

Does AI Qualify for the Job?

A Bidirectional Model Mapping Labour and AI Intensities

Fernando Martínez-Plumed
Universitat Politècnica de València
fmartinez@dsic.upv.es

Songül Tolan
Joint Research Centre
European Commission
songul.tolan@ec.europa.eu

Annarosa Pesole
Joint Research Centre
European Commission
annarosa.pesole@ec.europa.eu

José Hernández-Orallo
Universitat Politècnica de València
jorallo@upv.es

Enrique Fernández-Macías
Joint Research Centre
European Commission
enrique.fernandez-
macias@ec.europa.eu

Emilia Gómez
Joint Research Centre
European Commission
emilia.gomez-
gutierrez@ec.europa.eu

ABSTRACT

In this paper we present a setting for examining the relation between the distribution of research intensity in AI research and the relevance for a range of work tasks (and occupations) in current and simulated scenarios. We perform a mapping between labour and AI using a set of cognitive abilities as an intermediate layer. This setting favours a two-way interpretation to analyse (1) what impact current or simulated AI research activity has or would have on labour-related tasks and occupations, and (2) what areas of AI research activity would be responsible for a desired or undesired effect on specific labour tasks and occupations. Concretely, in our analysis we map 59 generic labour-related tasks from several worker surveys and databases to 14 cognitive abilities from the cognitive science literature, and these to a comprehensive list of 328 AI benchmarks used to evaluate progress in AI techniques. We provide this model and its implementation as a tool for simulations. We also show the effectiveness of our setting with some illustrative examples.

KEYWORDS

Labour market, tasks, AI intensity, AI impact, AI benchmarks, simulation

ACM Reference Format:

Fernando Martínez-Plumed, Songül Tolan, