

CHATGPT

Experimental evidence on the productivity effects of generative artificial intelligence

Shakked Noy* and Whitney Zhang

We examined the productivity effects of a generative artificial intelligence (AI) technology, the assistive chatbot ChatGPT, in the context of midlevel professional writing tasks. In a preregistered online experiment, we assigned occupation-specific, incentivized writing tasks to 453 college-educated professionals and randomly exposed half of them to ChatGPT. Our results show that ChatGPT substantially raised productivity: The average time taken decreased by 40% and output quality rose by 18%. Inequality between workers decreased, and concern and excitement about AI temporarily rose. Workers exposed to ChatGPT during the experiment were 2 times as likely to report using it in their real job 2 weeks after the experiment and 1.6 times as likely 2 months after the experiment.

Recent advances in generative artificial intelligence (AI) may have widespread implications for production and labor markets. Generative AI systems such as ChatGPT or DALL-E, which can be prompted to create new text or visual outputs from large amounts of training data, are qualitatively unlike most historical examples of automation technologies. Previous waves of automation predominantly affected “routine” tasks consisting of explicit sequences of steps that could be easily codified and programmed into a machine or computer, such as assembly-line manufacturing tasks or bookkeeping tasks (1, 2). By contrast, creative, difficult-to-codify tasks such as writing and image generation avoided automation, a pattern scholars have noted might change with the advent of deep learning, which now underpins generative AI systems.

The emergence of powerful generative AI technologies reintroduces a host of classic questions in a new context (3–5). Automation technologies by definition perform specific tasks in place of humans. But more broadly, these technologies may either displace humans completely from certain occupations or augment existing human workers by increasing their productivity (6–9). If automation technologies such as industrial robots mostly displace human workers, they can increase unemployment. Moreover, their impacts on aggregate productivity may be small or nonexistent to the degree that they mainly serve to redistribute income previously earned by displaced workers to the capital owners supplying their robot replacements (10). If automation technologies such as computers augment existing workers, they can simultaneously benefit workers, capital owners, and consumers by

raising wages, boosting productivity, and lowering prices (11–13).

A potent generative writing tool such as ChatGPT could conceivably either displace or augment human labor. ChatGPT could entirely replace certain kinds of writers, such as grant writers or marketers, by letting companies directly automate the creation of grant applications and press releases with minimal human oversight. Alternatively, instead of displacing workers, ChatGPT could substantially raise the productivity of grant writers and marketers, for example, by automating relatively routine, time-consuming subcomponents of their writing tasks, such as translating ideas into a rough draft. In this case, these services would become cheaper and demand could expand, resulting in higher employment and greater productivity for companies, cheaper products for consumers, and potentially higher wages for workers (14). Furthermore, inequalities between workers could either decrease if lower-ability workers are supported more by ChatGPT or increase if higher-ability workers have the skills necessary to take advantage of the new technology.

Which of these eventualities will generative AI systems bring about? The answer depends on a host of research questions (RQs). RQ1: How does access to generative AI systems affect workers’ productivity in existing tasks? Do workers choose to use these systems? Conditional on using these systems, how do workers interact with them and how do they affect productivity (15–18)? RQ2: Do these systems differentially affect low- and high-ability workers? RQ3: How do workers subjectively react to these technologies (19)?

Methods

This paper took the first step toward answering these questions (20). In a preregistered online experiment, we recruited 453 experienced, college-educated professionals on the survey platform Prolific and assigned each to

complete two occupation-specific, incentivized writing tasks (21). The experiment took place from 27 January to 21 February 2023 and involved GPT-3.5. The occupations that we drew from were marketers, grant writers, consultants, data analysts, human resource professionals, and managers. The tasks, which included writing press releases, short reports, analysis plans, and delicate emails, comprised 20- to 30-min assignments designed to resemble real tasks performed in these occupations. Indeed, most of our participants reported completing similar tasks before and rated the assigned tasks as realistic representations of their everyday work (see the supplementary materials).

Participants faced high-powered incentives in the form of large bonus payments to produce high-quality work. They received a base payment of \$10 plus up to \$14 in bonus payments for output quality, with the average overall rate of \$17/hour substantially exceeding the Prolific standard of \$12/hour. We cross-randomized the structure of bonus payments faced by participants to show the robustness of our results to different incentive schemes (see below for more details). Output quality was assessed by blinded experienced professionals working in the same occupations. Evaluators were asked to treat the output as if encountered in a work setting and were incentivized to grade outputs carefully on a scale of 1 to 7 (22). Each piece of output was seen by three evaluators, with an average within-essay cross-evaluator correlation of 0.44 (23).

We randomly assigned 50% of participants to the treatment group and 50% to the control group. The treatment group was instructed to register for ChatGPT between the first and second task, received guidance on using it, and were told they were permitted to use it on the second task if they found it useful. The control group was instead instructed to register for the LaTeX editor Overleaf in an attempt to hold the time and hassle costs of signup constant between the two groups. The control group was not told they could use Overleaf on the second task and <5% of participants subsequently reported using it.

In addition to output quality evaluations, we collected self-reported and objective measures of participants’ time spent on the tasks and took snapshots of participants’ outputs each minute while they performed the task to construct objective measures of activity and to detect ChatGPT usage (see the supplementary materials).

