

Medical Research as a Productivity Indicator

Maya M. Durvasula^{*} Sabri Eyuboglu[†] David M. Ritzwoller[‡]

ABSTRACT. Across fields, the quantity of research has increased substantially, without an attendant increase in output. We argue that, in medicine, this indicator of declining productivity reflects a compositional shift toward low-capital, low-productivity research. Using a fine-tuned, open-source large language model, we construct a novel census of capital-intensive, high-productivity medical investment—clinical trials. Since 2010, the annual quantity of clinical trials has been constant. By contrast, the quantity of other forms of clinical research has increased substantially. Within clinical trials, there is substantial heterogeneity in productivity. The distribution of this heterogeneity is stable over time.

Keywords: Medical Innovation, Productivity, Large Language Models

JEL: C81, O32

Date: May 15, 2024

^{*} Stanford Department of Economics and Stanford Law School, maya.durvasula@stanford.edu

[†] Stanford Department of Computer Science, eyuboglu@stanford.edu

[‡] Stanford Graduate School of Business, ritzvoll@stanford.edu

We thank Nicholas Bloom, Agnes Cameron, Jiafeng Chen, Matthew Gentzkow, Han Hong, Guido Imbens, Ramesh Johari, Charles Jones, Evan Munro, Lisa Larrimore Ouellette, Christopher Ré, Joseph Romano, Brad Ross, Bhaven Sampat, Dean Stratakos, Heidi Williams, Brandon Yang, Frank Yang, and James Zou for helpful comments and conversations. We are especially grateful to Arjun Desai and Karan Goel for their work on the technical infrastructure that enabled this project. Pamela Nelson Foster provided invaluable assistance with data and computational resource access. We gratefully acknowledge financial support from the National Science Foundation through grant DGE-1656518 (Durvasula, Eyuboglu, and Ritzwoller), the Knight Hennessy Scholars Program, the Stanford Law School John M. Olin Program in Law and Economics, the National Bureau of Economic Research Innovation Information Initiative Summer Fellows Program, and the OpenAI Researcher Access Program. The data produced in this paper are available at the link: https://github.com/DavidRitzwoller/pubmed_clinical_trials.

1. INTRODUCTION

Technological progress in medicine is widely credited with enabling large improvements in human health (Easterlin, 1999; Mokyr et al., 2002; Acemoglu and Johnson, 2007). Breakthroughs at the turn of the century—the completion of the Human Genome Project, the advent of DNA sequencing, the adoption of combinatorial chemistry—promised transformation of diagnosis and treatment or, at least, an expansion of the set of available scientific opportunities.¹ Yet, an emerging literature finds that biomedical research productivity is in a state of decline; see Goldin et al. (2024) for a review. This literature is structured around an incongruity identified by Bloom et al. (2020). The quantity of medical research—proxied by estimates of the number of clinical trials and scientific publications—has increased substantially, without an accompanying improvement in mortality. Many explanations center on characteristics of the “idea production function.” Perhaps low-hanging fruit have been depleted (the *fishing out* hypothesis, described in Jones, 1995),² or else the educational costs to a young innovator of reaching the frontier of knowledge are increasingly high (the *burden of knowledge* hypothesis, due to Jones, 2009). Though the specific parameter implicated differs, the mechanisms underlying these hypotheses suggest a structural change in research productivity.

This paper tests an alternative, compositional hypothesis. We document that increases in the quantity of medical research have been driven by increases in the quantity of unproductive, low-capital research with limited capacity to effect quantitatively significant changes to mortality. As a consequence, increased quantities of medical research need not indicate that the productivity of high-capital research investment is declining.

We construct a novel census of capital-intensive, high-productivity medical investment—clinical trials. We contrast this set with its complement in clinical research, proxied by publications that cite clinical trials. This comparison is natural and reflects the standard hierarchy of medical evidence (Jones and Podolsky, 2015). Decision-making in medicine—by regulators, firms, insurers, physicians, and patients—hinges on the causal, mortality-relevant data produced by costly, time-intensive randomized experiments (Chavez-MacGregor and Giordano, 2016; Alsan et al., 2024).

¹See Scannell et al. (2012) for details on these factors, in a broader discussion of declining productivity in the pharmaceutical industry. Combinatorial chemistry, in the 1980s and 1990s, increased the number of molecules (potential drug candidates) that a chemist could synthesize “perhaps 800-fold.” Similarly, since the first genome sequence was determined in the 1970s, DNA sequencing has become “over a billion times faster.”

²See, also, Cowen (2011) and Gordon (2017).

Although research engaging with trial evidence undoubtedly produces valuable information, it does so at comparatively low cost and with lesser impact (Kones et al., 2014). Since 2010, we find that the quantity, quality, and composition of published clinical trials has been stable. By contrast, the quantity of other clinical research has doubled.

No existing dataset contains sufficient information to assess the composition of clinical research. Administrative databases—such as ClinicalTrials.gov—aim to provide an index of trials, but suffer from widespread noncompliance with reporting requirements (DeVito et al., 2020). Proprietary databases take as an input the contents of ClinicalTrials.gov and other registries and, thus, capture changes in reporting patterns over time (e.g., Cortellis, 2024). Best practices for identifying trials in publication data do so imprecisely, with both high true *and* false positive rates (Thomas et al., 2021).

We construct a new census of clinical trials. We collect all records of scientific publications indexed in the National Library of Medicine’s PubMed / MEDLINE database, from 2010 forward, that disclose the results of a clinical trial studying a medicine in human subjects. Our definition is restrictive and intended to capture only those research inputs that are most relevant for human mortality. Our goal is to identify clinical trials in our sample of interest—from the set of roughly 35 million PubMed records—with an accuracy and precision close to that of a human labeller. Standard machine learning methods perform poorly and hand-labelling at scale is infeasible. Instead, we train a bespoke large language model, optimized for our task.

We proceed in stages. We devise an interface to iteratively revise model prompts. With the highest-performing prompts, we use proprietary language models, OpenAI’s GPT-3.5 and GPT-4, to classify a large number of records (Achiam et al., 2023; Bubeck et al., 2023). These noisy labels are used as training data to fine-tune a set of open-source models. The resulting model matches the performance of the best proprietary model and produces a sample of approximately 150,000 clinical trials.

We document two sets of facts. First, there is quantitatively significant heterogeneity in both the composition of medical research and the composition of clinical trials. While clinical trial counts are essentially constant over time, the quantity of papers citing trials has increased by nearly a factor of two. This trend reflects a large increase in the quantity of research originating from China, in addition to other countries who, individually, do not account for a large share of the total quantity of medical research. Within clinical trials, proxies for quality and informativeness

suggest considerable heterogeneity. Roughly half of trials in our sample are never cited by a leading medical journal and nearly 15 percent are never cited by another scientific publication. Second, these forms of heterogeneity within trials are stable over time. We find that—across measures of quality, importance, and geography—the distribution of clinical trial heterogeneity is unchanging. Taken together, we interpret these facts as evidence that the productivity of high-capital medical investment is stable over our time period.

Our choice of setting reflects two substantive considerations. First, the determinants of technological progress in medicine are interesting in their own right.³ Second, and more centrally, nearly all micro-level evidence on rising research investment, and declining productivity, comes from studies of clinical trials and scientific publications.⁴ Measured increases in the quantity of clinical trials (Bloom et al., 2020) and average “declining disruptiveness” of scientific papers (Park et al., 2023) are at the center of active debates about whether research in science and medicine has hit a point of diminishing returns (e.g., Dieppe and Kose, 2020).⁵ Concerns about declining productivity have influenced a wave of policy interest in new mechanisms to promote research and development, as well as the revival of “industrial policy” (see e.g., Council of Economic Advisers, 2021, 2024; Lindsey and Hammond, 2020; Hausmann, 2023). A precise characterization of compositional changes in this industry, then, is a directly policy-relevant object.

We propose and validate the existence of quantitatively significant heterogeneity in medical research.⁶ This is the primary contribution of this paper. We measure research inputs with a precision that is new to this literature, enabled by the existence of novel technology for data construction at scale (Achiam et al., 2023; Bubeck et al., 2023). Our decomposition suggests that measures of research investment that neglect these differences—and, instead, infer uniform productivity or capital for scientific investments—can yield misleading conclusions about industry-wide productivity.

³Cutler and McClellan (2001) observe that technological change has accounted for the bulk of increases in health spending over time, but that health has improved alongside these increases. In a series of case studies, they argue that the benefits of technological progress exceed their costs. See, also, Almond et al. (2010), who estimate that the returns to medical spending are positive on the margin.

⁴See Goldin et al. (2024) for a summary. Notable exceptions are studies of agriculture and semiconductors.

⁵See discussions in the popular press, e.g. Collison and Nielsen (2018); Broad (2023); Thompson (2021); Piper (2023).

⁶Park et al. (2023) examine “disruptiveness” in scientific research, using data on the universe of publications across six decades. They document that the absolute counts of highly disruptive papers remain constant, consistent with our finding that high-capital, high-productivity research remains stable. Our setting allows for a clear delineation between types of research investments and, importantly, allows us to take a stand on the value of information produced by these distinct types for welfare-relevant outcomes.

This paper contributes to a growing literature on the causes of declining productivity, in general and in medicine. Additional explanations include market imperfections and failures—including barriers to innovation (see e.g., [Aghion et al., 2019](#)), inadequate public sector support for research ([Gruber and Johnson, 2019](#)), and productivity imbalances across sectors ([Acemoglu et al., 2024](#)). Our work is related to a literature—which [Syverson \(2017\)](#) summarizes as putting forth a “mismeasurement hypothesis”—that views the measured decline in aggregate productivity as illusory (see e.g., [Mokyr, 2014](#)). In contrast, our hypothesis is not that productivity is mismeasured, per se, but that failure to account for heterogeneity can generate a misattribution of trends. More directly, [Myers and Pauly \(2019\)](#), [Scannell et al. \(2012\)](#), and [Cockburn \(2006\)](#) investigate the determinants of declining productivity in the pharmaceutical sector, with a focus on large increases in factor prices for clinical trials rather than quantities.

Finally, we draw on, and contribute to, a long literature on the measurement of innovation. Since [Pakes \(1986\)](#), economists have recognized that patents—the most commonly used measure of research *output*—differ greatly in economic and technological significance. [Griliches \(1990\)](#) provides a classic discussion of the difficulties associated with using patents as economic indicators, given this heterogeneity.⁷ Our work builds on these insights to investigate the implications of heterogeneity in *inputs* for the use of medical research as a productivity indicator.

2. PRODUCTIVITY, CAPITAL, AND MEDICAL RESEARCH

Studies of research productivity have, at their core, some notion of an “idea production function”—a mapping from research inputs (*effort*) to outputs (*ideas*). The canonical [Romer \(1990\)](#) specification is given by

$$\frac{dA_t}{dt} \frac{1}{A_t} = \alpha_t S_t . \quad (2.1)$$

Here, the variable A_t denotes total factor productivity. The left-hand side of (2.1) captures the rate at which “ideas” are produced. On the right-hand side, the factor S_t measures capital devoted to research. The structural parameter α_t is interpreted as research productivity—the efficiency with which inputs are transformed into outputs.⁸

⁷[Bryan and Williams \(2021\)](#) provide an updated survey and discuss the challenges of devising credible measures of innovative effort.

⁸See e.g., [Romer \(1990\)](#), [Jones \(1995\)](#), and [Aghion and Howitt \(1992\)](#), among many others.

The revival of interest in tools, such as industrial policy, to better allocate public resources to research and development is a reaction, in part, to two stylized facts (Council of Economic Advisers, 2021, 2024; Lindsey and Hammond, 2020; Hausmann, 2023). First, the growth rate of research output has been constant (*i.e.*, the left-hand side of (2.1) is constant). Second, the quantity of resources devoted to research (*i.e.*, the factor S_t) has increased substantially. The linear specification (2.1) then implies that research productivity α_t has declined. That is, ideas are getting harder to find. For a review of the relevant literature, see Goldin et al. (2024). Detailed characterization of these facts is due to Bloom et al. (2020).

2.1 Heterogeneity. We propose an alternative reconciliation of these facts. Our hypothesis hinges on heterogeneity. If measured increases in the quantity of research are accompanied by a compositional shift toward unproductive or low-capital research, then constant growth does not necessarily imply declining productivity for capital-intensive research. That is, large increases in the quantity of low-capital-intensity or low-productivity research need not have quantitatively important impacts on growth.

To fix ideas, let the quantity \mathcal{I}_t index the set of “research units” available at time t . There are two types of research units, high, H, and low, L. Each research unit of type H has a stock of research capital S_t^H and a productivity coefficient α_t^H . Similarly, research units of type L have research capital and productivity S_t^L and α_t^L , respectively. Consider the generalized production function

$$\frac{dA_t}{dt} \frac{1}{A_t} = |\mathcal{I}_t^H| \alpha_t^H S_t^H + |\mathcal{I}_t^L| \alpha_t^L S_t^L, \quad (2.2)$$

where \mathcal{I}_t^H and \mathcal{I}_t^L index the set of H-type and L-type research units and $|\cdot|$ denotes the cardinality operator. That is, relative to the canonical specification (2.1), we entertain heterogeneity across research units in both productivity and capital-intensity.

Suppose that the quantity, productivity, and capital of H-type units is constant over time. Thus, in each period t , the first term in (2.2) is given by $|\mathcal{I}^H| \alpha^H S^H$, where we have removed the time subscript. Suppose that at time $t = 0$, there are no L-type units, and that at time $t = 1$, the quantity of L-type units \mathcal{I}_1^L becomes non-zero. The impact of this change only affects growth through the factors S_1^L and α_1^L . If either quantity is negligible, relative to the contribution of the H-type units, then there is a negligible impact on growth.

This story is closely related to a strand of the economic growth literature that argues that measured increases in the quantity of research reflect increases in the number of product varieties (Dinopoulos and Thompson, 1998; Peretto, 1998; Young, 1998; Howitt, 1999). By contrast, we examine the idea that shifts in the relative shares of established varieties can rationalize the same trends.

2.2 Medical Research. The objective of this paper is to test the compositional hypothesis described in Section 2.1. We do so in the context of medical research. We argue that in this setting we are able to decompose research units into coarse types. Moreover, for each type, institutional context allows us to infer absolute and relative productivity and capital-intensity. Measures of quantity can be recovered from existing data. An additional advantage of this setting is that medical research has an intuitive notion of output: improvement in mortality. To keep our focus on the composition of inputs—the right-hand side of (2.2)—we note that it is a well-established empirical regularity that life expectancy has risen at a steady rate since 1840 (see e.g., Oeppen and Vaupel, 2002; Bloom et al., 2020). We proceed assuming that the left-hand side of (2.2) is constant, as is standard in this literature.⁹

Medical research takes many forms. We focus on two types: clinical trials, and other forms of clinical research that aim to inform the provision of medical care. Clinical trials are expensive, time-consuming, highly-regulated experiments, which have as their goal the production of causal and mortality-relevant evidence (Alsan et al., 2024; Wouters et al., 2020). The contours of the second set, other forms of clinical research, are less precise. This set includes observational studies, meta-analyses and systematic reviews, cost-effectiveness studies, etc. (Food and Drug Administration, 2024). In contrast to clinical trials, these other study types are less capital-intensive and less productive, where productivity is defined as their efficiency in turning research inputs into changes in mortality (Jones and Podolsky, 2015; Chavez-MacGregor and Giordano, 2016; Kones et al., 2014).¹⁰

⁹We make the assumption that life expectancy is improving at a steady rate, given the stability of these trends for more than 150 years. It is of course unlikely that medical innovations developed in the early 2010s are clearly discernible in recent mortality data, especially given the impact of the COVID-19 pandemic.

¹⁰Of course, one can identify specific studies that deviate from our categorization. See Kones et al. (2014) for discussion of a setting in which, the authors argue, observational data are as informative as randomized controlled trial evidence. However, clinical trials remain the “gold standard” for medical research (Jones and Podolsky, 2015). Non-randomized evidence is consistently excluded from prescribing guidelines, on the grounds that it provides evidence of lesser quality. Chavez-MacGregor and Giordano (2016) and Kones et al. (2014) discuss hierarchies of medical evidence in more detail.

3. IDENTIFYING CLINICAL TRIALS IN PUBLICATION DATA

We construct a census of clinical trials. No existing dataset, used off-the-shelf, is sufficient. Although dollar-denominated expenditures are an obvious metric, available records of investment in clinical trials are self-reported, in some form, by pharmaceutical firms and are essentially impossible to validate independently.¹¹ Administrative databases—such as ClinicalTrials.gov—aim to provide an index of trials, but suffer from widespread non-compliance with reporting requirements, with substantial variation in enforcement over time (DeVito et al., 2020). Proprietary databases pose two challenges: first, it is, again, difficult to validate their contents and, second, leading databases take as inputs data in public registries, thus adopting the same shifts in reporting over time (see e.g., the list of registry inputs for one frequently-used database, Cortellis, 2024). For each of these options, there is no symmetric way to identify the quantity of other clinical research.

3.1 Publication Data. Publication data provide a natural alternative. The database PubMed / MEDLINE indexes roughly 34 million publications, a near universe of published biomedical research.¹² Scientific publications provide an intuitive proxy for the quantity of research inputs.¹³ Of these 34 million records, which are clinical trials?

The existing literature takes several approaches. In health sciences—where data on the universe of clinical trials are relevant for the construction of meta-analyses and the development of prescribing guidelines—researchers typically employ over-inclusive methods with high true positive rates. To use one prominent example, Cochrane, a British organization that synthesizes the findings of medical research, developed the “Cochrane RCT Classifier,” a machine learning-based classifier for retrieving randomized controlled trials. The Cochrane approach successfully identifies trials with a 99 percent true positive rate. For every eight true positives, however, it returns 92 false positives. See Thomas et al. (2021) for details.

¹¹There are two frequently used sources of data on dollar-denominated expenditures. DiMasi et al. (2016) provide one set of estimates of drug research and development costs, based on confidential data obtained from 10 firms. Myers and Pauly (2019) and Scannell et al. (2012) use research and development expenditure figures published by the Pharmaceutical Research and Manufacturers of America, a trade organization representing the industry.

¹²See Appendix A.1 for further details on our treatment of the PubMed data.

¹³The existence of publication bias means that scientific publications do not perfectly capture research inputs. See Andrews and Kasy (2019) for a discussion. In a setting such as ours—where objects of interest are trends in research over time—we need only assume that changes in publication bias year-on-year do not perfectly offset changes in research investment.

Approaches used by researchers in social science, instead, largely draw on metadata contained in PubMed.¹⁴ Lichtenberg (2018) and Bloom et al. (2020) collect all records assigned a “publication type”—in the PubMed indexing process—of “clinical trial.” Feldman et al. (2019) collect all records with a publication type similar to clinical trial (e.g., also including “randomized controlled trial”).

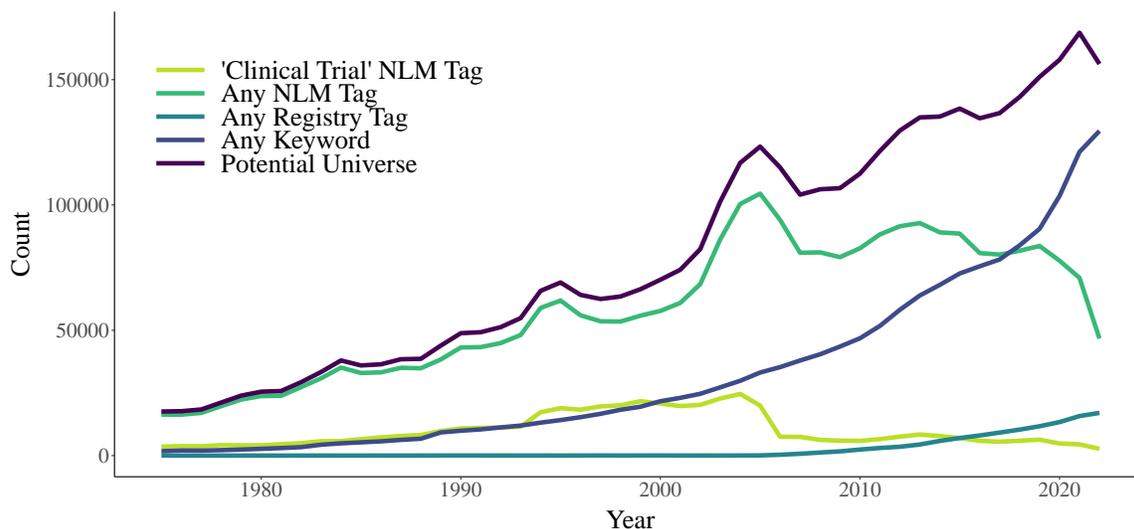
On inspection, these tags are imprecise. Both approaches retrieve large numbers of records that are not, in fact, clinical trials, while excluding potentially relevant records. Figure 1 contrasts counts of published trials constructed in various ways. The light green line (‘Clinical Trial’ NLM Tag) plots the number of publications, in each calendar year, indexed with the publication type “clinical trial.” The teal line (Any NLM Tag) plots counts of records indexed with any of 18 types that are likely to include clinical trials or related medical research.¹⁵ Alternative approaches to identifying trials are similarly imprecise. Medical journals increasingly require trials to be registered in open registries, e.g., ClinicalTrials.gov, as a condition of publication. The blue line (Any Registry Tag) counts records that report, in their abstract, an identifier associated with one of the four largest international trial registries. Language in a record’s abstract can indicate that it might be a clinical trial if, for example, it references a “treatment group” and “control group.” The dark blue line (Any Keyword) plots the count of records with such keywords over time. We define the “potential universe” of clinical trials as those publications with any of the following: a publication type variable similar to clinical trial, a registry identifier reported in its abstract, or a keyword reported in its abstract. The purple line (Potential Universe) plots this trend over time. Figure 1 highlights that alternative approaches to data construction, in this setting, yield meaningfully different conclusions about levels, trends, and composition. See Appendix A.2 for further details on the construction of these series.

3.2 Objective. What exactly are we trying to recover when we search for clinical trials? We argue that the object of interest is a census of clinical trials that study the effects of a medicine in human subjects.

¹⁴Ostrom (2024) and Kao et al. (2023) are two exceptions. Ostrom (2024) uses data from a meta-analysis of published clinical trials. Kao et al. (2023) collect trial records from the proceedings of scientific conferences.

¹⁵On inspection, it appears as though the NLM shifted from using the tag “Clinical Trial” to “Randomized Controlled Trial”, as well as more specific tags (e.g., “Clinical Trial, Phase 1”) .

FIGURE 1. Universe of Potential Clinical Trials



Notes: Figure 1 displays counts of the number of clinical trials indexed in PubMed / MEDLINE over time, constructed using alternative search strategies. The National Library of Medicine (NLM) categorizes each publication into a “pubtype.” The light green line (‘Clinical Trial’ NLM Tag) displays the number of publications in the “Clinical Trial” pubtype. The teal line (Any NLM Tag) displays the number of publication whose pubtype is an element of a set of 18 categories likely to include clinical trials. The blue line (Any Registry Tag) gives the number of publications that report, in their abstract, an identifier associated with one of the four largest international trial registries. The dark blue line (Any Keyword) indicates the number of publications whose abstract contains a keyword indicative of a clinical trial. The purple line (Potential Universe) displays the number of clinical trials in the union of the sets of publications identified with the other lines. See Appendix A.2 for further details.

Definition 3.1. The sample of interest is composed of all publications that report the results of a prospective, interventional clinical trial that evaluates the effects of investigational or approved drugs in a setting with exclusively human subjects.

Definition 3.1 embeds several restrictions, intended to yield a measure of new investment in information relevant to human-facing medical care. We view this set of studies as the set most likely to generate information relevant to human mortality. Studies involving animals, literature summaries, and re-analyses of existing data, for example, are excluded. We take the “potential universe” plotted in Figure 1 as our baseline sample of potential clinical trials. We further restrict attention to the 1,821,429 records published in or after 2010.¹⁶

¹⁶There is a trade-off between data coverage and quality. Around 2010, a set of policy changes—including trial registration requirements and harmonization in publication formatting—affected trial disclosure in scientific publications. See Laine et al. (2007) and Schulz et al. (2010) for examples of two such changes.

To quantify our ability to identify the clinical trials in this set of publications, we hand-label approximately 3,000 randomly selected publications based on the content of their abstracts.¹⁷ Of the hand-labelled publications, 11.2% meet the criteria for inclusion in our sample. The hand-labelled data are split into three subsets—validation, training, and testing—based on their eventual use.

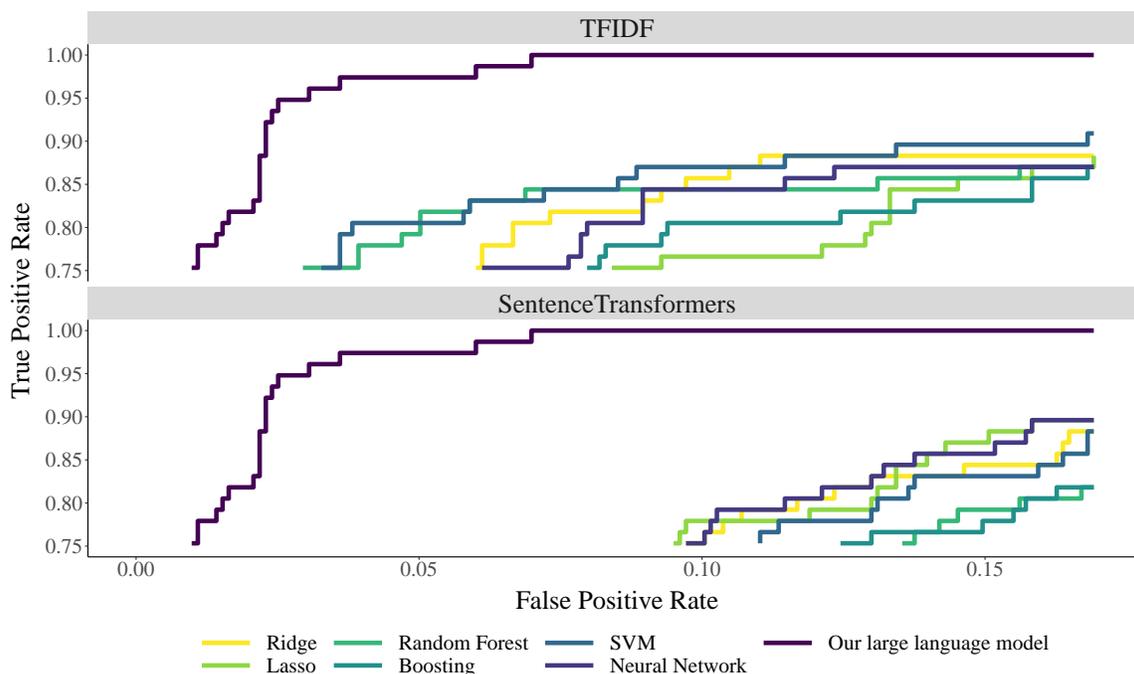
As a baseline, we assess the performance of several standard machine learning algorithms. The results are displayed in [Figure 2](#). Further details on the construction of this figure are given in [Appendix B.2](#). Each algorithm is estimated on the hand-labels assigned to the publications in the training and validation sample splits. For feature vectors, we use either TF-IDF embeddings computed in the corpus of abstracts in the hand-labeled sample or the embeddings of each abstract obtained from the SENTENCETRANSFORMER language model (Reimers and Gurevych, 2019). We find that the performance of the standard machine learning algorithms is again unsuitable for our application. At a 90% true positive rate, the best performing model identifies 50 true positives for every 50 false positives.

3.3 Prompt Design and Fine-Tuning. We construct a large language model optimized for our task. This is accomplished in three steps. First, we iteratively construct a set of prompts that exhibit good performance when posed to proprietary models—OpenAI’s GPT-3.5 and GPT-4. Second, we extract noisy labels for a large number of publications in our sample by querying the proprietary models. Third, the noisy labels are used to train an off-the-shelf large language model. The resulting model is then used to identify clinical trials in our baseline sample. This process of constructing specialized large language models is referred to as “fine-tuning” (see e.g., Taori et al., 2023, for a prominent application of this idea).

3.3.1 Prompt Design. We devise three types of prompts and measure their performance on samples of abstracts taken from our validation dataset. We measure model performance with the hand-labeled data, conduct error analyses, and iteratively revise prompts. When modifications to the prompt cease to yield improvements in performance, we finalize prompts best suited for GPT-3.5 and GPT-4, respectively. [Appendix B.3](#) provides more detail on this process. [Figure 3](#) displays estimates of the true and false positive rates computed in the test data, for both proprietary models, with black dots.

¹⁷We develop a custom labelling interface that reduces the time required to label 100 records, for the authors, by a factor of seven. [Appendix B.1](#) describes this tool and its application to this context in detail.

FIGURE 2. Performance of Standard Machine Learning Methods



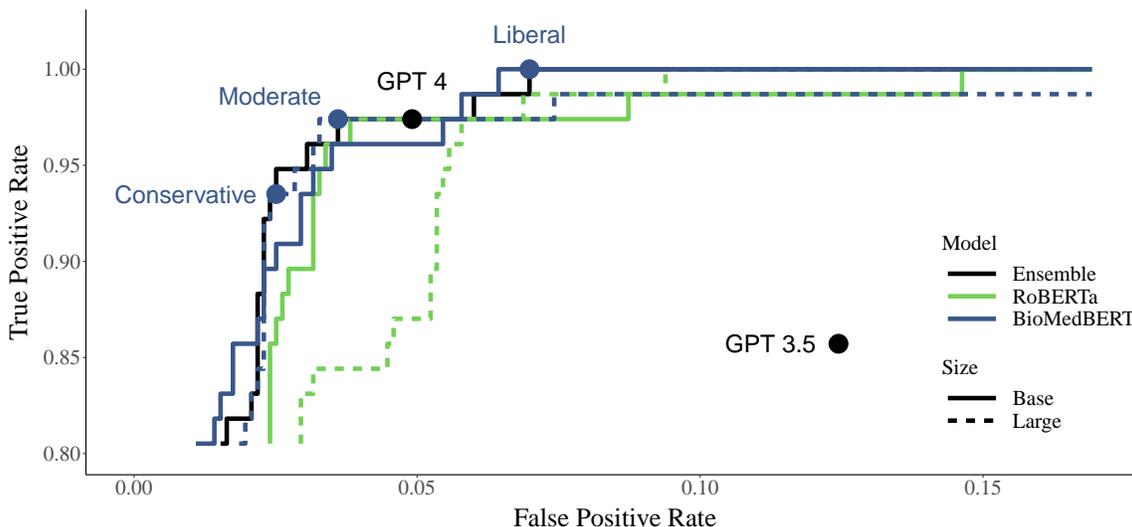
Notes: Figure 2 displays the receiver operating characteristic of several machine learning models, trained on two types of embeddings, for classifying whether a publication satisfies the restrictions enumerated in Definition 3.1. The hyper-parameters of each model are determined with 10-fold cross-validation. The models are trained in the training and validation data set. Error rates are estimated in the testing data set. In addition, we display the performance of our, ensemble, fine-tuned large language model in purple. This is the same curve labelled “Ensemble” in Figure 3.

GPT-4 performs remarkably well. For every 71 true positives, GPT-4 identifies 29 false positives, while achieving a true positive rate of approximately 0.974. This represents a substantial improvement over existing methods. However, practical and substantive considerations mitigate the applicability of proprietary models for classification of the complete baseline sample. Practically, it is—and likely will remain—prohibitively expensive to deploy GPT-4 at this scale.¹⁸ Substantively, proprietary models are black boxes. Their details and substance are not public and are known to change at regular intervals.

3.3.2 Fine-Tuning. We compute noisy labels for 64,000 randomly selected publications in our baseline sample using the best performing prompts for both GPT-3 and GPT-4. These noisy labels are used as data to train off-the-shelf language models from two architecture classes: (1) ROBERTA

¹⁸We incur a cost of roughly \$4,500 to extract noisy labels for 64,000 abstracts. Extrapolating this figure to the full sample of 1.8 million abstracts gives a price of approximately \$130,000.

FIGURE 3. Receiver Operating Characteristic



Notes: Figure 3 displays estimates of the receiver operating characteristic of several models used to classify the publications indexed by PubMed / MEDLINE according to whether they satisfy the restrictions enumerated in Definition 3.1. All metrics are computed in the testing split of the hand-labelled publications. Estimates of the true positive rate and false positive rate of OpenAI’s proprietary models GPT-3 and GPT-4 are indicated with black dots. The curves give the performance of four fine-tuned, open-source large language models, in addition to an ensemble model. The blue dots indicate the true positive rate and false positive rate of the models used to construct the “Conservative,” “Moderate,” and “Liberal” samples of clinical trials.

(Liu et al., 2019) and (2) BIOMEDBERT (Gu et al., 2021). Both models are available in either “baseline” or “large” varieties.¹⁹ The BIOMEDBERT model is an instance of a BERT model (Devlin et al., 2019), that has itself been fine-tuned on the abstracts of the publications indexed in PubMed. We systematically vary model architecture, size, and pre-training regimen. We observe no large differences in performance stemming from these hyper-parameters. Appendices B.4 and B.5 give further details the fine-tuning process and the results of these experiments, respectively.

3.4 Performance. Given the text of an abstract, our fine-tuned language models output a probability that the publication satisfies the restrictions enumerated in Definition 3.1. Publications whose probabilities fall above a chosen threshold are classified as belonging in our sample. Figure 3 displays estimates of the true and false positive rates, computed with the test data, as we vary this

¹⁹There are many open-source large language models (e.g., LLaMa, Mistral, Pythia). Many are several orders of magnitude (e.g., 7-70 billion parameters) larger than ROBERTA (~ 350 million parameters). We selected ROBERTA and BIOMEDBERT via trial-and-error that included testing the performance of these larger models. In our setting, ROBERTA outperforms LLaMA (both 7 and 70 billion parameter versions, trained with QLoRA). Both models exhibited comparable performance, but ROBERTA was substantially faster and simpler to train.

threshold for both sizes of the ROBERTA and BIOMEDBERT models.²⁰ The figure additionally displays the performance of an ensemble model estimated in the training data. This model is obtained by estimating a logistic regression of the hand-labels on the probabilities output by all four models displayed [Figure 3](#). The ensemble model is used to produce our final sample.

The fine-tuned models are able to match the performance of GPT-4 in the test data. We choose three thresholds according to the stringency with which they enforce the sample restrictions. The estimated true and false positive rates associated with these points are displayed in [Figure 3](#) and labeled as “Conservative,” “Moderate,” and “Liberal.” Our preferred sample is associated with the conservative threshold.²¹ For every 82 true positives, the conservative model identifies 18 false positives.²² We report results associated with the moderate and liberal thresholds in the appendix as tests of robustness. Our final, conservative, sample consists of 152,027 publications classified as satisfying the restrictions enumerated in [Definition 3.1](#).

[Appendix A](#) compares the contents of our sample to the counts of potential clinical trials plotted in [Figure 1](#). Although certain elements of PubMed metadata—including the union of all NLM tags—capture many of the records in our final sample, we confirm that they miss many records that we flag as trials and include many records that do not satisfy [Definition 3.1](#). No combination of existing search strategies, then, classifies records with the same accuracy or precision as our final sample.

4. MEDICAL RESEARCH AS A PRODUCTIVITY INDICATOR

[Figure 1](#) suggests that, across metrics, there has been a substantial increase in the quantity of medical research over time, including since 2010. This is consistent with documented increases

²⁰Each of the models whose results are displayed in [Figure 3](#) are trained with noisy labels obtained with GPT-4. We document in [Appendix B.5](#) that models trained with noisy labels obtained from GPT-3.5 or with the 1000 hand-labeled records in the training data perform significantly worse.

²¹The liberal model may be particularly useful for conducting literature reviews, where a near-perfect true positive rate is needed.

²²In the test data, the conservative model assigns incorrect labels to 27 of 993 papers. We conduct an error analysis. See [Appendix B.6](#) for further details. In 13 cases, there is a clear error. In 14 cases, however, errors are associated with records that were difficult to categorize for a human labeller. For example, PubMed record 32737793 is flagged as satisfying [Definition 3.1](#) with all three thresholds, but is a literature review. By contrast, PubMed record 27880726 was categorized as satisfying [Definition 3.1](#) twice by a human labeller. It is excluded from both the conservative and moderate model-generated samples. On inspection, the abstract does not explicitly state that the study enrolled human subjects, but hints that it may have been conducted in an animal model. Review of the associated full text confirms that this study did, in fact, enroll only rats.

in the quantity of scientific papers in Bloom et al. (2020) and Park et al. (2023), as well as economy-wide increases in research effort reported in Goldin et al. (2024).

4.1 Decomposing Medical Research. We decompose this pattern, into changes attributable to high-productivity, high-capital investments—clinical trials—and changes attributable to other forms of clinical research. Our census of clinical trials provides measures for the first category. Although there are many reasonable ways to construct samples of other clinical research, our preferred approach defines this category as those scientific publications that *cite* clinical trials.²³ These are records of research investment that necessarily engage with findings from trials (such as meta-analyses and re-analyses of existing data), but, as discussed in Section 2, are essentially always lower productivity and less capital-intensive.

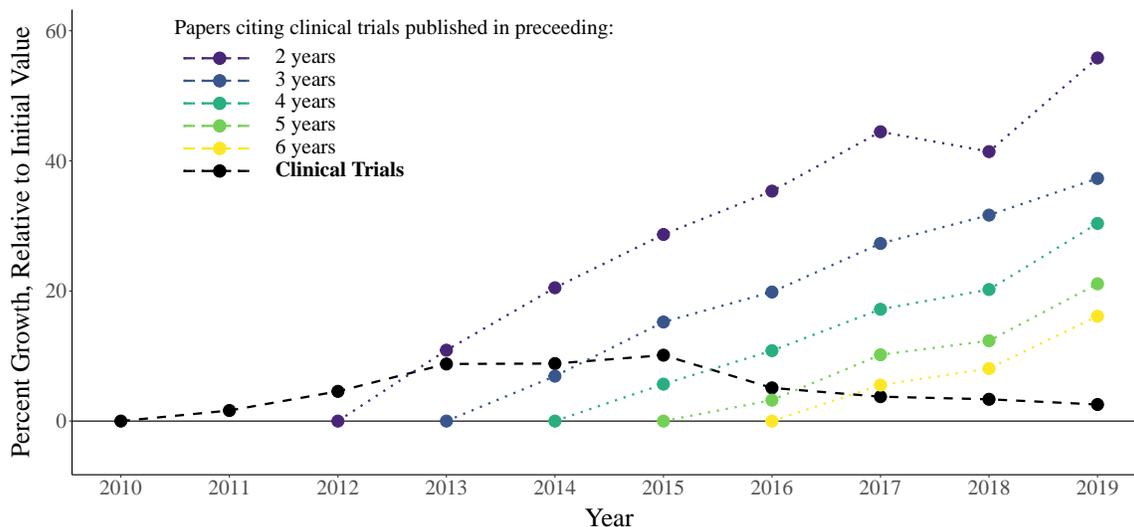
Figure 4 documents that the quantity of clinical trials—the black line—has remained constant over the past ten years. In contrast, other clinical research has increased substantially. Observe that we display the percent *growth* in the quantity of clinical trials in our sample published in each calendar year, alongside the *growth* in the set of papers that cite clinical trials. Each series is normalized to begin at zero. To account for concerns about citation delays and truncation, we construct various measures of papers that cite clinical trials, using 2–6 year citation counts. Two-year citation counts, thus, are available starting in the year 2012 (two years into our sample) and count publications that cite clinical trials published in the preceding two years.²⁴ This presentation masks substantial differences in the levels of these series. There are 10,903 clinical trials in our census in 2010, and 73,253 publications that cite clinical trials (using two-year citations) in 2012.

