. Most postings (90.42%) have no positive sequences. Of the remaining fraction, most have only one positive sequence.

Table B.2: Computing Hardware and Costs for the Three Stages of WHAM

	(1)	(2)	(3)
	Pretraining	Fine-Tuning	Full Sample Prediction
Computational setup	GCP Virtual Machine with 1 NVIDIA V100 GPU	GCP Virtual Machine with 1 NVIDIA V100 GPU	GCP Virtual Machine with 8 NVIDIA A100 GPUs
Total time (hours)	12	3	36
Job postings (per hour)	NA	NA	7,000,000
Cost per hour (USD)	\$3	\$3	\$40
Total Cost (USD)	\$36	\$9	\$1,440

Note: This table details the computational setup and time/money costs associated with the different stages of WHAM. All computations were performed on the Google Cloud Platform.

Figure B.2: The raw unweighted share of new job ads offering hybrid or fully remote work is highest in the UK, as UK has very high proportion of ‘white-collar’ jobs being advertised



Note: This figure shows the share of vacancy postings that say the job allows one or more remote workdays per week. We compute these monthly, country-level shares as the raw mean from the universe of new job vacancy postings in each country from each month. Our baseline approach presented in Figure 3 uses vacancy shares from the US to control for occupational composition across countries.

Table C.1: Keywords Used to Implement the Dictionary Approach to Remote-Work Classification

working remotely	working from home	work remotely
work from home	work at home	teleworking
telework	telecommuting	telecommute
smartworking	smart working	remote work teleworking
remote work	remote	remotely
homeoffice	home office	home based
homebased		

Note: These are the keywords that appear in Table A.2 of Adrian et al. (2021) for detecting the presence of remote work in the text of job postings. The three exceptions are `homebased`, `home based`, and `remotely` which we add to the original terms to improve accuracy.

Table C.2: Terms Used for Negation in the Dictionary Approach

aint	arent	cannot	cant	couldnt	darent	didnt	doesnt
ain't	aren't	can't	couldn't	daren't	didn't	doesn't	dont
hadnt	hasnt	havent	isnt	mightnt	mustnt	neither	don't
hadn't	hasn't	haven't	isn't	mightn't	mustn't	neednt	needn't
never	none	nope	nor	not	nothing	nowhere	oughtnt
shant	shouldnt	uhuh	wasnt	werent	oughtn't	shan't	shouldn't
uh-uh	wasn't	weren't	without	wont	wouldnt	won't	wouldn't
rarely	seldom	despite	no				

Note: This is a set of terms that the VADER sentiment analysis tool uses for negation, and which Shapiro et al. (2022) adopt. We add the term 'no' to the baseline negation set.

Figure C.1: GPT-3 Example Prompt

Playground Q&A × ∨ Save View code Share ...

Identify if the text offers the possibility of remote work at least one day per week: 🎤

Text: "We are looking for a Deputy Home Manager with domiciliary care experience to join our company. You will work from home care facilities with a strong track record of quality service."

Remote work: **No**

Submit ↶ ↷ 🔄 🗨 👍 65

Note: To compare our WHAM model to recent advances in generative AI, we analysed our audit sample of text sequences drawn from job vacancy postings using GPT-3. The above is an illustration of the prompt we used. We report the performance of this approach to classification in Table 4 which compares this measurement algorithm to WHAM. The example shown above highlights that, unlike some other widely used methods, this technology is similar to WHAM in its ability to process context.