A complete description of our experimental design, copies of relevant survey questionnaires, and additional figures validating our central measures and extending our main results are included in the supplementary materials. Descriptive statistics about the sample, as well as balance and selective attrition tests, are

Department of Economics, Massachusetts Institute of Technology, Cambridge, MA, USA.

*Corresponding author. Email: snoy@mit.edu

available in Table 1. The attrition rate was 6% in the control group and 11% in the treatment group. Balance tests indicate that across 13 pretreatment characteristics, the treatment and control groups exhibited a small but significant difference only for only two characteristics: employment status and being a human resources professional. Our partly within-person design, which controls for performance on the pretreatment task, should eliminate any influence of selective attrition on our results. In the supplementary materials, we also report Lee bounds (24) on our main results and versions of our results controlling for employment status and occupation, which confirmed that our results are highly robust to selective attrition.

Results

Take-up of ChatGPT

In the treatment group, 92% of treated participants successfully registered for ChatGPT, and 80% chose to use it on the second task (25). Users gave it an average self-assessed usefulness score of 4.4 out of 5.

Before treatment, 70% of our participants had heard of ChatGPT and 32% had used it before. Self-reported and objective measures indicated that ~10 to 20% of the control group used ChatGPT on the tasks (see the supplementary materials), meaning there was at least a 60-percentage point experimentally induced gap in usage between our treatment and control groups on the second task. Our estimates reflect the effects of ChatGPT on the average productivity of the 60 to 70% of participants whose usage was determined by their treatment assignment, and constitute lower bounds on the effects of ChatGPT usage on productivity. In the supplementary materials, we report two-stage least-squares results

adjusting our estimates upward for imperfect compliance.

Productivity

We first show results for our two productivity measures: time taken and evaluator grades (Fig. 1). The experimental intervention shifted both outcomes substantially. In the treatment group, time taken on the posttreatment task dropped by 11 min (0.75 SDs) relative to the control group, who took an average of 27 min ($P < 0.001$). Average evaluator grades in the treatment group increased by 0.45 SDs ($P < 0.001$), with similar increases for overall grades and specific grades for writing quality, content quality, and originality.

These effects are not limited to specific pockets of the time or grade distributions. As shown in Fig. 1, C and D, the entire time distribution shifted to the left (faster work) and the entire grade distribution shifted to the right (higher quality). At the individual worker level, as shown in Fig. 2, treated workers who received a low grade on the first task experienced both 1- to 2-point increases in grades and 10-min decreases in time spent, whereas workers who received a high grade maintained their grade level while also reducing their time spent by ~10 min.

These results are virtually identical across our two main incentive schemes, which covered 80% of respondents: a “linear” scheme in which respondents were paid \$1 for each point they received on each submission (each of which was graded on a 1- to 7-point scale), and a “convex” scheme in which respondents were additionally paid \$3 for earning a grade of 6 or 7. The results shown in Fig. 1 are based on these two incentive schemes. The fact that treated participants reduced their time spent by a similar amount even when faced

with strong incentives to produce high-quality output (under the convex scheme) demonstrates that the time-saving effects of ChatGPT are not specific to linear payment regimes and apply robustly across incentive structures.

In our third incentive arm involving 20% of participants, we required participants to spend exactly 15 min on each task, thereby holding effort fixed across the treatment and control groups and allowing us to interpret any difference in grades as a pure effect of ChatGPT access on productive capacity. In this arm, the treatment increased grades by a similar albeit not statistically significant 0.33 SDs (26).

In an additional intervention, after completing the second task, 30% of the treatment group were shown their first-task human-created output and given the opportunity to edit or replace it using ChatGPT. Of these participants, 19% chose to replace their response with ChatGPT’s output and another 17% used ChatGPT to edit their original response, suggesting that participants viewed ChatGPT as a means to improve output quality as well as save time.

Productivity inequality

The control group exhibited persistent productivity inequality: Participants who scored well on the first task also tended to score well on the second task. As Fig. 2A shows, there was a correlation of 0.41 ($P < 0.001$) between a control participant’s grade on the first task and their grade on the second task, holding the evaluator constant.

In the treatment group, initial inequalities were more than half-erased by the treatment: The correlation between first-task and second-task grades was only 0.14 ($P < 0.001$ for the difference in slopes). This reduction

Table 1. Descriptive statistics.

Variable	N (control)	Mean (control)	N (treatment)	Mean (treatment)	Difference
Annual salary in main job (\$)	234	67,764	213	71,938	4173
Years of tenure in occupation	234	10.49	215	10.07	-0.43
Employed	226	91%	210	96%	5.0%**
Occupation: HR professional	235	6%	218	11%	4.6%*
Occupation: business consultant	235	13%	218	11%	-1.3%
Occupation: data analyst	235	11%	218	11%	-0.0%
Occupation: grant writer	235	16%	218	17%	1.2%
Occupation: manager	235	43%	218	41%	-1.7%
Occupation: marketer	235	11%	218	9%	-2.3%
Time spent (task 1, min)	227	26.10	212	26.58	0.47
Average grade (task 1)	233	3.63	211	3.77	0.15
Job satisfaction (task 1, scale of 1 to 10)	234	6.30	215	6.34	0.04
Self-efficacy (task 1, scale of 1 to 10)	234	6.89	215	6.90	0.01

This table presents descriptive statistics for our sample. We recode salary reports of >\$500,000 to missing (affects two observations). “Employed” includes full-time and part-time employment. * $P < 0.10$; ** $P < 0.05$.

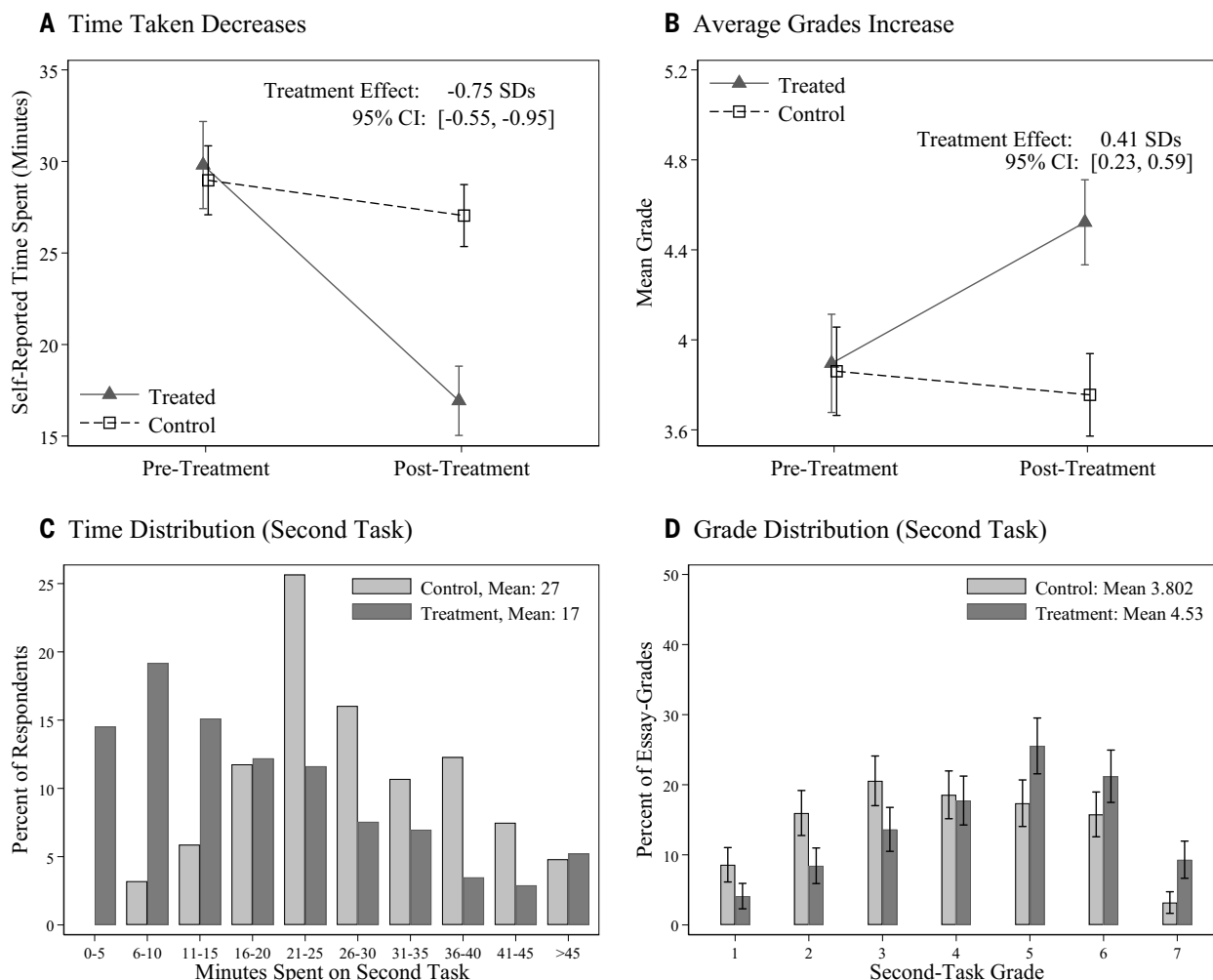


Fig. 1. Treatment effects on productivity. Effects shown are restricted to the linear and convex incentive groups. **(A)** and **(B)** Means (and 95% confidence intervals for those means) of self-reported time taken and average grades in the first and second task separately in the treatment and control groups. The results look very similar for the objective measure of time active (see the supplementary materials). Also shown are treatment effect coefficients and 95% confidence intervals, rescaled to be in terms of pretreatment SDs of the outcome variable. The coefficients are estimated from regressions

of the within-participant change in outcome from before to after treatment on a treatment dummy, occupation*task order fixed effects, and incentive arm fixed effects. In **(A)**, this is at the participant level and SEs are heteroskedasticity robust. In **(B)**, this is at the participant-evaluator level, the regression also includes grader fixed effects, and SEs are clustered at the participant level. **(C)** and **(D)** Raw graphs of the outcome distribution in the treatment versus control group on the second task; **(C)** is at the participant level and **(D)** is at the participant-evaluator level.

in inequality was driven by the fact that participants who scored lower on the first task benefited more from ChatGPT access. As Fig. 2A shows, the gap between treatment and control is much larger at the left end of the x axis.