4.2 Stable Composition of Clinical Trials. To contextualize the trends displayed in Figure 4, it is helpful to return to the simple model from Section 2. We document that the quantity of capital-intensive, high-productivity research inputs remains unchanged over our period of interest, while low-capital, low-productivity inputs increase. As changes in outputs in this sector (mortality) are assumed to remain constant (see Section 2), there are thus three ways of reconciling these facts

²³Alternative approaches to construction of this sample, as in Figure 1, also suggest an increase in the quantity of non-clinical trial research. We prefer to use this citation-based definition because it allows for consistency in sample definition. Other approaches that rely on keyword searches, NLM tags, etc., as Section 3 and its appendices make clear, capture different sets of records over time.

²⁴Figure C.3 and Figure C.4, displayed in the supplemental appendix, document that this trend is robust to alternative sample definitions, using either our liberal or moderate samples, and to citation-weighting count data.

FIGURE 4. Growth in Clinical Research, Stability of Clinical Trials



Notes: Figure 4 displays measurements of the number of clinical trials, and papers that cite clinical trials, published in each calendar year. Each series is reported in terms of the percent change relative to its initial value. The sample of published clinical trials is constructed with the conservative model. To address truncation, we report the number of publications that cite clinical trials published in the preceding t years for each t between 2 and 6.

in our stylized framework: (i) other clinical research is sufficiently low capital or low productivity, across the board, as to have essentially no impact on growth, (ii) increases in other clinical research have, heterogeneously, low capital or low productivity, or (iii) clinical trial research productivity is declining over this time period, fully offsetting the growth in other types of research. If either Case (i) or Case (ii) were an accurate characterization, rising observed research effort—increases in the number of scientific publications or the number of records in PubMed that the previous methods have identified as clinical trials—need not imply diminishing returns to high-capital investment.

Our data allow us to reject the hypothesis that we are in Case (iii). Figure 5 documents that—on four dimensions—measures of trial productivity, across clinical trials, have remained stable over the past decade. Panel A of Figure 5 considers the share of publications in our census with the following characteristics: any funding from a government agency (purple), any (three-year) citations from a leading journal in medicine (teal), or any citations from a leading journal *conditional* on

having public funding (green).²⁵ The trends in Panel A communicate two findings. First, all three proxies for research productivity are stable, essentially unchanging, over our time period. Second, nearly 55 percent of clinical trials are *never* cited by a leading journal, suggesting substantial, but time-invariant, heterogeneity across trials.

These proxy variables capture slightly different conceptual objects. Interpreting each, separately, provides suggestive evidence on different aspects of the research production function. Stability in public funding suggests minimal changes, at least in the cross section, in allocation of resources to clinical trials over time.²⁶ Stability in citation patterns, in contrast, suggests that the “usefulness” of information produced by trials also remains stable.

Panel B of Figure 5 provides a more granular view of this heterogeneity. For each publication in our census, we collect citations from publications in leading journals and from all of PubMed. We plot the total (three-year) citation counts to publications at each citation quantile over time.²⁷ Grey areas at the bottom of each plot represent publications that receive zero citations. This distribution is stable over our time period.²⁸ Specifically, we document that the tails of the citation distribution are unchanging. Observe that we elongate the y-axis, to highlight the right-tail. Citations received at the top of this distribution are constant across periods. The left-tail is similarly stable: the share of records receiving zero citations is essentially constant, year-on-year.

This result is consistent with the idea that, for clinical trials, the right-hand side of the production function, (2.1), in our simple framework—i.e., $\alpha_t S_t$ —is both heterogeneous and stable over time. For this pattern to be consistent with a systematic decrease in the research productivity of clinical trials, it must be the case that there is an exactly offsetting increase in the distribution of research

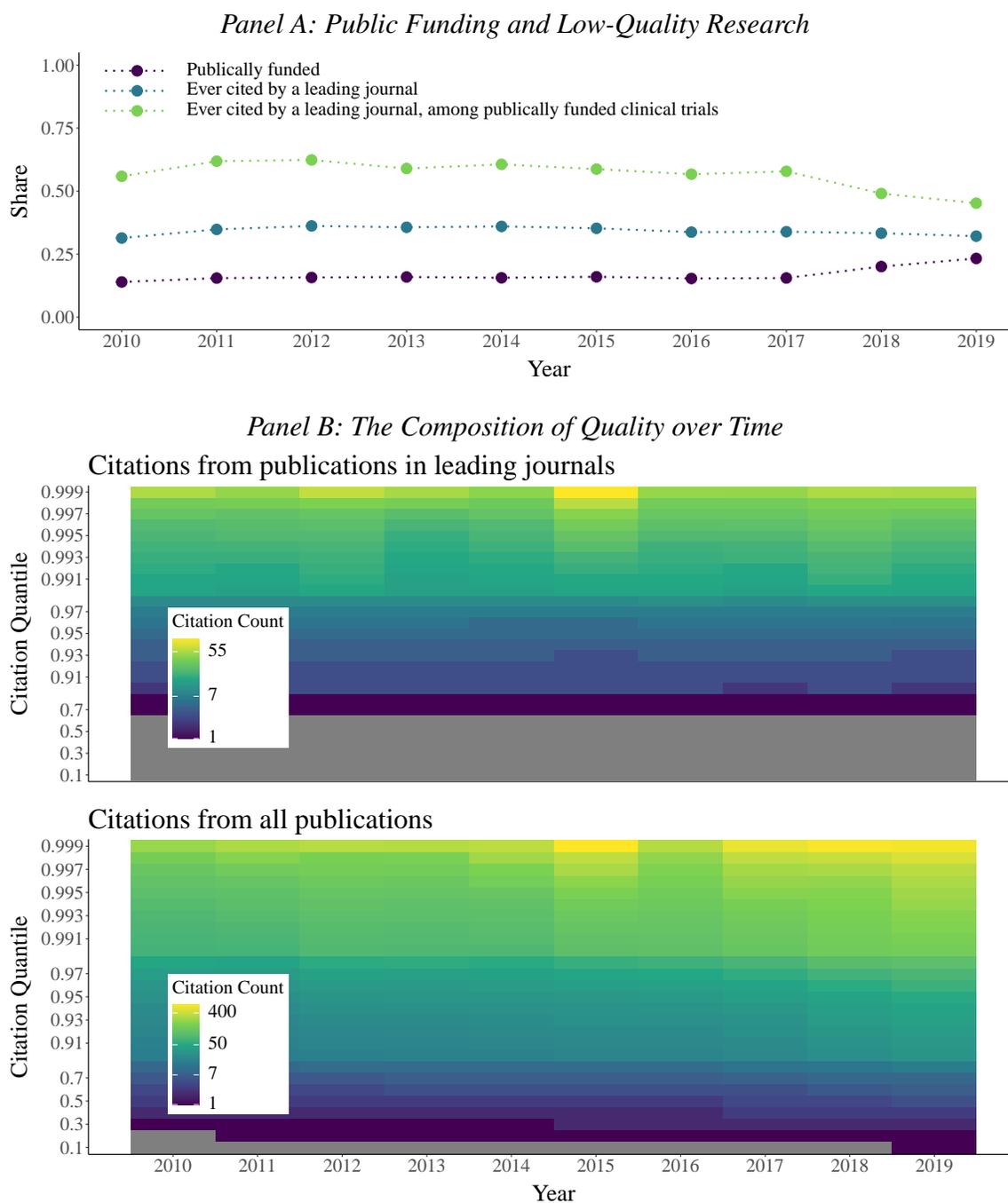
²⁵We designate a publication as having public funding if it appears in the National Institutes of Health RePORTER data, as linked to a funded research grant, or if PubMed reports a source of research funding, including funding from non-U.S. agencies. To designate a set of the “top” journals in medicine, we follow Angrist et al. (2020) We collect all citations originating from the “trunk journals” in medicine, the Journal of the American Medical Association and the New England Journal of Medicine, to records in PubMed between 2010 and 2022. We designate a journal as being “leading” if it received at least 100 citations from a trunk journal over this time period. This yields a list of 84 journals. See Appendix A.3 for details. Appendix C reproduces this figure using five-year citations. The trend is unchanged.

²⁶NIH RePORTER data allow us to tab total expenditures associated with grants acknowledged in these publications. These grant-based expenditure measures, too, are stable. Note, however, that linking research grant dollars to specific scientific papers is an imprecise exercise. See Li (2017) for an extended discussion of this disconnect.

²⁷Appendix C presents these results using alternative cuts of our data and five-year citations. The patterns remain unchanged.

²⁸Reassuringly, Panel B indicates that this exercise does capture large shifts in the research environment: citation counts appear especially high in the later years of data, when three-year citation counts (beginning for year 2017) capture research written during the COVID-19 pandemic.

FIGURE 5. Stability of Heterogeneity



Notes: Figure 5 illustrates the stability of the heterogeneity in the composition and quality of published clinical trials between 2010 and 2019. Panel A displays three time series: the proportion of clinical trials that are publicly funded, the proportion of clinical trials that are ever cited by a leading journal, and the proportion of publicly funded clinical trials that are ever cited by a leading journal. Panel B displays two heat maps measuring the distribution of citations received by clinical trials in each calendar year from leading journals and from all publications, respectively. The y -axis has been stretched to elongate the right-tail of the citation distribution. Colors are displayed in a log scale. For both panels, the sample of published clinical trials is constructed with the conservative model.

capital. Although we cannot reject that there are shifts in upstream, basic inputs into clinical trials or large changes in factor prices during this period, there is little credible evidence consistent with either during this time period. A small number of papers examine frictions in markets for basic research. See, for example, [Hill and Stein \(2021\)](#) and [Myers \(2020\)](#). To our knowledge, this literature has not provided evidence of substantial shifts in these frictions over time, which in turn would increase the cost or decrease the productivity of basic science.²⁹ On factor prices, as discussed in Section 3, nearly all estimates of research and development costs are self-reported by the pharmaceutical industry. Any measured increases in these expenditures that cannot be independently verified should be viewed with caution, especially in light of our findings of stability.

More generally, any type of cancellation—an argument against the conclusion that research productivity across clinical trials is stable—would appear to be a knife-edge case. In our view, it is far more likely that the distributions of both research productivity and research capital *within* clinical trials have been constant over this time period.

5. DISCUSSION

All research investments are not alike. Certain investments are more likely to yield innovations that have welfare-relevant impacts. Others, while potentially valuable contributions to a research ecosystem, may have insignificant impacts on growth. We argue that this type of heterogeneity is quantitatively significant in medical research. We use an open-source, fine-tuned large language model to construct a novel census of clinical trials from 2010 to the present. We contrast stable trends in clinical trials—in quantity, quality, and composition—with substantial growth in the quantity of other forms of medical research. Our findings rationalize, in this field, two stylized facts about research investments and innovative output: the quantity of research inputs has grown substantially, while improvements in mortality continue at the same stable rate.

A second dimension of heterogeneity is helpful in contextualizing our findings: geography. [Figure 6](#) displays trends in the annual quantity of published clinical trials, and papers citing clinical

²⁹Anecdotal evidence suggests that there may have been substantial decreases in the cost of basic science during this time period. See [Scannell et al. \(2012\)](#).

trials, disaggregated by the location of the first-listed author.³⁰ In our sample, the top four producers of published clinical trials, and papers citing clinical trials, are, in order, the United States, China, Germany, and Japan. Clinical trial investment—across geographies—is stable.³¹ Roughly 30 percent of published clinical trials originate in the United States. This share is unchanging across periods. By contrast, Panel B indicates a large shift in the medical research ecosystem. Since 2013, there has been a small increase in the quantity of other medical research published by authors in the United States. This small change is dwarfed by increases—on the order of 30-50 percent—in the quantity of other medical research from China and from countries outside of the top four producers of medical research.

These compositional facts are, in many ways, unsurprising. In recent years, editorials in prominent journals in medicine have observed that changes in the incentive structures for physicians, scientists, and even medical students, across countries, have led to “an extreme of quantity at the expense of quality” in medical publishing (Siegel et al., 2018).³² Recent work documents a substantial increase in the quantity of meta-analyses and systematic reviews—summaries of existing clinical trials—driven by scientists from outside of the United States (Ioannidis, 2016a). In parallel, there has been a rapid increase in the number, and scale, of “mega-journals”—scientific journals published frequently, with a large and growing number of pages (Ioannidis et al., 2023).

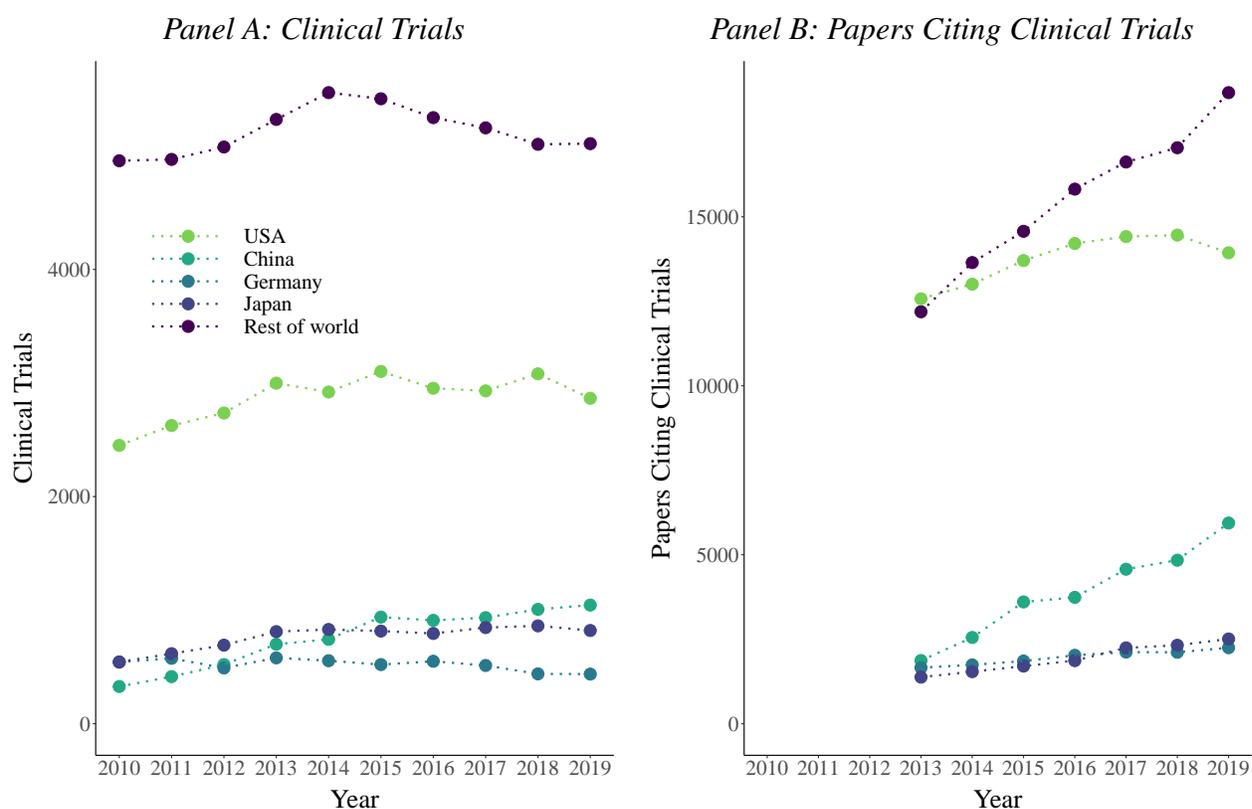
Nonetheless, substantial constraints—financial, ethical, and practical—keep firms, physicians, and scientists from increasing the frequency of clinical trials. Clinical trials are infrastructure-intensive. Trials are run at specific sites, which must identify patients who satisfy all relevant eligibility criteria and monitor them for periods ranging from six months to twenty years (Piantadosi, 2024). Best available estimates suggest that trials can cost on the order of 500 million to one

³⁰In most domains of scientific and medical publication, last-listed authors are senior investigators. First-listed authors are typically junior investigators. Agha and Molitor (2018) observe that, for large-scale clinical trials, first-authors are more likely to be the principal investigator. In our context, the first- and last-author have the same listed country for 85 percent of clinical trial records and 91 percent of non-clinical trial records. We collect details on author location from Clarivate Analytics’ Web of Science. See Appendix A.4 for further details.

³¹Note that Panel A of Figure 6 displays measurements of the location of clinical trial publication authors, not the location of clinical trial sites. We cannot rule out that, for example, trials authored by researchers with United States mailing addresses were conducted elsewhere. Thus, we do not interpret these as facts on clinical trial “offshoring.” See Petryna (2007) and Durvasula (2023) for longer discussions of the geography of clinical trial sites.

³²See, as one example, program certification requirements, which describe publication expectations for medical students and clinical faculty, from the Accreditation Council for Graduate Medical Education in Accreditation Council for Graduate Medical Education (2017).

FIGURE 6. Composition of Medical Research Across Countries



Notes: Figure 6 displays measurements of the number of clinical trials, and papers that cite clinical trials, published in each calendar year, by country. The sample of published clinical trials is constructed with the conservative model. To address truncation, we report the number of publications that cite clinical trials published in the preceding three years.

billion dollars to complete (Qiao et al., 2019). Perhaps most importantly, sites must be capable of producing data sufficient to persuade regulators. Certain regulators require that sponsors of new drug applications submit trial evidence collected from a domestic population.³³ In recent years, however, the U.S. Food and Drug Administration (FDA)—responding to concerns from patient groups—has indicated a limited willingness to approve drugs on the basis of evidence collected from exclusively foreign sites (for a longer discussion, see Alsan et al., 2024). In one especially high-profile example, a cancer drug was rejected by the FDA after being tested exclusively in China (Kolata, 2022). As long as the U.S. market continues to have outsized value for pharmaceutical firms—a consequence

³³Until December 2023, Japan required domestic phase I clinical trials of drugs developed overseas before Japanese individuals could participate in international phase III trials for pharmaceutical regulatory approval. This policy reflected, in part, differences in disease burden in Japan, relative to countries where drugs are routinely tested. This may explain why Japan is a top producer of clinical trials over our sample period. See Namba et al. (2024) for one discussion.

of the especially high prices paid by American consumers—it may be unsurprising to find that evidence disproportionately originates in the United States. Observational studies, meta-analyses, and case reports are, by contrast, cheap and not subject to the same financial and regulatory pressures. Digital tools, including search engines and language models, will likely further depress these costs.

Existing work documents patterns that are not inconsistent with those in this paper. However, we offer an alternative interpretation. Like [Park et al. \(2023\)](#) and [Bloom et al. \(2020\)](#), we find a substantial increase in the size of the medical research ecosystem. [Park et al. \(2023\)](#) document that this increase in scale brings with it a decline in the average importance of published scientific articles. Technically, we find the same. As increases in non-clinical trials outpace (non-existent) changes in clinical trials, the average productivity of studies in Figure 3 mechanically decreases. In debates on productivity, however, we argue that this measure of central tendency is not the policy-relevant object. To be sure, these findings raise a host of other important questions—about misaligned incentives, wasteful spending, and the productive efficiency of science. But canonical frameworks, such as [Romer \(1990\)](#), center on the simple idea that it is high-intensity investments that generate meaningful shifts in outcomes. An explosion of low-capital research that has a negligible impact on outcomes, then, is not cause for concern *from a growth perspective*. That is, growth-oriented studies may prefer to focus—like much of the research on science and innovation—on the right-tail of the productivity distribution (e.g., [Kelly et al., 2021](#)).

Our findings do not speak to shifts in upstream, basic science inputs into clinical trials, or to shifts in factor prices. In our view, there is little credible evidence indicating changes on either dimension during this time period. This setting allows us, only, to note that any shift in inputs or factor prices must be perfectly offset by a countervailing force, to account for the complete stability of clinical trial counts, citations, journal placements, and reliance on public funding in our data.

Our findings, also, apply to roughly one decade: 2010 to 2022. We impose this restriction because policy changes in the early 2000s substantially shifted incentives for both publication and standardized formatting of clinical trial results. See [Laine et al. \(2007\)](#) and [Schulz et al. \(2010\)](#) for details. Given our methods, these changes mean that data construction over a longer period risks decreasing the quality of our data. Thus, we cannot speak to trends that precede this period.

That clinical trial investment is remarkably stable in the past decade, however, is likely the relevant consideration for contemporary public policy.

What our findings, and the simple framework borrowed from the growth literature, do suggest is that there is likely value in efforts to increase the productivity and decrease the costs of clinical trials. A long-standing discussion in medicine and statistics, spurred by [Altman \(1994\)](#), centers on the argument that “we need less research, better research, and research done for the right reasons.” Our data indicate that nearly half of all clinical trials published in our sample are never cited by leading medical journals, and roughly 15 percent are never cited at all. This heterogeneity is consistent with findings in [Ioannidis \(2016b\)](#) and [van Zwet et al. \(2023\)](#), who imply that there may exist trade-offs between quantity and data quality in this setting. If increasing the scale of research infrastructure is a policy objective, such efforts should invest in tools, like patient registries, that enable the production of high-quality, representative information ([Alsan et al., 2024](#); [Durvasula, 2023](#)).

REFERENCES

- Accreditation Council for Graduate Medical Education (2017). ACGME common program requirements. Accessed: 2024-05-05.
- Acemoglu, D., Autor, D., and Patterson, C. (2024). Bottlenecks: Sectoral imbalances and the us productivity slowdown. *NBER Macroeconomics Annual*, 38(1):153–207.
- Acemoglu, D. and Johnson, S. (2007). Disease and development: The effect of life expectancy on economic growth. *Journal of Political Economy*, 115(6):925–985.
- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. (2023). GPT-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Agha, L. and Molitor, D. (2018). The local influence of pioneer investigators on technology adoption: Evidence from new cancer drugs. *Review of Economics and Statistics*, 100(1):29–44.
- Aghion, P., Akcigit, U., Bergeaud, A., Blundell, R., and Hémous, D. (2019). Innovation and top income inequality. *The Review of Economic Studies*, 86(1):1–45.
- Aghion, P. and Howitt, P. (1992). A model of growth through creative destruction. *Econometrica: Journal of the Econometric Society*, pages 323–351.
- Almond, D., Doyle Jr, J. J., Kowalski, A. E., and Williams, H. (2010). Estimating marginal returns to medical care: Evidence from at-risk newborns. *The quarterly journal of economics*, 125(2):591–634.
- Alsan, M., Durvasula, M., Gupta, H., Schwartzstein, J., and Williams, H. (2024). Representation and extrapolation: Evidence from clinical trials. *The Quarterly Journal of Economics*, 139(1):575–635.
- Altman, D. G. (1994). The scandal of poor medical research. *Bmj*, 308(6924):283–284.
- Andrews, I. and Kasy, M. (2019). Identification of and correction for publication bias. *American Economic Review*, 109(8):2766–2794.
- Angrist, J., Azoulay, P., Ellison, G., Hill, R., and Lu, S. F. (2020). Inside job or deep impact? extramural citations and the influence of economic scholarship. *Journal of Economic Literature*, 58(1):3–52.
- Bloom, N., Jones, C. I., Van Reenen, J., and Webb, M. (2020). Are ideas getting harder to find? *American Economic Review*, 110(4):1104–1144.
- Broad, W. J. (2023). What happened to all of science’s big breakthroughs? *The New York Times*.
- Bryan, K. A. and Williams, H. L. (2021). Innovation: Market failures and public policies. In *Handbook of industrial organization*, volume 5, pages 281–388. Elsevier.
- Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., Lee, P., Lee, Y. T., Li, Y., Lundberg, S., et al. (2023). Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*.
- Chavez-MacGregor, M. and Giordano, S. H. (2016). Randomized clinical trials and observational studies: Is there a battle? *Journal of Clinical Oncology*, 34(8):772–773.

- Cockburn, I. M. (2006). Is the pharmaceutical industry in a productivity crisis? *Innovation policy and the economy*, 7:1–32.
- Collison, P. and Nielsen, M. (2018). Science is getting less bang for its buck. *The Atlantic*.
- Cortellis (2024). Cortellis Labs. <https://www.cortellislabs.com/page/?api=api-CLI>. Accessed on: 2024-05-04.
- Council of Economic Advisers (2021). Economic report of the president. U.S. Government Printing Office.
- Council of Economic Advisers (2024). Economic report of the president. U.S. Government Printing Office.
- Cowen, T. (2011). *The Great Stagnation: How America Ate All the Low-Hanging Fruit of Modern History, Got Sick, and Will (Eventually) Feel Better*. Penguin. A Penguin eSpecial from Dutton.
- Cutler, D. M. and McClellan, M. (2001). Is technological change in medicine worth it? *Health affairs*, 20(5):11–29.
- DeVito, N. J., Bacon, S., and Goldacre, B. (2020). Compliance with legal requirement to report clinical trial results on clinicaltrials.gov: A cohort study. *The Lancet*, 395(10221):361–369.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Dieppe, A. and Kose, M. A. (2020). The global productivity slump: What policies to rekindle? *Brookings*.
- DiMasi, J. A., Grabowski, H. G., and Hansen, R. W. (2016). Innovation in the pharmaceutical industry: New estimates of R&D costs. *Journal of health economics*, 47:20–33.
- Dinopoulos, E. and Thompson, P. (1998). Schumpeterian growth without scale effects. *Journal of Economic Growth*, 3(4):313–335.
- Durvasula, M. M. (2023). Inclusive infrastructure for clinical trials. *Brookings: Building a Better NIH*.
- Durvasula, M. M., Ouellette, L. L., and Williams, H. L. (2021). Private and public investments in biomedical research. In *AEA Papers and Proceedings*, volume 111, pages 341–345. American Economic Association 2014 Broadway, Suite 305, Nashville, TN 37203.
- Easterlin, R. A. (1999). How beneficent is the market? A look at the modern history of mortality. *European Review of Economic History*, 3(3):257–294.
- Feldman, S., Ammar, W., Lo, K., Trepman, E., van Zuylen, M., and Etzioni, O. (2019). Quantifying sex bias in clinical studies at scale with automated data extraction. *JAMA network open*, 2(7):e196700–e196700.
- Food and Drug Administration (2024). What are the different types of clinical research? Accessed: 2024-05-04.
- Goldin, I., Koutroumpis, P., Lafond, F., and Winkler, J. (2024). Why is productivity slowing down? *Journal of Economic Literature*, 62(1):196–268.
- Gordon, R. J. (2017). *The Rise and Fall of American Growth: The US Standard of Living since the Civil War*. Princeton University Press.

- Griliches, Z. (1990). Patent statistics as economic indicators: A survey. *Journal of Economic Literature*, 28(4).
- Gruber, J. and Johnson, S. (2019). *Jump-starting America: How breakthrough science can revive economic growth and the American dream*. Hachette UK.
- Gu, Y., Tinn, R., Cheng, H., Lucas, M., Usuyama, N., Liu, X., Naumann, T., Gao, J., and Poon, H. (2021). Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23.
- Guu, K., Lee, K., Tung, Z., Pasupat, P., and Chang, M.-W. (2020). REALM: Retrieval-augmented language model pre-training. *arXiv preprint arXiv:2002.08909*.
- Hausmann, R. (2023). Why industrial policy is back. *Project Syndicate*.
- Hill, R. and Stein, C. (2021). Race to the bottom: Competition and quality in science. *Northwestern University and UC Berkeley*.
- Howitt, P. (1999). Steady endogenous growth with population and R& D inputs growing. *Journal of Political Economy*, 107(4):715–730.
- Ioannidis, J. P. (2016a). The mass production of redundant, misleading, and conflicted systematic reviews and meta-analyses. *The Milbank Quarterly*, 94(3):485–514.
- Ioannidis, J. P. (2016b). Why most clinical research is not useful. *PLoS medicine*, 13(6):e1002049.
- Ioannidis, J. P., Pezzullo, A. M., and Boccia, S. (2023). The rapid growth of mega-journals: Threats and opportunities. *Jama*, 329(15):1253–1254.
- Jones, B. F. (2009). The burden of knowledge and the “death of the renaissance man”: Is innovation getting harder? *The Review of Economic Studies*, 76(1):283–317.
- Jones, C. I. (1995). R&D-based models of economic growth. *Journal of Political Economy*, 103(4):759–784.
- Jones, D. S. and Podolsky, S. H. (2015). The history and fate of the gold standard. *The Lancet*, 385(9977):1502–1503.
- Kao, J., Ross, J. S., and Miller, J. E. (2023). Transparency of results reporting in cancer clinical trials. *JAMA Network Open*, 6(8):e2328117–e2328117.
- Kelly, B., Papanikolaou, D., Seru, A., and Taddy, M. (2021). Measuring technological innovation over the long run. *American Economic Review: Insights*, 3(3):303–320.
- Kolata, G. (2022). FDA Panel Rejects Lilly’s Cancer Drug Tested Only in China. *The New York Times*.
- Kones, R., Rumana, U., and Merino, J. (2014). Exclusion of ‘nonRCT evidence’ in guidelines for chronic diseases—is it always appropriate? the look ahead study. *Current Medical Research and Opinion*, 30(10):2009–2019.
- Laine, C., Horton, R., DeAngelis, C. D., Godlee, F., Drazen, J. M., Frizelle, F. A., Haug, C., Hébert, P. C., et al. (2007). Update on trials registration: Clinical trial registration: Looking back and moving ahead. *International Committee of Medical Journal Editors*.
- Li, D. (2017). Expertise versus bias in evaluation: Evidence from the NIH. *American Economic Journal: Applied Economics*, 9(2):60–92.

- Lichtenberg, F. R. (2018). The impact of biomedical research on US cancer mortality. *Measuring and Modeling Health Care Costs*, 76:475.
- Lindsey, B. and Hammond, S. (2020). Faster growth, fairer growth: Policies for a high road, high performance economy. *Niskanen Center*.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Mokyr, J. (2014). Secular stagnation? not in your life. *Secular stagnation: facts, causes and cures*, 83.
- Mokyr, J. et al. (2002). Innovation in an historical perspective: Tales of technology and evolution. *Technological Innovation and Economic Performance*, 23:36.
- Myers, K. (2020). The elasticity of science. *American Economic Journal: Applied Economics*, 12(4):103–134.
- Myers, K. and Pauly, M. (2019). Endogenous productivity of demand-induced R&D: Evidence from pharmaceuticals. *The RAND Journal of Economics*, 50(3):591–614.
- Namba, M., Kaneda, Y., Ozaki, A., and Tanimoto, T. (2024). Clinical trials: Japan’s opt-out policy raises risks of adverse drug responses. *Nature*, 626(7998):261–261.
- Nori, H., King, N., McKinney, S. M., Carignan, D., and Horvitz, E. (2023). Capabilities of GPT-4 on medical challenge problems. *arXiv preprint arXiv:2303.13375*.
- Oeppen, J. and Vaupel, J. W. (2002). Broken limits to life expectancy. *Science*, 296(5570):1029–1031.
- Oostrom, T. (2024). Funding of clinical trials and reported drug efficacy. *Journal of Political Economy*, Forthcoming.
- Ouellette, L. L. and Sampat, B. N. (2024). The feasibility of using Bayh-Dole march-in rights to lower drug prices: An update. Technical report, National Bureau of Economic Research.
- Pakes, A. (1986). Patents as options: Some estimates of the value of holding european patent stocks. *Econometrica: Journal of the Econometric Society*, pages 755–784.
- Park, M., Leahey, E., and Funk, R. J. (2023). Papers and patents are becoming less disruptive over time. *Nature*, 613(7942):138–144.
- Peretto, P. F. (1998). Technological change and population growth. *Journal of Economic Growth*, 3(4):283–311.
- Petryna, A. (2007). Clinical trials offshored: On private sector science and public health. *BioSocieties*, 2(1):21–40.
- Piantadosi, S. (2024). *Clinical trials: a methodologic perspective*. John Wiley & Sons.
- Piper, K. (2023). Why is science slowing down? *Vox*.
- Qiao, Y., Alexander, G. C., and Moore, T. J. (2019). Globalization of clinical trials: Variation in estimated regional costs of pivotal trials, 2015-2016. *Clinical Trials*, 16(3):329–333.
- Reimers, N. and Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using siamese bert-networks.

- Romer, P. M. (1990). Endogenous technological change. *Journal of Political Economy*, 98(5):2.
- Scannell, J. W., Blanckley, A., Boldon, H., and Warrington, B. (2012). Diagnosing the decline in pharmaceutical R&D efficiency. *Nature Reviews Drug Discovery*, 11(3):191–200.
- Schulz, K. F., Altman, D. G., and Moher, D. (2010). CONSORT 2010 statement: Updated guidelines for reporting parallel group randomised trials. *Journal of Pharmacology and pharmacotherapeutics*, 1(2):100–107.
- Siegel, M. G., Brand, J. C., Rossi, M. J., and Lubowitz, J. H. (2018). “Publish or perish” promotes medical literature quantity over quality. *Arthroscopy: The Journal of Arthroscopic & Related Surgery*, 34(11):2941–2942.
- Syverson, C. (2017). Challenges to mismeasurement explanations for the US productivity slowdown. *Journal of Economic Perspectives*, 31(2):165–86.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. (2023). Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- Thomas, J., McDonald, S., Noel-Storr, A., Shemilt, I., Elliott, J., Mavergames, C., and Marshall, I. J. (2021). Machine learning reduced workload with minimal risk of missing studies: Development and evaluation of a randomized controlled trial classifier for cochrane reviews. *Journal of Clinical Epidemiology*, 133:140–151.
- Thompson, D. (2021). America is running on fumes. *The Atlantic*.
- Trinh, T. H. and Le, Q. V. (2019). A simple method for commonsense reasoning. *arXiv preprint arXiv:1806.02847*.
- van Zwet, E., Gelman, A., Greenland, S., Imbens, G., Schwab, S., and Goodman, S. N. (2023). A new look at p values for randomized clinical trials. *NEJM evidence*, 3(1):EVIDoA2300003.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2023). Attention is all you need.
- Wikimedia Foundation (2023). Wikimedia downloads.
- Wouters, O. J., McKee, M., and Luyten, J. (2020). Estimated research and development investment needed to bring a new medicine to market, 2009-2018. *Jama*, 323(9):844–853.
- Young, A. (1998). Growth without scale effects. *Journal of political economy*, 106(1):41–63.
- Zaheer, M., Guruganesh, G., Dubey, A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., and Ahmed, A. (2021). Big Bird: Transformers for longer sequences. *arXiv preprint arXiv:2007.14062*.
- Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., and Fidler, S. (2015). Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. *arXiv preprint arXiv:1506.06724*.

Supplemental Appendix to:
Medical Research as a Productivity Indicator*

Maya M. Durvasula*

Sabri Eyuboglu†

David M. Ritzwoller‡

Contents

Appendix A. Publication Data	1
A.1. The PubMed / MEDLINE Database	1
A.2. The Universe of Potential Clinical Trials	2
A.3. Determining the Set of Leading Medical Journals	6
A.4. The Web of Science	6
Appendix B. Prompt Design, Fine-Tuning, and Performance Assessment	6
B.1. Hand-Labeling	6
B.2. Benchmark Comparison to Standard Machine Learning Methods	7
B.3. Prompt Design and Error Analysis	8
B.4. Fine-tuning	13
B.5. Performance Assessment	13
B.6. Final Model: Error Analysis	15
Appendix C. Additional Figures and Further Analyses	15
Appendix D. Prompt Repository	21

*Date: May 15, 2024

★ Stanford Department of Economics and Stanford Law School, maya.durvasula@stanford.edu

† Stanford Department of Computer Science, eyuboglu@stanford.edu

‡ Stanford Graduate School of Business, ritzwo11@stanford.edu

APPENDIX A. PUBLICATION DATA

In this Appendix, we detail our treatment of the PubMed / MEDLINE database (hereafter, “PubMed”). We give a general overview of the PubMed data in [Appendix A.1](#). In [Appendix A.2](#), we describe the process of identifying our initial sample of abstracts, i.e., the “Potential Universe” displayed in [Figure 4](#). We discuss the process of determining the set of leading medical journals in [Appendix A.3](#).

A.1 The PubMed / MEDLINE Database. Our version of the PubMed / MEDLINE database contains roughly 34 million records and is current through December 2022. These records were constructed by parsing bulk MEDLINE XML files.³⁴ Technically, PubMed and MEDLINE are different products. MEDLINE, a subset of PubMed, is the U.S. National Library of Medicine’s (NLM) bibliographic database, which contains references to journal articles in life sciences, with a primary focus on biomedicine. A committee at the NLM determines the set of indexed journals, meaning that only journals that meet certain quality and content standards are indexed. In practice, this means that so-called predatory journals are excluded, as are pre-prints and non-peer reviewed articles. Informal conversations with staff at the NLM suggest that MEDLINE *should* contain the universe of peer-reviewed publications in legitimate journals. PubMed includes a broader set of records, including pre-prints and publications deposited through alternative processes. Official documentation for MEDLINE is available at the link: https://www.nlm.nih.gov/databases/download/pubmed_medline_documentation.html. Throughout the text, and from this point forward in the Appendix, we refer to this database as “PubMed.”

A natural question for researchers who rely on these data as a census of research investments is whether PubMed records are complete and accurate. To our knowledge, there is no paper that reports such validation exercises for each field of PubMed. We focus on those data elements most relevant to our work. There are 2,335,653 records in PubMed that have no associated year of publication, which we drop from our baseline sample. There are 16 records, across all years of data, that are not indexed with tags describing their contents (NLM tags). From 2010 forward, roughly 92 percent of records have associated abstracts. We randomly inspect roughly 150 records with no abstracts. None correspond to publications that satisfy our definition of a clinical trial.

We assume, throughout this paper, that records with missing citation data have zero associated citations. We confirm that this is generally accurate in two ways. First, we search for Google Scholar records corresponding to roughly 100 randomly selected records with zero citations. In each case, Google Scholar indicates that the paper has no more than two citations. Second, we link records in PubMed to Clarivate Analytics’ Web of Science database. We compare citation counts constructed using information in the two databases. We find that the two measures have a correlation of 0.9.

³⁴We are grateful to Heidi Williams for sharing this processed data.

Researchers interested in different cuts of the data may need to conduct additional validation exercises. Figure 1 highlights that the frequencies of certain NLM tags have changed over time, often sharply. Researchers should account for such changes in any use of these flags. For other assessments of the completeness of certain PubMed fields, see [Durvasula et al. \(2021\)](#) and [Ouellette and Sampat \(2024\)](#).

A.2 The Universe of Potential Clinical Trials. We construct our sample using records drawn from the universe of publications indexed in PubMed. The initial sample consists of 34,957,127 unique publications indexed in PubMed, as of 15 April 2023. We drop 2,335,653 records that are missing information on publication year, to yield a base sample of 32,621,474 records. From 2010 forward, at least 92 percent of scientific publications published in each year have associated abstracts. We also drop 14,179 publications with publication year 2023, as we have incomplete data for 2023.

A.2.1 NLM Tags. The National Library of Medicine (NLM) assigns each publication a ‘pubtype.’ In the entirety of PubMed, there are 16 records missing pubtype tags. To the best of our knowledge, there have been no efforts to validate the PubMed indexing process used to generate these flags. We flag records with each of the following pubtype (or associated unique identifiers, reported in the field ‘pubtypeUI’) tags:

Adaptive trial; Clinical conference; Clinical study; Clinical trial; Clinical trial protocol; Clinical trial, Phase 1; Clinical trial, Phase 2; Clinical trial, Phase 3; Clinical trial, Phase 4; Comparative study; Controlled clinical trial; Equivalence trial; Evaluation study; Observational study; Pragmatic clinical trial; Randomized controlled trial; Twin study; Validation study.

These categories are chosen to include all pubtypes with the potential to contain a publication satisfying the restriction enumerated in [Definition 3.1](#). In particular, we follow [Feldman et al. \(2019\)](#), who use multiple pubtype fields to retrieve a sample of records. Here, two authors reviewed the list of potential pubtypes to identify categories likely to include records of interest. [Table A.1](#) reports the frequency of each NLM tag across all records in PubMed and subset to those published after 1 January 2010.

A.2.2 Clinical Trial Registry Identifiers. In 2004, the International Committee of Medical Journal Editors recommended that research journals decline to publish outcomes associated with trials not pre-registered in some repository. [Laine et al. \(2007\)](#) summarizes these policies, based on a 2007 revision, in more detail. The U.S. Food and Drug Administration Amendments Act of 2007, Section 801, mandates registration of all clinical trials regulated by the FDA in [ClinicalTrials.gov](#). Many countries have adopted similar guidance. Several countries and international organizations now maintain registries of trials.

TABLE A.1. Composition of PubMed by NLM Tag

	A. All records		B. 2010-2022		C. Conserv. Sample	
	Frequency	%	Frequency	%	Frequency	%
adaptive trial	36	0.00	36	0.00	11	0.01
clinical conference	7,045	0.02	1,724	0.01	12	0.01
clinical study	5,053	0.02	5,050	0.04	641	0.42
clinical trial	498,722	1.53	79,029	0.55	13,969	9.19
clinical trial protocol	9,542	0.03	9,542	0.07	124	0.08
clinical trial, phase 1	23,704	0.07	14,139	0.10	10,951	7.20
clinical trial, phase 2	37,623	0.12	22,749	0.16	17,166	11.29
clinical trial, phase 3	20,597	0.06	15,647	0.11	11,112	7.31
clinical trial, phase 4	2,251	0.01	1,790	0.01	1,250	0.82
comparative study	1,752,380	5.37	424,681	2.97	22,521	14.82
controlled clinical trial	88,132	0.27	13,839	0.10	2,181	1.43
equivalence trial	1,047	0.00	1,047	0.01	505	0.33
evaluation study	243,290	0.75	123,897	0.87	998	0.66
observational study	127,461	0.39	127,370	0.89	5,996	3.94
pragmatic clinical trial	2,112	0.01	2,112	0.01	253	0.17
randomized controlled trial	549,711	1.69	284,774	1.99	79,907	53.57
twin study	9,030	0.03	4,964	0.03	14	0.01
validation study	101,817	0.31	62,261	0.44	495	0.33
	<i>N</i> =32,621,474		<i>N</i> =14,316,494		<i>N</i> =151,997	

Notes: Table A.1 reports the frequency and percentage of records indexed in PubMed that have been categorized by the NLM as falling into each of 18 categories. These categories are selected for their potential to contain a publication satisfying the restriction enumerated in Definition 3.1.

Registry identifiers are distinctive strings of letters and numbers. We flag records that include clinical trial registry identifiers in their abstract text. In particular, we search for records containing acronyms associated with the following registries:

ClinicalTrials.gov (NCT); European Union Drug Regulating Authorities Clinical Trials Database (EUDRACT); International Traditional Medicine Clinical Trial Registry (ISRCTN); Australian New Zealand Clinical Trials Registry (ACTRN).

We note the potential to overcount records using these searches. We focus on instances where each trial identifier prefix is followed by numbers, letters, or punctuation (e.g., NCT12345, ISRCTN: 12345). However, we may collect records that include these characters in other settings (e.g., the world distinct). Table A.2 reports the frequency of each registry identifier flagged across all records in PubMed and subset to those published after 1 January 2010. Figure 1 plots the number of trials with any registry identifier over time. Reassuringly, we observe registry identifiers in PubMed data

TABLE A.2. Composition of PubMed by Trial Registry Identifiers in Abstract Text

	A. All records		B. 2010-present		C. Conserv. Sample	
	Frequency	%	Frequency	%	Frequency	%
NCT	97,923	0.30	94,279	0.66	30,656	20.17
EUDRACT	2,835	0.0	2,816	0.00	1,564	1.03
ISRCTN	12,261	0.03	11,325	0.07	1,275	0.84
ACTRN	5,818	0.01	5,656	0.03	553	0.36
	<i>N</i> =32,621,474		<i>N</i> =14,316,494		<i>N</i> =151,997	

Notes: Table A.2 reports the frequency and percentage of records indexed in PubMed that include a string associated with a clinical trial registry identifier in their abstract text. That is, the first row counts the number and percentage of publications that contain the string “nct.”

beginning around 2010, as registration mandates were implemented, and find a steady, smooth increase over time.

A.2.3 Keywords. We flag records that include any keyword likely to indicate that the record reports the results of a clinical trial in their abstract text. We selected these keywords after several reviewing roughly 200 abstracts flagged through the pubtype process, described above.

These keywords are:

Randomized; Controlled trial; Control trial; Clinical trial; Treatment group; Control group; Intervention; Clinical study.

Table A.3 reports the frequency of each keyword flagged across all records in across all records in PubMed and subset to those published after 1 January 2010.

A.2.4 Intersection. There are 3,925,958 records with at least one of the following attributes: a clinical-trial indicative NLM tag; a clinical trial registry identifier; a clinical-trial indicative keyword in the abstract text. The three categories overlap. Table A.4 records the size of the overlap of each category. In the main text, we restrict attention to this sample for the years 2010-2022. This sample includes 1,821,429 publications.

A.2.5 Novel Census. We add columns to Tables A.1 to A.4 that document the frequency of each flag in the data constructed for this paper (we use the conservative sample). Several facts are worth noting. First, Table A.2 suggests that—between 2010 and 2022—less than 25 percent of records identified as reporting the results of a clinical trial reported a registry identifier in their abstract. Although, in principle, registry identifiers may be reported in publication full-texts, rather than abstracts, the low frequency of identifiers in our sample suggests non-compliance with requirements, consistent with patterns in ClinicalTrials.gov documented in DeVito et al. (2020). Second, Table A.1 sheds light on the differences in trial composition that we observe in Figure 1: NLM tags that,

TABLE A.3. Composition of PubMed by Keywords in Abstract Text

	A. All records		B. 2010-present		C. Conserv. Sample	
	Frequency	%	Frequency	%	Frequency	%
randomized	550,165	1.69	366,378	2.56	65,637	43.18
controlled trial	107,935	0.33	80,789	0.56	11,535	7.59
control trial	6,049	0.02	4,846	0.03	371	0.24
clinical trial	129,510	0.40	92,035	0.64	17,644	11.61
treatment group	44,627	0.14	27,826	0.19	4,296	2.63
control group	447,901	1.37	286,250	2.00	14,245	9.37
intervention	629,237	1.93	463,373	3.24	10,460	6.88
clinical study	32,652	0.10	18,589	0.13	2,411	1.59
	<i>N</i> =32,621,474		<i>N</i> =14,316,494		<i>N</i> =151,997	

Notes: Table A.3 reports the frequency and percentage of records indexed in PubMed that include each of a collection of keywords in their abstract text. That is, the first row counts the number and percentage of publications that contain the string “randomized.”

TABLE A.4. Overlap of Publication Attributes

Records with any:	that have any:	A. All records	B. 2010-present	C. Conserv. Sample
Registry ID		116,564	111,837	33,040
	NLM tag	81,894	78,366	26,885
	Keyword	73,164	70,414	19,103
NLM tag		2,858,924	1,804,081	124,260
	Registry ID	81,894	3,528	26,885
	Keyword	515,755	220,681	67,153
Keyword		1,564,659	1,044,461	91,349
	Registry ID	73,164	70,414	19,103
	NLM tag	515,755	295,074	67,153

Notes: Table A.4 reports the size of the intersection between three sets of records indexed in PubMed. These categories are defined by the properties: Possession of a clinical-trial indicative NLM tag, inclusion of a clinical trial registry identifier in abstract text, and inclusion of a clinical-trial indicative keyword in abstract text. These properties are defined in Appendices A.2.1 to A.2.3, respectively.

ostensibly, capture clinical trials have little overlap with one another and do not fully capture the records in our sample. Of the roughly 150,000 records in our census, approximately 100,000 could be identified using NLM tags for “clinical trial” (and its variants) and “randomized controlled trial.” Observe, however, that use of such flags captures a much larger number of records than one would intend to include. Finally, Table A.4 suggests that although various text-based flags that may indicate a record is a clinical trial, in the sense of Definition 3.1, are correlated, there is less overlap than one might expect.

A.3 Determining the Set of Leading Medical Journals. Angrist et al. (2020) propose a strategy to identify leading journals in a field, which we adapt to this context. In medicine, the “trunk journals,” per Angrist et al. (2020), are the Journal of the American Medical Association and the New England Journal of Medicine. We collect all citations originating from these journals to records in PubMed between 2010 and 2022. Angrist et al. (2020) collect initial journal lists for their citation exercise of the journal is one of the top fifty most cited by a trunk journal. As medicine includes a large number of subfields, we instead collect a list of journals that receive at least 100 citations from a trunk journal over this time period, which yields a list of 84 journals.

A.4 The Web of Science. We collect supplementary information about papers in our sample from Clarivate Analytics’ Web of Science.³⁵ Web of Science data allow us to examine more detailed information about each author in our sample. While PubMed includes some details on authors, including limited institutional affiliation and address information, the Web of Science provides standardized author addresses for roughly 75 percent of observations in our sample.

For each paper in our sample, we extract countries from available mailing addresses for first- and last-listed authors. First- and last-authors have the same listed country for 85 percent of clinical trial records and 91 percent of non-clinical trial records. Thus, we impute to each record the country associated with the first-listed author. In Figure 6, we plot geographic trends over time.

APPENDIX B. PROMPT DESIGN, FINE-TUNING, AND PERFORMANCE ASSESSMENT

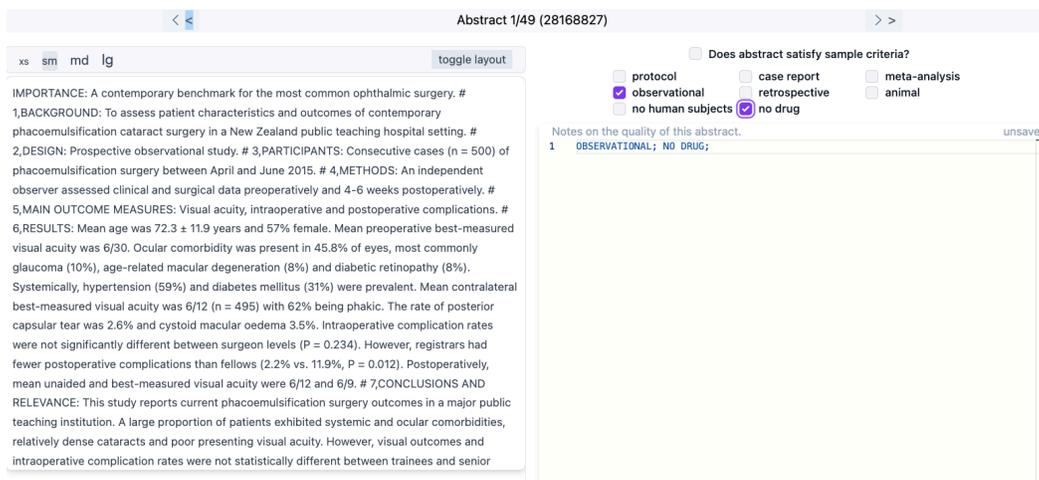
In this appendix, we detail the construction of the sample of clinical trials studied in Section 4. We begin in Appendix B.1 by describing the process we used to hand-label a sample of publications according to whether they satisfy the restrictions enumerated in Definition 3.1. In Appendix B.2, we document poor performance of several approaches for classifying publications according to whether they belong in our sample based on standard machine learning methods. Appendix B.3 overviews the process of constructing prompts that perform well when used to query proprietary language models. Details associated with fine-tuning an open-source language model are given in Appendix B.4. Additional measurements of the performance of our model are given in Appendix B.5.

B.1 Hand-Labeling. A standard concern when using large language models is the quality of the output. We hand-label a set of 3,000 abstracts, selected randomly from our candidate set of records, which allow us to quantify the performance of our classification procedure.

We devise an interface for abstract labeling, which allows us to review records quickly. Figure B.1 provides an example of this interface. A labeller reads in a set of publication records and is shown one PubMed identifier (PMID) and one abstract at a time. The labeller can select that the abstract satisfies our sample criteria, or else can indicate a reason that the abstract should be excluded. The

³⁵We use a copy of the Web of Science licensed to Stanford University.

FIGURE B.1. Abstract Labeling Interface



Notes: Figure B.1 displays a screenshot of the application we developed to hand-label abstracts. Each abstract was labelled according to whether it satisfied the sample criteria enumerated in Definition 3.1. If a publication did not satisfy our sample restrictions, we stipulated a reason. The user interface was developed with the “Meerkat” Python package available at the link: <https://github.com/HazyResearch/meerkat>.

buttons for exclusion mirror the inclusion and exclusion criteria in Definition 3.1, using short-hand convenient for labellers. We exclude records that report, only, the protocol for a clinical trial without results (“protocol”), that report the results of an observational (“observational”) or retrospective (“retrospective”) study, that report the results of a study involving animals (“animal”), that do not include human subjects (“no human subjects”),³⁶ that do not study a drug (“no drug”), that report the results of a specific clinical case or patient experience (“case report”), or that summarize the findings of existing work (“meta-analysis”). By our own estimates, the labelling interface increased the speed of labelling by a factor of seven. Each record was reviewed, and labelled, twice.

The hand-labeled data are split into three subsets—validation, training, and testing—based on their eventual use. Splits are assigned randomly. We assign 1000 records to a test set, 1000 to a validation set, and 1082 to a training set. Of these records, we assign labels only to records that have abstracts. Our final training data includes 1082 records, our validation dataset includes 1000 records, and our test set includes 993 records. The loss of seven records in our test set reflects a coding error. We add additional records, randomly selected, to the training and validation datasets when abstracts are missing. We do not add extra records to the testing dataset.

B.2 Benchmark Comparison to Standard Machine Learning Methods. We measure the performance of a variety of standard machine learning algorithms for classifying publications

³⁶In practice, “animal” and “no human subjects” catch different sets of records. We flag records that include *any* animal with “animal.” Many studies in our sample of candidate records occur *in vitro*—in test tubes or other laboratory-based settings. There, no living subjects of any kind are enrolled and the “no human subjects” button is used.

according to whether they meet our sample restrictions. Models are trained in the training and validation samples and tested in the testing sample.

We compute two embeddings for each abstract in our hand-labeled sample. First, we compute Term Frequency, Inverse Document Frequency (TF-IDF) embeddings based on the corpus of abstracts in the hand-labeled sample. Second, we extract embeddings associated with the abstract of each publication in the hand-labeled sample with the SENTENCETRANSFORMERS language model (Reimers and Gurevych, 2019). The SENTENCETRANSFORMERS language model is a standard source for embeddings, and is a modification of the BERT language model (Devlin et al., 2019).

We consider six classes of prediction algorithms: logistic ridge regression, logistic lasso regression, Support Vector Machines, Random Forest Regression, Boosted Regression, and Convolutional neural networks. Each model is obtained from the “sklearn” Python package. The hyper-parameters of each model are chosen with 10-fold cross-validation implemented in the training and validation samples. A balanced loss is used to measure the performance of each hyper-parameter.

Figure 2 displays receiver operating characteristics for the each class of models estimated with both types of embedding. Error rates are estimated in the testing data. The performance of the receiver operating characteristic for the “ensemble” model whose performance is displayed in Figure 3 is given in purple. The fine-tuned large language model developed in this paper substantially out-paces the performance of standard machine learning methods.

B.3 Prompt Design and Error Analysis. We extract weak labels for 64,000 randomly selected abstracts with two proprietary large language models: OpenAI’s GPT-3.5 and GPT-4 (Nori et al., 2023; Bubeck et al., 2023). *A priori*, we expect GPT-4 to output labels of slightly lower quality than hand-labeled output and GPT-3.5 to produce slightly noisier labels than GPT-4. The extraction of GPT-3.5 labels is substantially faster and cheaper.

To produce these labels at scale, each of the two models must be appropriately prompted. That is, each pre-trained model must be provided with a block of text as an input and asked to return the appropriate completion of this input. *A priori*, it is not clear what prompt structure will work well. The relative infancy of this area of research renders it difficult to identify a set of “best practices.” The proprietary nature of these models makes them difficult to study. See, for details, a guide to prompt engineering from OpenAI here: <https://platform.openai.com/docs/guides/prompt-engineering>. All of the prompts that we consider are displayed as figures in Appendix D.

We identify three general prompt formats, which differ both in the amount of detail provided about our classification task and in the structure of the requested model completion. The simplest prompt provides a version of our sample definition, Definition 3.1, and the text of an abstract, and asks the model to return ‘TRUE’ if the abstract satisfies these criteria. Otherwise, it returns false.

The text of this prompt is displayed in [Figure D.8](#). We refer to this prompt as Prompt 1.0. A second, more complicated prompt provides the same definition, with a set of examples of publication characteristics that do and do not satisfy the definition. The text of this prompt is displayed in [Figure D.12](#). We refer to this prompt as Prompt 2.0. A third, more complex prompt provides the same definitions and examples, but asks the model to return an explanation of why the record does or does not satisfy these criteria. The text of this prompt is displayed in [Figure D.14](#). We refer to this prompt as Prompt 3.0. We devise initial language for each prompt iteratively, using small numbers of records from our hand-labelled validation data.

We test each of these three prompts in our 1000-record hand-labelled validation dataset, using both GPT-3.5 and GPT-4. We conduct a detailed error analysis, reported in [Table B.2](#). For each instance in which a model returns a label (TRUE/FALSE) that differs from that in the hand-labelled dataset, we inspect the record. We categorize errors into types and subtypes. To make this concrete, suppose that GPT-3.5, given Prompt 1.0 and an abstract that reports the results of a meta-analysis, erroneously classifies the record as clinical trial, per our definition. In [Table B.2](#), we flag this error under Prompt 1 for GPT-3.5, as an error of type “meta-analysis.” Within meta-analysis, we assign the error to a sub-type. If the abstract explicitly includes terms such as “meta-analysis,” “literature search,” or “literature review,” we categorize this as an explicit error. If the abstract references that the publication is summarizing existing studies, or searching a database for records and collecting their findings, we categorize this as an implicit error. We devise error types based on our inclusion/exclusion criteria, in [Definition 3.1](#), and select sub-types based on common categories of errors.

This exercise is instructive on three margins. First, it highlights the highest performing prompts. Second, it indicates particular types of errors in categorization—which suggest opportunities for more precise language in a prompt. Third, it draws attention to differences in the performance of the two models. Observe that, in the “Other” error type, we include a subtype called “overly literal interpretation of inclusion criteria.” We primary record errors of this type for Prompt 3, which asked the model to return a completion that described why, or why not, a record was classified as being in our sample. Here, GPT-3.5 makes 30 such errors—classifying records as FALSE by recapitulating the sample definition provided. GPT-4 makes two errors. Insight that emerges from this exercise, then, is that different prompt structures may be preferable depending on the model used.

We revise each class of prompt based on these findings. For Prompt 1.0, the simplest true/false prompt, we consider three variants. We refer to these revised prompts as Prompts 1.1, 1.2, and 1.3. The text of these prompts is displayed in [Figures D.9 to D.11](#). For Prompts 2 and 3, we consider one variant each. We refer to these revised prompts as Prompts 2.1 and 3.1. The text of these prompts is displayed in [Figures D.13 and D.15](#). These changes reflect the differences in performance catalogued in [Table B.2](#).

TABLE B.1. Prompt Performance in Validation Data

<i>Panel A: GPT-3.5</i>			
Prompt Type	Prompt Sub-Type	False Positive Rate	True Positive Rate
1	0	0.202	0.971
	1	0.065	0.949
	2	0.049	0.934
	3	0.056	0.912
2	0	0.081	0.964
	1	0.072	0.964
3	0	0.167	0.971
	1	0.068	0.956
<i>Panel B: GPT-4</i>			
Prompt Type	Prompt Sub-Type	False Positive Rate	True Positive Rate
1	0	0.247	0.788
	1	0.172	0.898
	2	0.037	0.584
	3	0.037	0.489
2	0	0.162	0.876
	1	0.176	0.905
3	0	0.248	0.722
	1	0.171	0.883

Notes: [Table B.1](#) records the performance of each of eight prompt variants in a sample of 1,000 validation dataset records. Panel A reports performance associated with the proprietary model GPT-3.5. Panel B reports analogous statistics for model GPT-4. Prompt 1 asks the model to return TRUE or FALSE. Prompt 2 asks the model to return TRUE or the name of a specific, excluded category. Prompt 3 asks the model to return TRUE or an explanation of why the record should be excluded. Prompt sub-types correspond to various iterations. Sub-type 0 is the initial version of the prompt. Sub-types 1-3, where applicable, are subsequent iterations. [Appendix D](#) records the text of each prompt. We use Prompt 2, Sub-Type 0, to extract weak labels using GPT-3.5, and Prompt 2, Sub-Type 1 to extract weak labels using GPT-4.

[Appendix B.3](#) reports estimates of the true positive rate and false positive rate for each prompt in both models computed in the validation sample. Observe—for example, with Prompts 1.1, 1.2, and 1.3 queried to GPT-4—that even small changes in the text of a prompt yield substantial differences in performance. We attempt a second round of prompt iteration (Prompt 1.3). We observe that performance deteriorates with even small modification. Thus, we select the two highest performing prompts—one for each proprietary model. We use Prompt 2.0 to extract weak labels using GPT-3.5, and Prompt 1.2 to extract weak labels using GPT-4. These prompts are used to extract 64,000 labels—TRUE/FALSE flags—using GPT-3.5 and GPT-4. We refer to these as our set of “noisy labels.”

TABLE B.2. Baseline Prompt Error Analysis

Model	GPT-3			GPT-4		
	Prompt			Prompt		
Error Type	1.0	2.0	3.0	1.0	2.0	3.0
No Drug						
vitamin/supplement	12	12	9	15	4	15
surgical/medical procedure/diagnostic	40	27	29	35	21	29
abstract not specific about name of medicine	3	4	4	6	5	6
food/beverages	8	3	6	13	1	12
surgical/medical device	8	4	5	7	2	6
surgical/medical imaging	1	1	1	2	0	2
behavioral/physical therapy/exercise	40	22	36	25	8	13
misses the mention of a drug	25	10	5	2	1	0
surgical/medical material (e.g., resin/dental filling)	1	0	1	3	1	1
non-medical pollution/chemical/drug	1	1	1	1	2	1
other	1	0	1	2	1	2
Meta-Analysis						
explicit mention of meta-analysis/literature search/literature review	13	5	13	8	0	6
mention of database searched, explicit mention of meta-analysis	3	3	3	3	0	2
summary existing studies, without reference to data search / meta-analysis	6	7	4	3	4	3
Retrospective						
re-analysis of previously collected data without explicitly mentioning “retrospective”	8	2	6	6	5	4
explicitly described as “retrospective” or “retrospective analysis”	6	1	6	0	0	0
miscategorizes a study conducted in the past as “retrospective”	0	1	0	0	0	0
other						
Observational						
no active intervention described; does not explicitly say “observational”	10	9	11	6	0	4

no active intervention described; explicitly says “observational”	14	13	13	13	2	11
misses reference to active intervention	1	1	1	0	0	0
Protocol						
no results reported about current study; explicitly says “protocol”	10	9	11	6	0	4
no results reported about current study; does not explicitly say “protocol”	14	13	13	13	2	11
study is based on a simulation, not real world data	1	1	1	0	0	0
No Human Subjects						
only references to laboratory experiments, cells drawn from humans	0	1	0	0	0	0
Animal						
study conducted on animals	5	10	6	2	3	2
Other						
overly literal interpretation of inclusion criteria	0	2	30	0	0	2
studying outcomes/objects unrelated to a randomized trial	1	1	1	1	1	1

Notes: [Table B.2](#) categorizes errors documented in the first round of prompt iteration, using Prompts of sub-type 0. See [Table B.1](#) for details on model performance. For each model error—an instance in which a model returned a classification that deviated from hand-labelled data—we examined the associated record. We categorized errors by Type and Sub-Type. Types are drawn from the exclusion restrictions in [Definition 3.1](#). Sub-types describe consistent characteristics of publications that resulted in such an error. We report counts of errors, by model and prompt, of each Type and Sub-type.

B.4 Fine-tuning. This appendix gives technical details related to fine-tuning open-source large language models with noisy labels extracted from GPT-3.5 and GPT-4. We fine-tune pre-trained BERT models from two different architecture classes: (1) BIGBIRD and (2) BIOMEDBERT (Gu et al., 2021; Zaheer et al., 2021). Both models use medium-scale Transformer architectures, i.e., between 100 and 300 million parameters, pre-trained with the masked-language modeling (MLM) objective (Vaswani et al., 2023; Devlin et al., 2019). The models differ in their architectural details and pre-training corpora.

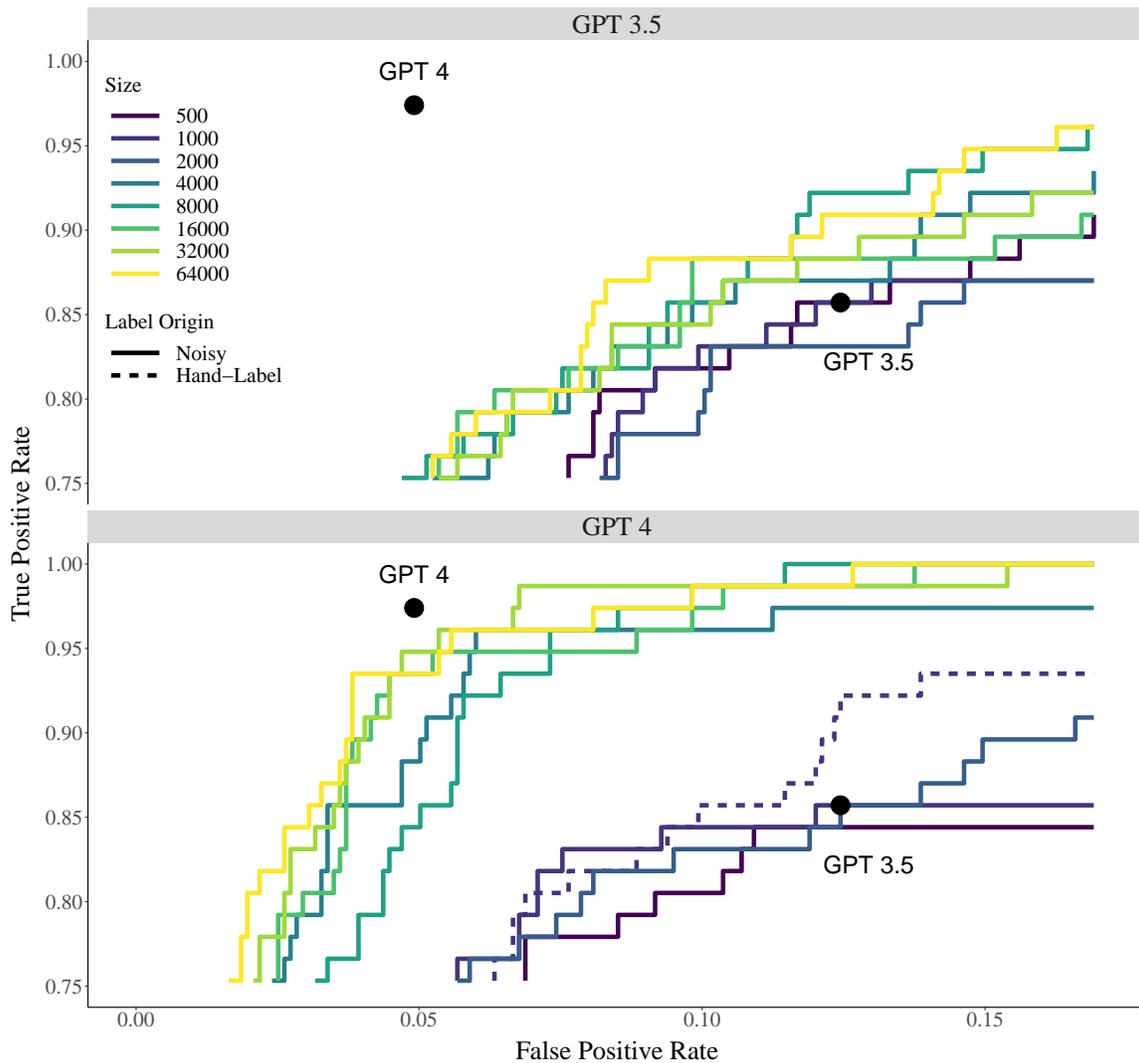
BIGBIRD uses a sparse attention mechanism to reduce the computational cost of processing long text sequences. The computational complexity of regular attention is quadratic in sequence length, while BIGBIRD reduces this to linear complexity. This enables efficient handling of long documents up to 4,096 tokens. In contrast, standard BERT can only process 512 tokens. BIGBIRD is useful because abstracts of clinical trials regularly exceed 512 tokens. Prior to pre-training, the model was warm-started from the ROBERTA checkpoint (Liu et al., 2019) and then pre-trained on the standard BERT corpus consisting of the Books (Zhu et al., 2015), CC-News (Guu et al., 2020), Stories (Trinh and Le, 2019), and Wikipedia (Wikimedia Foundation, 2023) datasets. We evaluate two model sizes: BIGBIRD Base (125 million parameters) and BIGBIRD Large (355 million parameters).

BIOMEDBERT uses a domain-specific pretraining corpus sourced from articles on PubMed Gu et al. (2021). It uses the same model architecture as RoBERTa, but a different, domain-specific token vocabulary optimized for PubMed articles. Like BIGBIRD, we evaluate two model sizes: BIOMEDBERT Base (125 million parameters) and BIOMEDBERT Large (355 million parameters).

We fine-tune the pre-trained architectures to classify abstracts according to the inclusion and exclusion criteria enumerated in Definition 3.1. We replace the language modeling classification head with a binary classification head consisting of a single linear layer and softmax activation. We fine-tune the full model (*i.e.* not just the classification head) with cross-entropy loss. We use the Adam optimizer with learning rate 1×10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$ (determined via hyperparameter sweep). We use a maximum sequence length of 4,096 tokens. If an abstract does not fit into the maximum sequence length provided by a model, the abstract is truncated. All models were trained on a single V100 GPU with 16GB of HBM.

B.5 Performance Assessment. In this appendix, we assess the performance of various fine-tuned language models for classifying whether a publication satisfies the restrictions enumerated in Definition 3.1. Figure B.2 displays estimates of the receiver operating characteristic of several fine-tuned large language models. Here, we consider only the performance of the ROBERTA base model. The top and bottom panels are trained on noisy labels extracted from GPT-3.5 and GPT-4, respectively. We vary the quantity of noisy labels used to train the language model. Additionally, in the bottom panel, we display the performance of a fine-tuned model trained with the 1000

FIGURE B.2. Receiver Operating Characteristic by Data Size and Origin



Notes: Figure B.2 displays estimates of the the receiver operating characteristic of several fine-tuned language models. The models differ in terms of the source of their training data. Solid lines display the performance of models trained on noisy labels extracted from either GPT-3.5 or GPT-4. The dotted line displays the performance of a model trained on the hand-labeled data in the training dataset. Estimates of the true positive rate and false positive rate of GPT-3.5 and GPT-4 are indicated with black dots.

hand-labelled observations in the training dataset. The performance of each language model is measured in the testing dataset.

The fine-tuned models trained with noisy labels extracted from GPT-4 significantly outperform the models trained with noisy labels extracted from GPT-3.5 or with the hand-labels. There appears to be a threshold, where performance dramatically improves, at around 8,000 training labels. The models displayed in Figure 3 were each obtained by using all 64,000 training labels.

B.6 Final Model: Error Analysis. The Conservative model incorrectly labels 27 papers. We inspect each of these errors. We review each abstract and aim to understand what might have generated an error. In roughly half of cases—13 of 27—there is a clear error. These clear errors include publications that explicitly report the results of observational studies (one error), that re-analyze existing data (eight errors), and that report literature reviews (one error). This set also includes clinical trials that satisfy all criteria in [Definition 3.1](#), except the treatment being studied is not a drug (three errors). In 14 cases, however, errors are associated with records that are difficult for human labellers to categorize, either because the content of the publication does not fit neatly into the inclusion and exclusion criteria implied by [Definition 3.1](#) or because the publication is written in a way that makes it difficult to determine details of the study. PubMed record 27880726 is illustrative. This publication studies the effect of chrysophanic acid (CA) on benign prostatic hyperplasia. The abstract is unclear about whether this treatment is studied in human subjects or in an animal model—though the publication title and full-text make clear that this was conducted in a sample of rats. Based on the abstract alone, however, nothing suggests that this is an animal study. It is excluded from both the conservative and moderate model-generated samples, but a human-labeller flagged it as satisfying our criteria.

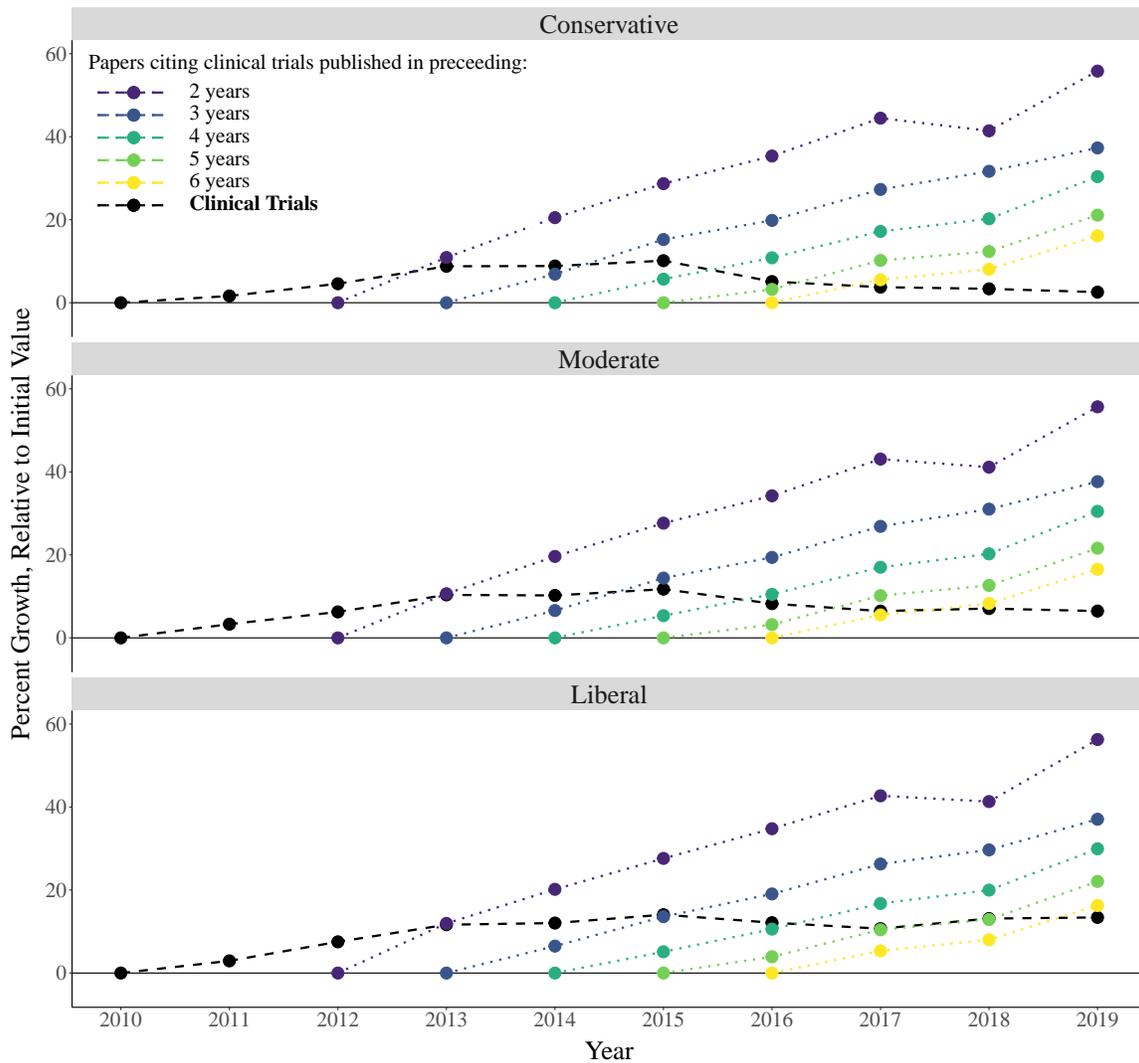
APPENDIX C. ADDITIONAL FIGURES AND FURTHER ANALYSES

This appendix collects additional figures and analyses. [Figure C.3](#) displays series analogous to the series displayed in [Figure 4](#), constructed with the Moderate and Liberal sample restrictions. The results are very similar.

One concern with the comparison made in [Figure 4](#) is that researchers may cite a growing number of publications over time. That is, more papers may cite clinical trials over time because citations are becoming, in some sense, “cheaper.” The relevant clinical literature, then, may not be any larger. To alleviate such concerns, [Figure C.4](#) reports counts, where we scale each citation assigned to a paper by the total number of citations in the originating paper. That is, a paper that receives two citations, each from a paper that cites two papers, will have a weighted citation count of 1.