Human-machine interactions

What kinds of human-machine interactions underlie the productivity results documented above? Did workers paste the task prompt into ChatGPT and immediately submit its output, minimizing their time spent and increasing their grades because ChatGPT's writing abilities exceeded theirs? Or did they treat ChatGPT as a helpful but imperfect tool, for example, using it to create a rough draft and then spending

time editing and improving the draft or using it to brainstorm or edit?

Our evidence supports the first possibility. Almost everyone submitted lightly edited or unedited ChatGPT output, and we observed small time expenditures on editing and no resulting improvement in respondents' grades. In the treatment group, 33% of participants reported submitting ChatGPT's initial output without editing it, and 53% reported editing before submitting. However, those who reported editing were active on the task for only 3.3 min on average after we first observed them pasting in a large quantity of text (presumably from ChatGPT), with most active for 0 to 2 min (27). Qualitative examination suggests that

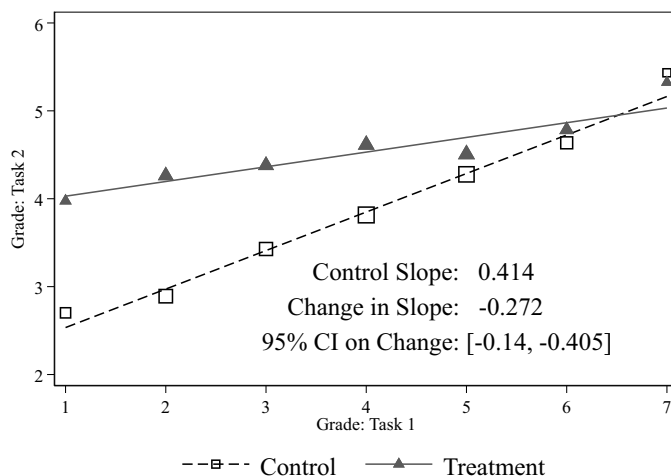
most of this editing was superficial, such as changing a placeholder or rearranging a sentence. Evaluator grades also suggest that this editing was ineffectual. There was no correlation between how long a participant was active after pasting in the ChatGPT text and the grade they ultimately received, and treated respondents who used ChatGPT did not receive higher average grades than the raw ChatGPT output that we gave to evaluators to grade (see the supplementary materials).

It is not obvious whether these dynamics should be interpreted as evidence that ChatGPT will displace human workers or evidence that it will augment them. Although ChatGPT directly substituted for participants' effort with little need for human input, it also enabled

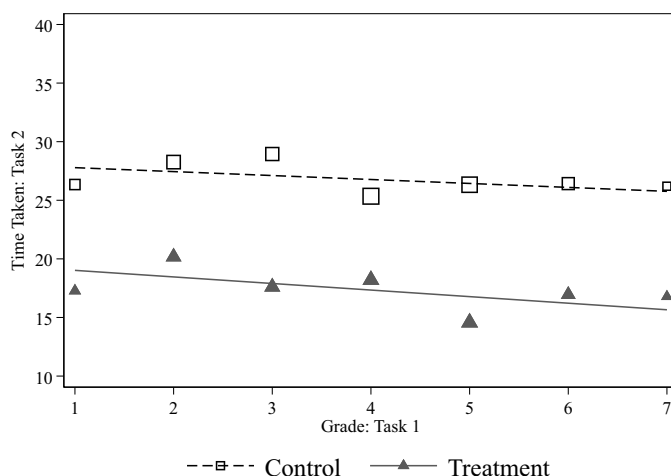
Fig. 2. Effects on grades and time across the initial grade distribution.

Participant-evaluator observations are binned together according to the task 1 grade given to this participant by this evaluator. (A and B) Within each bin, the panels plot the average task 2 grade (A) or task 2 time taken (B) for the observations in the bin separately by treatment versus control. Also shown are the control-group slope, control-treatment difference in slopes, and the 95% confidence interval for the difference. These latter results were calculated from a participant-evaluator level regression of the outcome variable on the task 1 grade, treatment status, treatment*task 1 grade, and grader fixed effects, clustering SEs at the participant level. The control group slope is the coefficient on the task 1 grade, and the difference in slopes is the coefficient on the treatment*task 1 grade interaction. Note that this difference in slopes will not match up exactly with the difference between the raw slopes plotted in the graph because these raw slopes do not use grader fixed effects.

A Grade Inequality Decreases



B Time Spent Drops Across the Initial Grade Distribution



participants to complete tasks much faster. We reflect on this further in the Discussion.

Subjective outcomes: Job satisfaction, self-efficacy, and beliefs about automation

Many of our treated participants had never heard of (30%) or never used (68%) ChatGPT before participating in the experiment. We used a battery of questions to assess their subjective reactions to encountering the technology. As depicted in Fig. 3, participants enjoyed the tasks more by 0.47 SDs when given access to ChatGPT ($P < 0.001$). Treated participants' concern for ($P < 0.01$) and excitement about ($P < 0.001$) future effects of AI on their occupations rose, and their overall optimism increased by 0.2 SDs ($P < 0.05$). These effects disappeared in the 2-week and 2-month follow-up surveys, indicating that they are best interpreted as short-run phenomena reflecting

respondents' first experiences with the technology (28).

Two-week and 2-month follow-up surveys

One powerful indication of the value of ChatGPT to participants is whether they continued to use it after the experiment, in their actual jobs. To track this, we resurveyed participants 2 weeks and 2 months after their completion of the initial survey, with response rates of 92% and 83%, respectively, and no treatment-control imbalance in response rates.

In the 2-week follow-up, 34% of former treatment group participants reported using ChatGPT in their job in the past week, compared with 18% of control group participants ($P < 0.001$). This large gap in usage fully persisted into the 2-month follow-up, when 42% of the treatment group and 27% of control respondents reported using ChatGPT in their jobs in the

past week ($P < 0.01$). The persistence of this gap suggests that the dissemination of ChatGPT into real professional activity is still in very early stages, with usage held back by a lack of knowledge about or experience with the technology.

In the 2-week follow-up, ChatGPT users gave the technology an average usefulness score of 3.66 out of 5.00, somewhat lower than in our main experiment, likely because of the greater length and complexity of real-world tasks. The participants reported using it for a broad range of tasks such as generating recommendation letters for employees, responding to customer service requests, brainstorming, rough-drafting emails, and editing.

Nonusers were divided into three roughly equal-sized groups, reporting that: (i) ChatGPT was not useful in their job, (ii) they did not know about it or did not have an account, or (iii) it was not allowed in their workplace or was usually unavailable during the day. The one-third of nonusers who claimed that it was not useful in their job mostly said that this was because the chatbot lacks context-specific knowledge that forms an important part of their writing. For example, they reported that their writing was "very specifically tailored to [their] customers and involves real time information" or "unique [and] specific to [their] company products."

Discussion

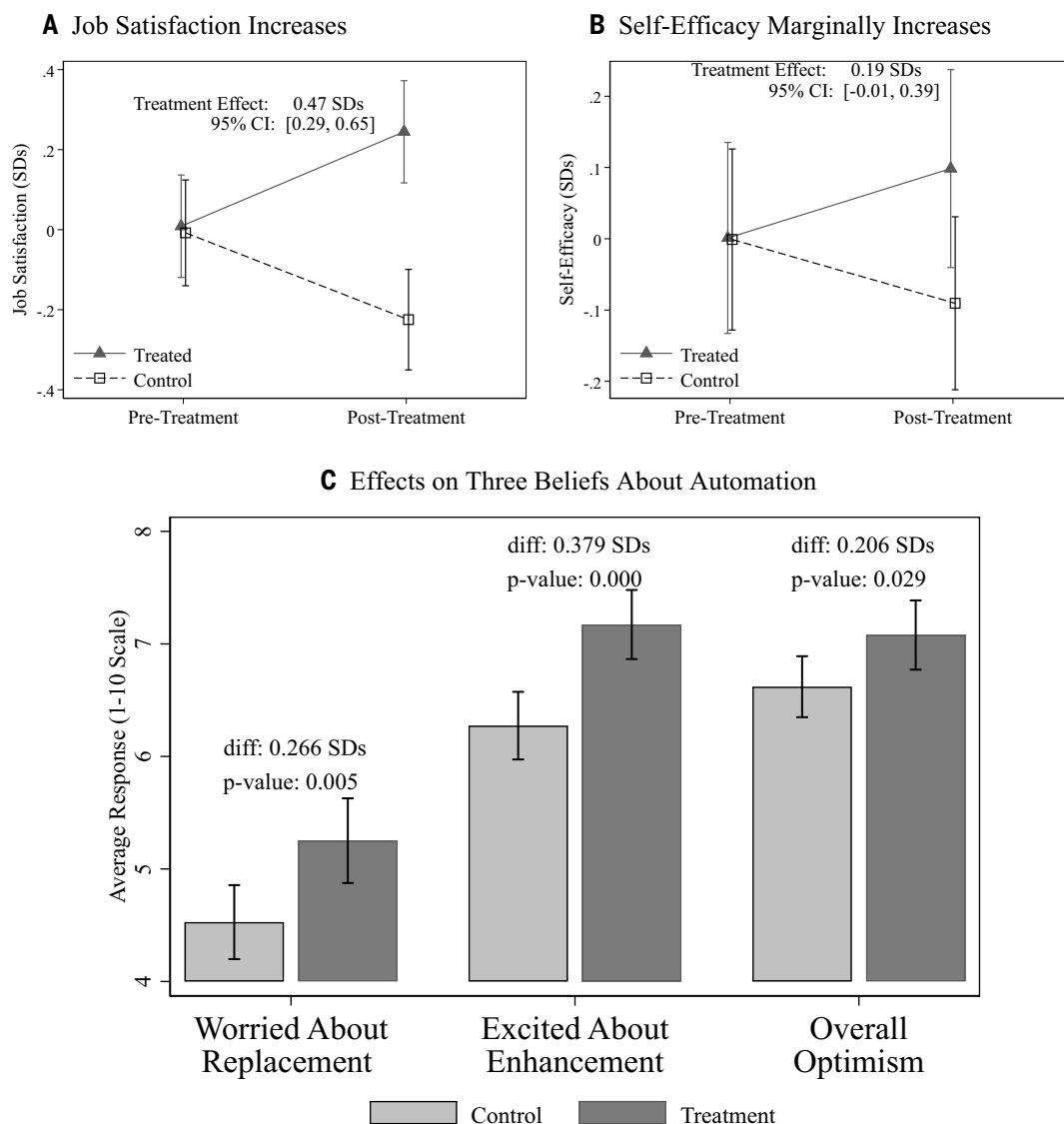
College-educated professionals performing midlevel professional writing tasks substantially increased their productivity when given access to ChatGPT. The generative writing tool increased the output quality of low-ability workers and reduced time spent on tasks for workers of all ability levels. At the aggregate level, ChatGPT reduced inequality. It is already being used by many workers in their real jobs.

These results are consistent with other studies showing productivity-enhancing and equalizing effects of recent AI technologies (8, 15, 16, 18). Relative to these studies, we analyzed productivity effects across several occupations and tasks, examined how workers use ChatGPT, measured subjective reactions to the technology, and documented persistent effects of our treatment on ChatGPT usage in real jobs.

Limitations

The experiment had several important limitations. We examined a limited range of occupations and tasks in which ChatGPT may be unusually useful. The tasks demanded clear, persuasive, relatively generic writing, which are arguably ChatGPT's central strengths. They did not require context-specific knowledge or precise factual accuracy. The version of ChatGPT used in this experiment cannot, by its

Fig. 3. Job satisfaction, self-efficacy, and beliefs about automation. (A) and (B) Job satisfaction and self-efficacy (originally elicited on scales of 1 to 10, normalized to have mean = 0 and SD = 1) before and after treatment in the treatment and control groups. Dots are means and error bars are 95% confidence intervals for means. Also shown are the coefficients on the treatment effect of a regression specified as in Fig. 1A. (C) Cross-sectional comparison of beliefs about automation in the treatment and control group, all on a scale of 1 to 10. The first question was “How worried are you about workers in your occupation being replaced by AI?” The second was “How optimistic are you that AI may make workers in your occupation more productive?” The third was “How do you feel about the impacts of future advances in AI (where 1 = very pessimistic and 10 = very optimistic).”



nature, access or supply context-specific knowledge and is not a reliable source of precise factual information.

The tasks could also be described through short, self-contained prompts, making use of ChatGPT easy, whereas many real-world tasks involve vaguer objectives and instructions, requiring workers to exercise initiative in determining what to do. Finally, participants in our tasks faced direct incentives in the form of bonus payments scaling with output quality, which encouraged them to maximize generic output quality and minimize time spent. White-collar workers are typically incentivized through longer-run promotion and firing incentives, which might instead encourage conspicuous exertion of effort or the development of a consistent personal style, both of which make ChatGPT less useful.

The tasks and incentive schemes were chosen to meet the constraints of the experimen-

tal design. We required short tasks that could be explicitly described for and performed by a range of anonymous workers online, and we needed to incentivize serious effort. Our judgment was that building factual accuracy requirements into the tasks would either result in tasks that felt artificial and unnatural (e.g., requiring participants to Google and report one or two specific facts) or overwhelm our budget for evaluators (e.g., giving participants an open-ended research task and exhaustively fact checking their assertions).

The aforementioned factors limit but do not eliminate the generalizability of our results. In real-world tasks, the need to fact-check ChatGPT's output will reduce its time-saving benefits, but the speed and writing quality increases observed in our experiment are sufficiently large that we suspect that ChatGPT will still often be useful. Moreover, newer versions of ChatGPT are more consistently factually ac-

curate, and some versions can access the internet to fact check themselves. We speculate that in more open-ended, real-world tasks, workers may find iterative rounds of prompting and discussion with ChatGPT useful even if they cannot immediately prompt out a final product. In these contexts, ChatGPT and human workers may be more strongly complementary than in our experiment. The importance of context-specific knowledge will also limit ChatGPT's utility, but there are plausible work-arounds. ChatGPT can be instructed to incorporate lists of context-specific factors, and organizations may be able to build customized ChatGPT-like models. Our follow-up surveys show that many workers do find ChatGPT useful in their real jobs.

Overall, we speculate that, relative to our experimental findings, the direct productivity effects of ChatGPT in the real economy will be somewhat lower and the technology will be

more strongly complementary to human workers. To what extent either of these is true remains an open question.

Implications

Our experiment captured only direct, immediate effects of ChatGPT on worker productivity. We could not examine the complex labor market dynamics that will arise as firms and workers adapt to ChatGPT. Several factors will mediate how the direct productivity impacts of ChatGPT affect wages and employment in exposed occupations. The first is the degree to which demand for goods produced by ChatGPT could expand as ChatGPT-fueled productivity increases make those goods much cheaper. For example, demand for programming services could plausibly expand massively if the price of those services fell. Aggregate programming employment might consequently increase. It is less clear whether demand for advertising or communication could expand as much, potentially entailing a reduction of employment in those sectors as fewer workers are needed to meet the same static demand. As an additional complication, ChatGPT might directly affect the composition of demand. For example, before ChatGPT, a piece of writing signaled that a company had invested at least some human labor, thought, and judgment into a message, which consumers might have appreciated; with this no longer being the case, demand for these messages could decrease (29).

The second factor is the nature and scarcity of the human skills best complemented by ChatGPT. Consider, for example, the use of ChatGPT to produce advertising content. Is this best accomplished by one senior advertising manager directly providing high-level guidance to ChatGPT or by 10 junior advertisers carefully designing prompts and editing ChatGPT's output? The answer will determine the structure of employment in the advertising sector. Similarly, suppose ChatGPT is highly complementary to human labor in programming tasks. If ChatGPT's human copilot needs to be an expert programmer capable of directly proofreading its output, then this could raise programmers' wages by boosting their productivity while their expertise remains scarce. If, by contrast, the complementary human role requires only basic programming knowledge and mainly involves checking output and refining natural-language prompts, then the pool of potential programmers would vastly increase and wages could fall even as productivity rises. More generally, tools such as ChatGPT could make expertise more accessible by facilitating learning (30).

Finally, the diffusion and effects of ChatGPT will also depend on organizational considerations that our experiment treating isolated individual workers does not address. ChatGPT might interact with traditional promotion and

hiring systems based partly on conspicuous exertion of effort. Large language models may be used to monitor or evaluate workers and avoid paying higher wages (31). Organizational and societal norms around the acceptability of using tools such as ChatGPT may take time to cohere and may significantly affect adoption of the technology (32–35).

Overall, the arrival of ChatGPT ushers in an era of vast uncertainty about the economic and labor market effects of AI technologies (36–38). Our experiment takes the first step toward answering the many questions that have arisen.

REFERENCES AND NOTES

1. D. H. Autor, *J. Econ. Perspect.* **29**, 3–30 (2015).
2. D. H. Autor, D. Dorn, *Am. Econ. Rev.* **103**, 1553–1597 (2013).
3. T. Eloundou, S. Manning, P. Mishkin, D. Rock, *arXiv:2303.10130* [econ.GN] (2023).
4. E. W. Felten, M. Raj, R. Seamans, SSRN [Preprint] (2023); <http://dx.doi.org/10.2139/ssrn.4375268>.
5. M. R. Frank et al., *Proc. Natl. Acad. Sci. U.S.A.* **116**, 6531–6539 (2019).
6. L. P. Boustan, J. Choi, D. Clingsmith, "Automation after the assembly line: computerized machine tools, employment and productivity in the United States" (National Bureau of Economic Research, 2022); <https://www.nber.org/papers/w30400>.
7. D. Acemoglu, P. Restrepo, "Robots and jobs: Evidence from US labor markets" (National Bureau of Economic Research, 2020); <https://www.nber.org/papers/w23285>.
8. K. Kanazawa, D. Kawaguchi, H. Shigeoka, Y. Watanabe, "AI, skill, and productivity: The case of taxi drivers" (National Bureau of Economic Research, 2022); <https://www.nber.org/papers/w30612>.
9. R. Arakawa, H. Yakura, M. Goto, *CatAlyst*: *arXiv:2302.05678* [cs.HC] (2023).
10. D. Acemoglu, P. Restrepo, *Am. Econ. Rev.* **108**, 1488–1542 (2018).
11. A. Agrawal, J. S. Gans, A. Goldfarb, *J. Econ. Perspect.* **33**, 31–50 (2019).
12. M. Hoffman, L. B. Kahn, D. Li, *Q. J. Econ.* **133**, 765–800 (2018).
13. J. Kleinberg, H. Lakkaraju, J. Leskovec, J. Ludwig, S. Mullainathan, *Q. J. Econ.* **133**, 237–293 (2018).
14. Productivity gains will translate into higher wages for workers if worker bargaining power or competition between employers for workers is sufficiently high to force employers to share part of the productivity gains with workers, and if the influx of workers into affected occupations is not large enough to completely offset these wage gains (see Discussion section).
15. E. Brynjolfsson, D. Li, L. R. Raymond, "Generative AI at work" (National Bureau of Economic Research, 2023); <https://www.nber.org/papers/w31161>.
16. S. Peng, E. Kalliamvakou, P. Cihon, M. Demirel, The impact of AI on developer productivity: Evidence from GitHub Copilot. *arXiv:2302.06590* [cs.SE] (2023).
17. A. Calderwood, V. Qiu, K. Ilonka Gero, L. B. Chilton, in *Joint Proceedings of the Workshops on Human-AI Co-Creation with Generative Models and User-Aware Conversational Agents co-located with 25th International Conference on Intelligent User Interfaces (IUI 2020)*, Cagliari, Italy, March 17, 2020. W. Geyer, Y. Khazaeni, M. Shmueli-Scheuer, Eds. (CEUR, 2018), vol. 2848 of *CEUR Workshop Proceedings*.
18. A. Campero et al., *arXiv:2206.12390* [cs.HC] (2022).
19. H. Schwabe, F. Castellacci, *PLOS ONE* **15**, e0242929 (2020).
20. A nascent literature has studied applications of machine learning to predictive tasks (consisting of yes/no diagnoses) (II).
21. E. Peer, D. Rothschild, A. Gordon, Z. Evernden, E. Damer, *Behav. Res. Methods* **54**, 1643–1662 (2022).
22. In each grading session, evaluators graded up to 14 responses. Evaluators received a base payment of \$16, plus up to \$8 in bonus payments depending on the correlation of their grades with the grades of other evaluators seeing the same responses.
23. This is the average correlation, across every pair of evaluators who saw the same essay, between the grade given by the first evaluator and the grade given by the second evaluator.
24. D. S. Lee, *Rev. Econ. Stud.* **76**, 1071–1102 (2009).
25. The choice to use ChatGPT is uncorrelated with treated respondents' salary, tenure, or grade on the first task.

26. Because of the small sample size in this group, the confidence interval includes zero and the treatment and control group differ in terms of average pretreatment grades (see the supplementary materials).
27. See the supplementary materials for a full description of this analysis.
28. The effects on beliefs about automation may also have dissipated because rising global awareness of the technology led to an equalization of familiarity between the treatment and control groups.
29. Y. Liu, A. Mittal, D. Yang, A. Bruckman, in *CHI '22: Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, New Orleans, LA, 29 April 2022 – 5 May 2022, S. Barbosa, C. Lampe, C. Appert, D. A. Shamma, S. Drucker, J. Williamson, K. Yatani, Eds. (Association for Computing Machinery, 2022); <https://dl.acm.org/doi/10.1145/3491102.3517731>.
30. J. Qadir, *TechRxiv* [Preprint] (2022); <https://doi.org/10.36227/techrxiv.21789434.v1>.
31. D. Acemoglu, A. F. Newman, *Eur. Econ. Rev.* **46**, 1733–1756 (2002).
32. J. Ayling, A. Chapman, *AI Ethics* **2**, 405–429 (2022).
33. J. Hohenstein et al., *arXiv:2102.05756* [cs.HC] (2021).
34. M. Suh, E. Youngblom, M. Terry, C. J. Cai, in *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, Yokohama, Japan, 8–13 May 2021, Y. Kitamura, A. Quigley, K. Isbister, T. Igarashi, P. Bjorn, S. Drucker, Eds. (Association for Computing Machinery, 2021); <https://dl.acm.org/doi/10.1145/3411764.3445219>.
35. H. Zohny, J. McMillan, M. King, *J. Med. Ethics* **49**, 79–80 (2023).
36. D. Autor, "The labor market impacts of technological change: From unbridled enthusiasm to qualified optimism to vast uncertainty" (National Bureau of Economic Research, 2022); <https://www.nber.org/papers/w30074>.
37. N. Kshetri, *IT Prof.* **22**, 63–68 (2020).
38. A. Agrawal, J. Gans, A. Goldfarb, *The Economics of Artificial Intelligence: An Agenda* (Univ. of Chicago Press, 2019).
39. Data and code for: S. Noy, W. Zhang, Experimental evidence on the productivity effects of generative artificial intelligence, *Open Science Framework* (2023); https://osf.io/xd7qw/?view_only=b8dc58dd6f44b979bf81069022fa392D01017605/OSF.IO/XD7QW.

ACKNOWLEDGMENTS

We thank the editor; three anonymous referees; and D. Acemoglu, N. Agarwal, D. Autor, L. Barros, T. Benheim, A. Finkelstein, J. Horton, S. Jaeger, A. Leslie, J. Mejia, I. Noy, L. Noy, E. Partridge, C. Raskin, A. Rao, N. Rousselle, C. Roth, F. Schilbach, B. Schoefer, L. Schubert, A. Shreekrumar, V. Vilfort, S. Wu, and participants at the MIT Labor Lunch for helpful comments and conversations. This experiment was approved by the MIT Committee on the Use of Humans as Experimental Subjects (protocol no. 2212000849). **Funding:** This work was supported by an Emergent Ventures grant, the Mercatus Center, George Mason University (S.N.), a George and Obie Shultz Fund grant, the MIT Economics Department (S.N.), and National Science Foundation Graduate Research Fellowship Grant 745302 (W.Z.). **Author contributions:** Conceptualization: S.N., W.Z.; Funding acquisition: S.N., W.Z.; Investigation: S.N., W.Z.; Methodology: S.N., W.Z.; Project administration: S.N., W.Z.; Supervision: S.N., W.Z.; Visualization: S.N., W.Z.; Writing – original draft: S.N., W.Z.; Writing – review and editing: S.N., W.Z. **Competing interests:** The authors declare no competing interests. **Data and materials availability:** This experiment was preregistered at the American Economic Association's registry for randomized controlled trials (<https://www.socialscienceregistry.org/trials/10882>). All data and code used in the analysis are available at the Open Science Framework (39). **License information:** Copyright © 2023 the authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original US government works. <https://www.science.org/about/science-licenses-journal-article-reuse>

SUPPLEMENTARY MATERIALS

science.org/doi/10.1126/science.adh2586
Materials and Methods
Supplementary Text
Figs. S1 to S21
Tables S1 to S3
MDAR Reproducibility Checklist

Submitted 20 February 2023; accepted 2 June 2023
10.1126/science.adh2586