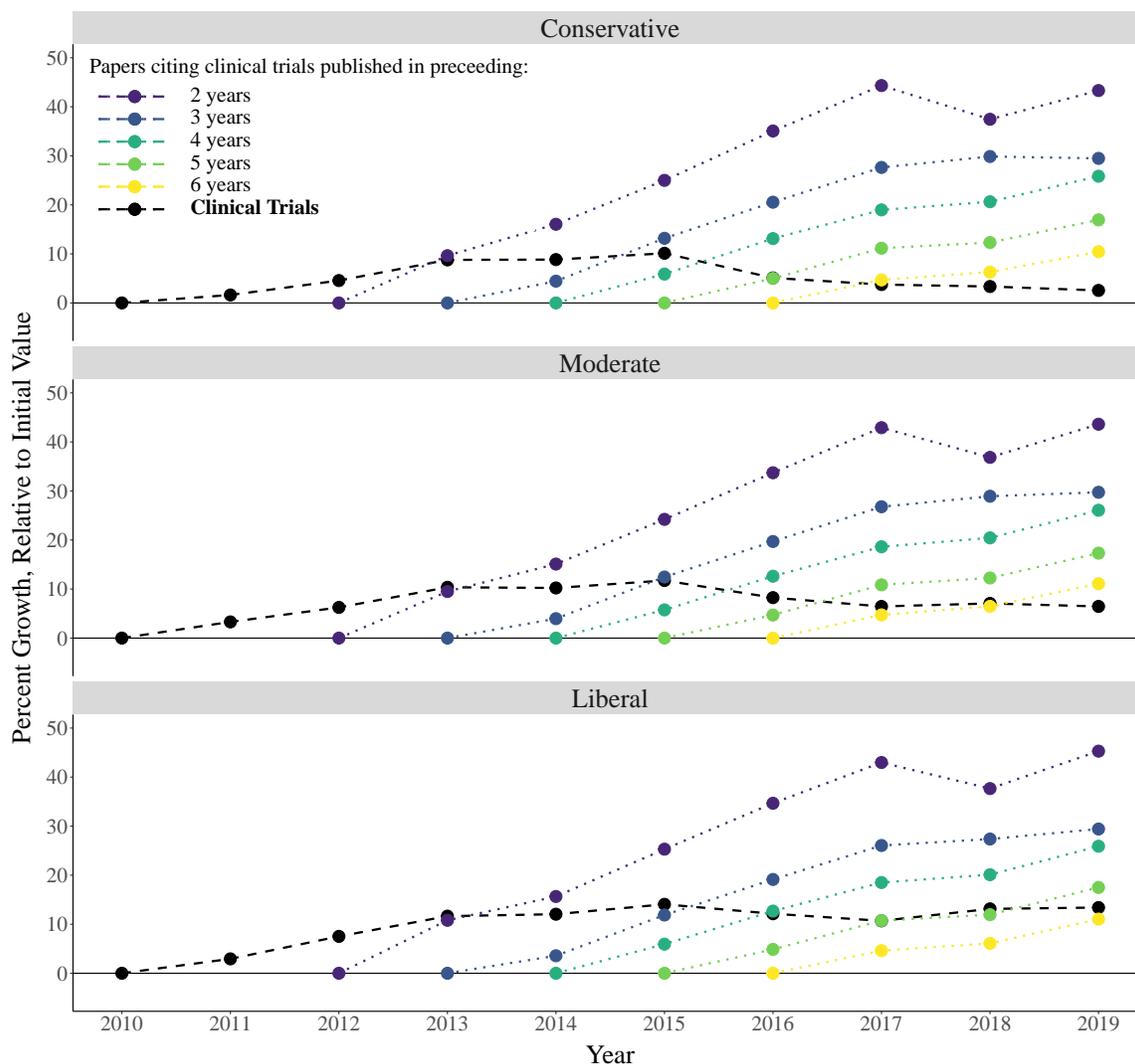
[Figure C.5](#) displays series analogous to Panel A of [Figure 5](#) with the Moderate and Liberal sample restrictions. [Figures C.6](#) and [C.7](#) display heat maps analogous to Panel B of [Figure 5](#), with the Moderate and Liberal sample restrictions and both 3-Year and 5-Year citation counts. In each case, the results are very similar across specifications.

FIGURE C.3. Growth in Clinical Research, by Sample Stringency



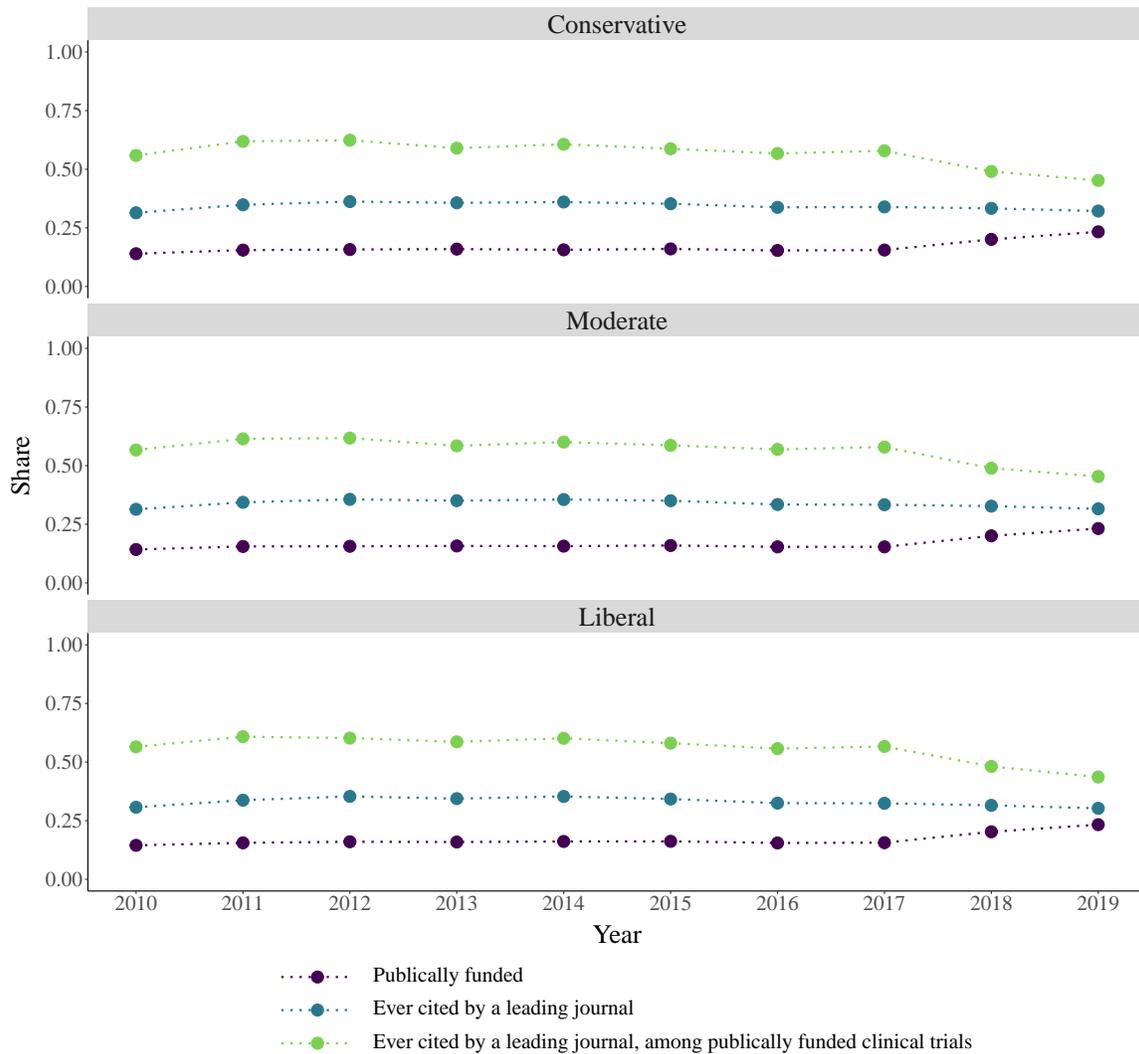
Notes: Figure C.3 displays measurements of the number of clinical trials, and papers that cite clinical trials, published in each calendar year, in all three cuts of our data—which we term “conservative,” “moderate,” and “liberal.” Each series is reported in terms of the percent change relative to its initial value. To address truncation, we report the number of publications that cite clinical trials published in the preceding t years for each t between 2 and 6.

FIGURE C.4. Growth in Cost-of-Citation Weighted Clinical Research



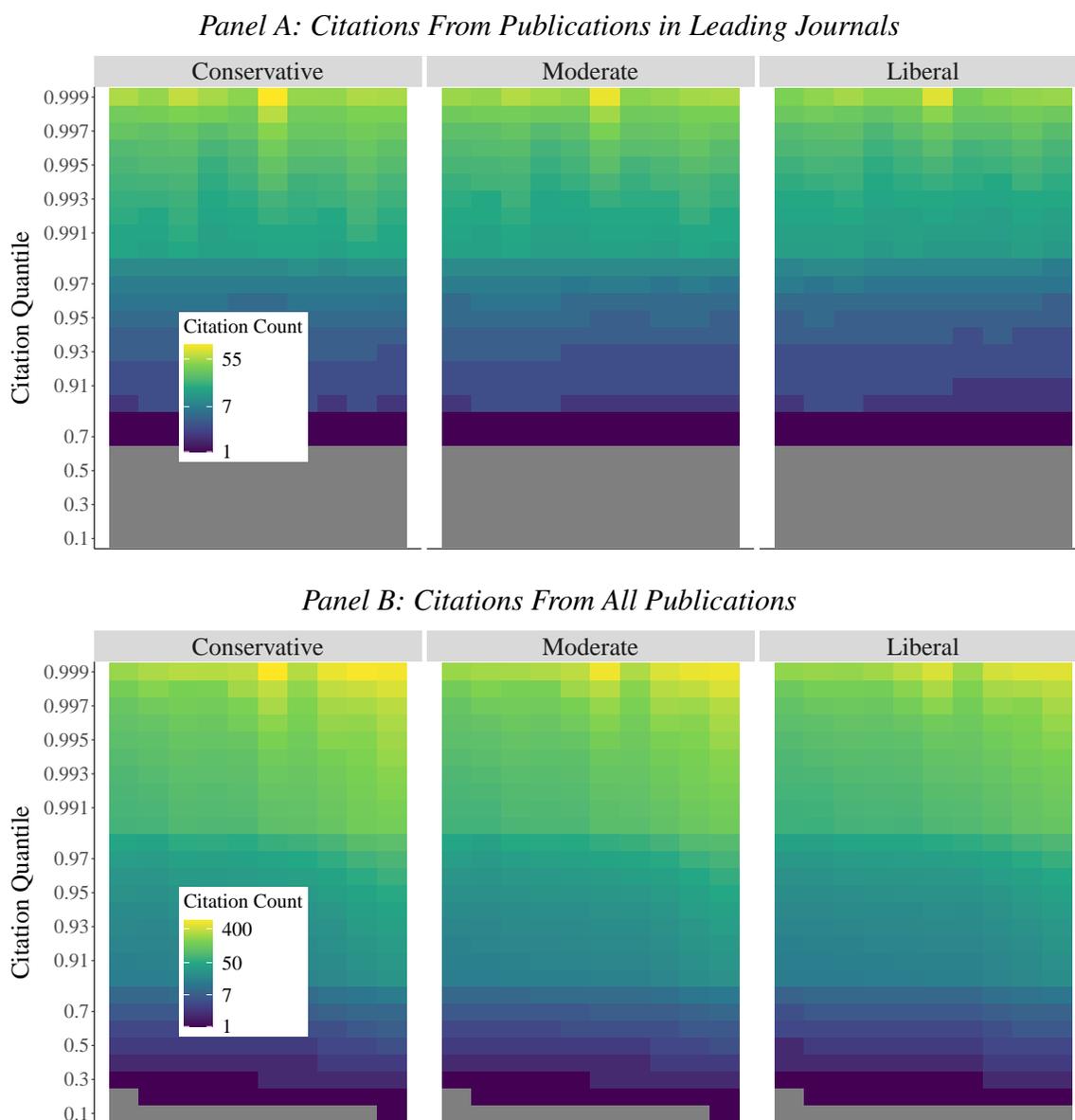
Notes: Figure C.4 displays measurements of the number of clinical trials, and papers that cite clinical trials, published in each calendar year, in all three cuts of our data—which we term “conservative,” “moderate,” and “liberal.” Here, we report the number of scientific papers using a citation-weighted metric. In particular, we scale each citation assigned to a paper by the total number of citations in the originating paper. Each series is reported in terms of the percent change relative to its initial value. To address truncation, we report the number of publications that cite clinical trials published in the preceding t years for each t between 2 and 6.

FIGURE C.5. Public Funding and Low-Quality Research, by Sample Stringency



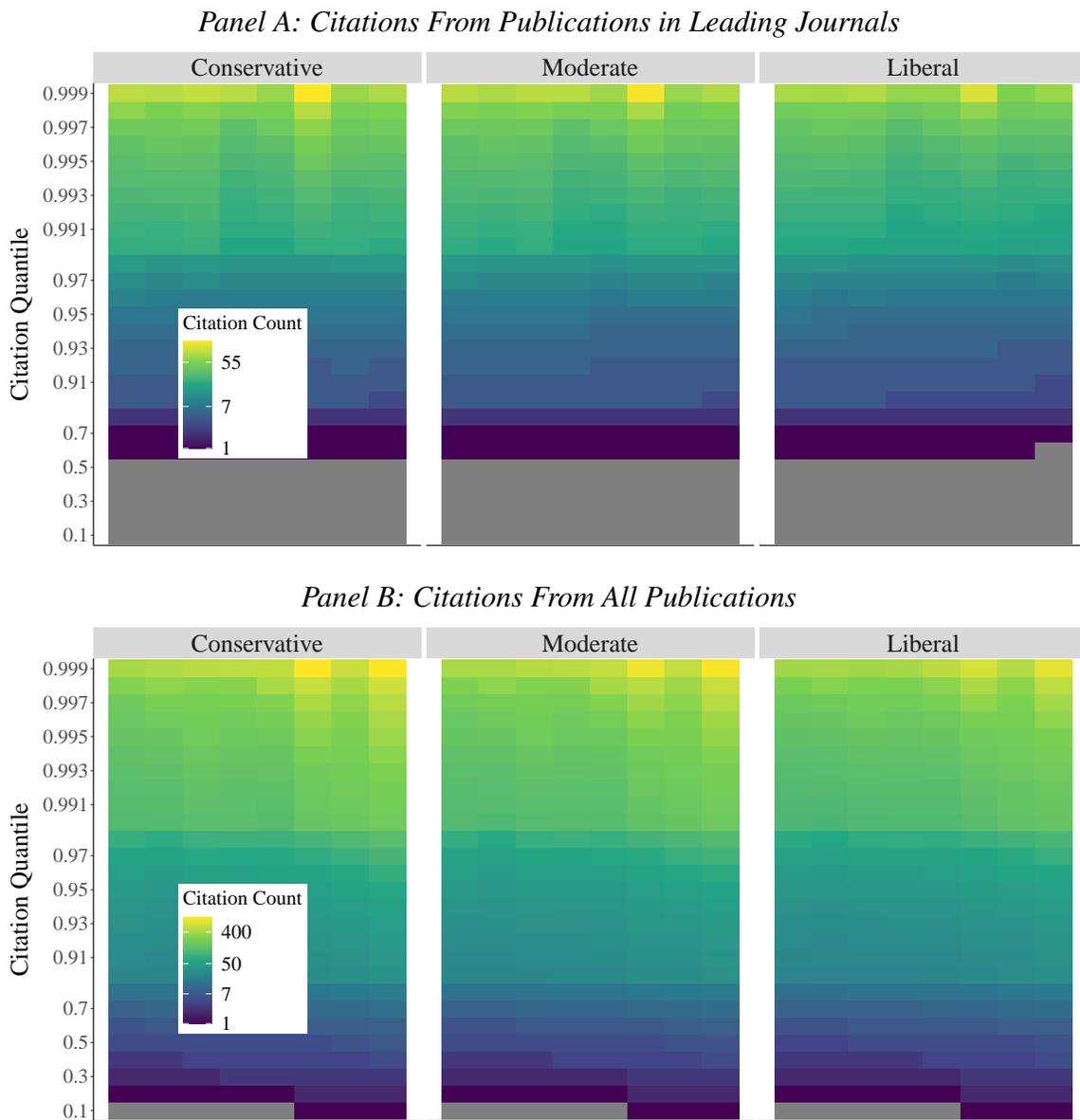
Notes: Figure C.5 illustrates the stability of the heterogeneity in the composition and quality of published clinical trials between 2010 and 2019 in all three cuts of our data—which we term “conservative,” “moderate,” and “liberal.” See Figure 3 for the true positive and false positive rates of these models. The figure displays three time series: the proportion of clinical trials that are publicly funded, the proportion of clinical trials that are ever cited by a leading journal, and the proportion of publicly funded clinical trials that are ever cited by a leading journal.

FIGURE C.6. Quality over Time, 3-Year Cites, by Sample Stringency



Notes: Figure C.6 displays heat maps measuring the distribution of three-year citations received by clinical trials in each calendar year from leading journals and from all publications, respectively. The y -axis has been stretched to elongate the right-tail of the citation distribution. Colors are displayed in a log scale. We consider the sample of clinical trials constructed with the conservative, moderate, and liberal model.

FIGURE C.7. Quality over Time, 5-Year Cites, by Sample Stringency



Notes: Figure C.7 displays heat maps measuring the distribution of five-year citations received by clinical trials in each calendar year from leading journals and from all publications, respectively. The y -axis has been stretched to elongate the right-tail of the citation distribution. Colors are displayed in a log scale. We consider the sample of clinical trials constructed with the conservative, moderate, and liberal model.

APPENDIX D. PROMPT REPOSITORY

This appendix serves as a repository for figures displaying each of the prompts considered in [Appendix B.3](#). Throughout, the dummy text “{abstract}” is placed in the location where the abstract of a publication would be input. The baseline “True/False” Prompt 1.0 is displayed in [Figure D.8](#). Refinements to this prompt are displayed in [Figures D.9 to D.11](#). The prompt that asks the model to categorize excluded publications, i.e., Prompt 2.0, is displayed in [Figure D.12](#). A refinement to this prompt is given in [Figure D.13](#). The prompt that asks the model to explain why it has excluded a publication, i.e., Prompt 3.0, is displayed in [Figure D.14](#). In turn, a refinement to this prompt is given in [Figure D.15](#).

FIGURE D.8. Prompt 1.0: True/False, Base

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE`. If the publication does not satisfy all criteria, return `FALSE`. Do not return any extraneous text. You must return either `TRUE` or `FALSE`.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.9. Prompt 1.1: True/False, Initial Examples

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE`. If the publication does not satisfy all criteria, return `FALSE`. Do not return any extraneous text. You must return either `TRUE` or `FALSE`.

For example, if the publication evaluates the effects of a surgical or medical device, procedure, or diagnostic, you should return `FALSE`. This means that if the publication studies any procedure used in the course of clinical care without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of ventricular assist device placement and thrombolysis may not satisfy the criteria.

If the publication evaluates the effects of a vitamin, supplement, food/diet program, or beverage, you should return `FALSE`. This means that if the publication studies any intervention that distributes nutritional aids and supplementary sources of vitamins without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of iron supplementation, omega-3 oil, infant formula, and diets that include healthy foods may not satisfy the criteria.

If the publication evaluates the effects of a behavioral therapy, physical therapy, exercise program, or other lifestyle intervention, you should return `FALSE`. This means that if the publication studies any intervention that aims to alter human behavior without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of smoking cessation programs, peer counseling, parenting courses, and cardiovascular fitness programs may not satisfy the criteria.

If the publication summarizes findings of other studies only, you should return `FALSE`. This means that if the publication reports findings from a meta-analysis, systematic review, literature review, literature search, or case review without also presenting novel findings on the effects of a specific, named drug, you should return `FALSE`. Studies that reference database searches in, for example, MEDLINE or EMBASE may not satisfy the criteria.

If the publication reports results from a retrospective or observational study, you should return `FALSE`. This means that if the publication reports findings that use previously collected data, or reports results from a study in which the investigators had no control over assignment to treatment or other experimental conditions, you should return `FALSE`. Studies that reference certain designs--such as case-control studies, cohort studies, or propensity score matching studies--may not satisfy the criteria. Similarly, studies that re-analyze previously collected data, often drawn from existing databases or registries, may not satisfy the criteria.

If the publication describes a clinical trial protocol, without also reporting results from the study, you should return `FALSE`. This means that if the publication describes the design, recruitment strategy, and intervention plan for a study, but does not report any findings, you should return `FALSE`. Publications written in the future tense, which describe interventions that will happen in the future, may not satisfy the criteria. Note that a protocol publication need not explicitly use the word "protocol."

If the publication evaluates the effects of a drug on animals, you should return `FALSE`. This means that if the publication describes a study that enrolled any non-human participants, you should return `FALSE`. Studies of rats, swine, primates--and studies on rat, swine, and primate models or cell lines--may not satisfy the criteria.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.10. Prompt 1.2: True/False, Extended Examples, Short

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific, named investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE`. If the publication does not satisfy all criteria, return `FALSE`. Do not return any extraneous text. You must return either `TRUE` or `FALSE`.

For example, if the publication summarizes or aggregates findings of other studies only, you should return `FALSE`. This means that if the publication reports findings from a meta-analysis, systematic review, literature review, literature search, case review, or post hoc evaluation of the findings of other trials without also presenting novel findings on the effects of a specific, named drug, you should return `FALSE`. Studies that reference database searches in, for example, MEDLINE or EMBASE may not satisfy the criteria.

If the publication reports results from a retrospective or observational study, you should return `FALSE`. This means that if the publication reports findings that use previously collected data, or reports results from a study in which the investigators had no control over assignment to treatment or other experimental conditions, you should return `FALSE`. Studies that reference certain designs--such as case-control studies, cohort studies, or propensity score matching studies--may not satisfy the criteria. Similarly, studies that re-analyze previously collected data, often drawn from existing databases or registries, may not satisfy the criteria.

If the publication describes a clinical trial protocol, without also reporting results from the study, you should return `FALSE`. This means that if the publication describes the design, recruitment strategy, and intervention plan for a study, but does not report any findings, you should return `FALSE`. Publications written in the future tense, which describe interventions that will happen in the future, may not satisfy the criteria. Note that a protocol publication need not explicitly use the word "protocol."

If the publication evaluates the effects of a drug on animals, you should return `FALSE`. This means that if the publication describes a study that enrolled any non-human participants, you should return `FALSE`. Studies of rats, swine, primates--and studies on rat, swine, and primate models or cell lines--may not satisfy the criteria.

If the publication evaluates the effects of a surgical or medical device, procedure, or diagnostic, you should return `FALSE`. This means that if the publication studies any procedure used in the course of clinical care without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of ventricular assist device placement and thrombolysis may not satisfy the criteria.

If the publication evaluates the effects of a vitamin, supplement, food/diet program, or beverage, you should return `FALSE`. This means that if the publication studies any intervention that distributes nutritional aids and supplementary sources of vitamins without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of iron supplementation, omega-3 oil, infant formula, probiotic supplements, and diets that include healthy foods may not satisfy the criteria. Note that these supplements may have different names in different settings (e.g., "ferrous" supplements or "iron" complexes).

If a publication only mentions a drug or medicine, but does not evaluate its effects, you should return `FALSE`. This means that if the publication studies an intervention in patients who happen to be taking a drug or measures the effects of a naturally-occurring substance in the human body, you should return `FALSE`.

If the publication evaluates the effects of a behavioral therapy, physical therapy, exercise program, or other lifestyle intervention, you should return `FALSE`. This means that if the publication studies any intervention that aims to alter human behavior without also studying the effects of a specific, named drug, you should return `FALSE`. Studies of smoking cessation programs, peer counseling, parenting courses, and cardiovascular fitness programs may not satisfy the criteria.

FIGURE D.10. Prompt 1.2: True/False, Extended Examples, Short (Cont.)

To determine if a drug is "specific" and "named," you should consider whether it is referenced by a unique brand, chemical, generic, or internal company name. General classes of drugs or therapies do not meet this requirement. For example, studies of "statins", "antidepressants", or "anesthesia" refer to generic categories of drugs, not specific drugs and would not meet this requirement. Similarly, "antihypertensive therapies", "antiretrovirals", and "dopaminergic therapies" are generic classes and would not meet this requirement. For this requirement to be met, you should be able to name at least one drug or medicine that is being evaluated in this study.

Please consider these guidelines very carefully before returning an answer. Before returning an answer, please identify the specific name of at least one drug or medicine being studied. Confirm that this named medicine/drug is not a food, vitamin, or supplement. Confirm, also, that this is not a "class" or "type" of medicine without a specific name.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.11. Prompt 1.3: True/False, Extended Examples, Long

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific, named investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE`. If the publication does not satisfy all criteria, return `FALSE`. Do not return any extraneous text. You must return either `TRUE` or `FALSE`.

For example, if the publication summarizes, synthesizes, or aggregates findings of other studies only, you should return `FALSE`. This means that if the publication reports findings from a meta-analysis, systematic review, literature review, literature search, case review, or post hoc evaluation of the findings of other trials without also presenting novel findings on the effects of a specific, named drug, you should return `FALSE`. For example, studies that reference database searches in, for example, MEDLINE or EMBASE may not satisfy the criteria. Articles intended to summarize the state of treatment guidelines or provide guidelines for clinicians may not satisfy the criteria.

If the publication reports results from a retrospective or observational study, you should return `FALSE`. This means that if the publication reports findings that use previously collected data, reports a post-hoc analysis, or reports results from a study in which the investigators had no control over assignment to treatment or other experimental conditions, you should return `FALSE`. For example, studies that reference certain designs--such as case-control studies, cohort studies, or propensity score matching studies--may not satisfy the criteria. Similarly, studies that re-analyze previously collected data, often drawn from existing databases or registries, may not satisfy the criteria. Finally, studies that explicitly describe their design as "retrospective" or "observational" do not satisfy the criteria.

If the publication describes a clinical trial protocol, without also reporting results from the study, you should return `FALSE`. This means that if the publication describes the design, recruitment strategy, and intervention plan for a study, but does not report any findings, you should return `FALSE`. For example, publications written in the future tense, which describe interventions that will happen in the future, may be protocols. Similarly, articles that describe results that "will be" collected are likely to be protocols. Note that a protocol publication need not explicitly use the word "protocol."

If the publication evaluates the effects of a drug on any animal subjects, you should return `FALSE`. This means that if the publication describes a study that enrolled ANY non-human participants, you should return `FALSE`. For example, studies of rats, swine, primates--and studies on rat, swine, and primate models or cell lines--may not satisfy the criteria. Any other study on animals, similarly, will not satisfy the criteria.

If the publication evaluates the effects of a surgical or medical device, procedure, or diagnostic, you should return `FALSE`. This means that if the publication studies any procedure used in the course of clinical care without also studying the effects of a specific, named drug, you should return `FALSE`. For example, studies of ventricular assist device placement and thrombolysis may not satisfy the criteria.

If the publication evaluates the effects of a vitamin, supplement, food/diet program, or beverage, you should return `FALSE`. This means that if the publication studies any intervention that distributes nutritional aids and supplementary sources of vitamins without also studying the effects of a specific, named drug, you should return `FALSE`. For example, studies of iron supplementation, omega-3 oil, infant formula, probiotic supplements, and diets that include healthy foods may not satisfy the criteria. Note that these supplements may have different names in different settings (e.g., "ferrous" supplements or "iron" complexes).

If the publication only reports the results of a scientific study conducted in a laboratory setting (i.e., not on live human subjects), you should return `FALSE`. This means that if the publication studies a setting with no human subjects--including a laboratory investigation of cell lines or biological materials drawn from human subjects--you must return `FALSE`.

FIGURE D.11. Prompt 1.3: True/False, Extended Examples, Long (Cont.)

If a publication only mentions a drug or medicine, but does not evaluate its effects, you should return `FALSE`. This means that if the publication studies an intervention in patients who happen to be taking a drug or measures the effects of a naturally-occurring substance in the human body, you should return `FALSE`.

If the publication evaluates the effects of a behavioral therapy, physical therapy, exercise program, or other lifestyle intervention, you should return `FALSE`. This means that if the publication studies any intervention that aims to alter human behavior without also studying the effects of a specific, named drug, you should return `FALSE`. For example, studies of smoking cessation programs, peer counseling, parenting courses, and cardiovascular fitness programs may not satisfy the criteria.

To determine if a drug is "specific" and "named," you should consider whether it is referenced by a unique brand, chemical, generic, or internal company name. General classes of drugs or therapies do not meet this requirement. For example, studies of "statins", "antidepressants", or "anesthesia" refer to generic categories of drugs, not specific drugs and would not meet this requirement. Similarly, "antihypertensive therapies", "antiretrovirals", and "dopaminergic therapies" are generic classes and would not meet this requirement. For this requirement to be met, you should be able to name at least one drug or medicine that is being evaluated in this study. Note, however, that simply because a medicine name is mentioned does not mean that these criteria are satisfied. The remainder of the definition must still be satisfied.

Please consider these guidelines very carefully before returning an answer. Before returning an answer, please identify the specific name of at least one drug or medicine being studied. Confirm that this named medicine/drug is not a food, vitamin, or supplement. Confirm, also, that this is not a "class" or "type" of medicine without a specific name.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.12. Prompt 2.0: Categorize Exclusion Restriction

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE`. If the publication does not satisfy all criteria, you must provide at least one reason for your choice. To do so, you must follow the following directions.

Reasons:

- Return "NO DRUG" if the publication does not report a clinical trial that studies at least one medicine/drug dispensed in the course of medical care. Drugs MAY include the following: investigational compounds, biological therapies (including vaccines), and small-molecule, synthetic products. Drugs DO NOT include the following: dietary supplements, diets, behavioral interventions, medical and surgical procedures, medical and surgical devices, diagnostic tests, and dental substances. Drugs also do not include chemicals, substances, and toxins that are not dispensed in the course of medical care.
- Return "ANIMAL" if the publication reports the results of a study conducted, in part or wholly, on animals.
- Return "NO HUMAN SUBJECTS" if the publication only reports the results of a scientific study conducted in a laboratory setting (i.e., not on live human subjects). Studies with no human subjects include laboratory investigations of cell lines or biological materials drawn from human subjects if no living human subjects are actively involved.
- Return "META-ANALYSIS" if the publication only reports the results of a meta-analysis, literature review, literature search, narrative review, or summary of individuals' careers. Summaries of previously published clinical trials--without new data collection or analysis--should be excluded for this reason.
- Return "OBSERVATIONAL" if the publication only reports results from an observational study (i.e., there was no active intervention). Observational studies may include studies using data drawn from existing datasets, registries, and electronic health records.
- Return "RETROSPECTIVE" if the publication only reports results from a retrospective study. A retrospective study is a non-interventional study that compares two groups of people: those with the disease or condition under study (cases) and a very similar group of people who do not have the disease or condition (controls). The allocation of people to the cases or controls is not chosen by the researcher, and data may be drawn from a dataset not explicitly intended for the research being reported.
- Return "PROTOCOL" if the publication only reports the protocol for a clinical trial. A protocol is a formal analysis plan for a clinical trial that describes the design and attributes of the study, but does not include clinical results (e.g., no statistical results or analyses are reported).
- Return "CASE REPORT" if the publication only reports information collected from clinical cases. A case report is a detailed report of the symptoms, signs, diagnosis, treatment, and follow-up of an individual patient or set of patients.
- Return "NOT PUBLISHED IN ENGLISH" if the publication is published in a language other than English.

If the publication does not satisfy the criteria for multiple reasons, return a complete list of reasons separated by semicolons, e.g., "OBSERVATIONAL; NO DRUG". If the publication satisfies all classification criteria, only return `TRUE`. Do not return any extraneous text. You must return either `TRUE` or a complete list of reasons described above.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.13. Prompt 2.1: Categorize Exclusion Restriction, Examples

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE'. If the publication does not satisfy all criteria, you must provide at least one reason for your choice. To do so, you must follow the following directions.

Reasons:

- Return "NO DRUG" if the publication does not report a clinical trial that studies at least one medicine/drug dispensed in the course of medical care. Drugs MAY include the following: investigational compounds, biological therapies (including vaccines), and small-molecule, synthetic products. Drugs DO NOT include the following: dietary supplements, diets, behavioral interventions, medical and surgical procedures, medical and surgical devices, diagnostic tests, and dental substances. Drugs also do not include chemicals, substances, and toxins that are not dispensed in the course of medical care.
- Return "ANIMAL" if the publication reports the results of a study conducted, in part or wholly, on animals.
- Return "NO HUMAN SUBJECTS" if the publication only reports the results of a scientific study conducted in a laboratory setting (i.e., not on live human subjects). Studies with no human subjects include laboratory investigations of cell lines or biological materials drawn from human subjects if no living human subjects are actively involved.
- Return "META-ANALYSIS" if the publication only reports the results of a meta-analysis, literature review, literature search, narrative review, or summary of individuals' careers. Summaries of previously published clinical trials--without new data collection or analysis--should be excluded for this reason.
- Return "OBSERVATIONAL" if the publication only reports results from an observational study (i.e., there was no active intervention). Observational studies may include studies using data drawn from existing datasets, registries, and electronic health records.
- Return "RETROSPECTIVE" if the publication only reports results from a retrospective study. A retrospective study is a non-interventional study that compares two groups of people: those with the disease or condition under study (cases) and a very similar group of people who do not have the disease or condition (controls). The allocation of people to the cases or controls is not chosen by the researcher, and data may be drawn from a dataset not explicitly intended for the research being reported.
- Return "PROTOCOL" if the publication only reports the protocol for a clinical trial. A protocol is a formal analysis plan for a clinical trial that describes the design and attributes of the study, but does not include clinical results (e.g., no statistical results or analyses are reported).
- Return "CASE REPORT" if the publication only reports information collected from clinical cases. A case report is a detailed report of the symptoms, signs, diagnosis, treatment, and follow-up of an individual patient or set of patients.
- Return "NOT PUBLISHED IN ENGLISH" if the publication is published in a language other than English.

If the publication does not satisfy the criteria for multiple reasons, return a complete list of reasons separated by semicolons, e.g., "OBSERVATIONAL; NO DRUG". If the publication satisfies all classification criteria, only return `TRUE'. Do not return any extraneous text. You must return either `TRUE' or a complete list of reasons described above.

For example, if the publication evaluates the effects of a surgical or medical device, procedure, or diagnostic, you should return `NO DRUG'. This includes publications that study any procedure used in the course of clinical care without also studying the effects of a specific, named drug. This may include, for example, studies of ventricular assist device placement and thrombolysis.

FIGURE D.13. Prompt 2.1: Categorize Exclusion Restriction, Examples (Cont.)

If the publication evaluates the effects of a vitamin, supplement, food/diet program, or beverage, you should also return `NO DRUG'. This includes publications that study any intervention that distributes nutritional aids and supplementary sources of vitamins without also studying the effects of a specific, named drug. This may include, for example, studies of iron supplementation, omega-3 oil, infant formula, and diets that include healthy foods.

If the publication evaluates the effects of a behavioral therapy, physical therapy, exercise program, or other lifestyle intervention, you should also return `NO DRUG'. This includes publications that study any intervention that aims to alter human behavior without also studying the effects of a specific, named drug. This may include, for example, studies of smoking cessation programs, peer counseling, parenting courses, and cardiovascular fitness programs.

If the publication summarizes findings of other studies only, you should return `META-ANALYSIS'. This includes publications that report findings from a meta-analysis, systematic review, literature review, literature search, or case review without also presenting novel findings on the effects of a specific, named drug. This may include, for example, studies that reference database searches in, for example, MEDLINE or EMBASE.

If the publication reports results from a retrospective study, you should return `RETROSPECTIVE'. This includes publications that report findings that use previously collected data. This may include, for example, studies that re-analyze previously collected data, often drawn from existing databases or registries.

If the publication reports results from an observational study, you should return `OBSERVATIONAL'. This includes publications that report findings that use previously collected data, or report results from a study in which the investigators had no control over assignment to treatment or other experimental conditions. This may include, for example, studies that reference certain designs--such as case-control studies, cohort studies, or propensity score matching studies.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.14. Prompt 3.0: Provide Reason for Exclusion

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE.` If the publication does not satisfy all criteria, you must provide a reason for your choice. In particular, you must return `FALSE` followed by an explanation for why it does not meet the criteria. Your response should take the form: `FALSE: [xxx].` Do not return any extraneous text. You must return either `TRUE` or `FALSE: `, followed by your explanation.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer:

FIGURE D.15. Prompt 3.1: Provide Reason for Exclusion, Examples

At the end of this prompt, you will be shown an abstract from an academic publication indexed in the PubMed/MEDLINE database.

Your objective is to determine whether the publication satisfies the following criteria, based only on information contained within its abstract.

Criteria:

The publication reports the results of a prospective clinical trial. The clinical trial may be of any phase. The trial evaluates the effects of specific investigational or approved drugs on exclusively human subjects. The abstract is written in English.

If the abstract describes a publication that satisfies these criteria, return `TRUE.` If the publication does not satisfy all criteria, you must provide a reason for your choice. In particular, you must return `FALSE` followed by an explanation for why it does not meet the criteria. Your response should take the form: `FALSE: [xxx].` Do not return any extraneous text. You must return either `TRUE` or `FALSE: `, followed by your explanation.

For example, if the publication evaluates the effects of a surgical or medical device, procedure, or diagnostic, you should return `FALSE: ` followed by an explanation that identifies the device, procedure, or diagnostic being studied and specifies that it is not a drug. This includes publications that study any procedure used in the course of clinical care without also studying the effects of a specific, named drug. This may include, for example, studies of ventricular assist device placement and thrombolysis.

If the publication evaluates the effects of a vitamin, supplement, food/diet program, or beverage, you should return `FALSE: ` followed by an explanation that identifies the vitamin, supplement, food/diet program, or beverage being studied and specifies that it is not a drug. This includes publications that study any intervention that distributes nutritional aids and supplementary sources of vitamins without also studying the effects of a specific, named drug. This may include, for example, studies of iron supplementation, omega-3 oil, infant formula, and diets that include healthy foods.

If the publication evaluates the effects of a behavioral therapy, physical therapy, exercise program, or other lifestyle intervention, you should return `FALSE: ` followed by an explanation that identifies the behavioral therapy, physical therapy, exercise program, or lifestyle intervention being studied and specifies that it is not a drug. This includes publications that study any intervention that aims to alter human behavior without also studying the effects of a specific, named drug. This may include, for example, studies of smoking cessation programs, peer counseling, parenting courses, and cardiovascular fitness programs.

If the publication summarizes findings of other studies only, you should return `FALSE: ` followed by an explanation with two components. First, specify that the publication does not report the results of a specific prospective clinical trial and, second, provide a brief summary of what the publication does report. This includes publications that report findings from a meta-analysis, systematic review, literature review, literature search, or case review without also presenting novel findings on the effects of a specific, named drug. This may include, for example, studies that reference database searches in, for example, MEDLINE or EMBASE.

If the publication reports results from a retrospective or observational study, you should return `FALSE: ` followed by an explanation with two components. First, specify that the publication does not report the results of a prospective, interventional clinical trial and, second, provide a brief summary of what the publication does report. This includes publications that report findings that use previously collected data, or report results from a study in which the investigators had no control over assignment to treatment or other experimental conditions. This may include, for example, studies that reference certain designs--such as case-control studies, cohort studies, or propensity score matching studies. Similarly, it may include studies that re-analyze previously collected data, often drawn from existing databases or registries.

If the publication describes a clinical trial protocol, without also reporting results from the study, you should return `FALSE: ` followed by an explanation that the publication does not report the results of a clinical trial. This includes publications that describe the design, recruitment strategy, and intervention plan for a study, but do not report any findings. This may include publications written in the future tense, which describe interventions that will happen in the future. Note that a protocol publication need not explicitly use the word "protocol."

FIGURE D.15. Prompt 3.1: Provide Reason for Exclusion, Examples (Cont.)

If the publication evaluates the effects of a drug on animals, you should return `FALSE: ' followed by an explanation that describes that the study was conducted on animals and identifies the animals being studied. This includes publications that describe studies that enrolled any non-human participants. This may include studies of rats, swine, primates--and studies on rat, swine, and primate models or cell lines.

The abstract that you will consider is as follows:

Abstract: {abstract}

Answer: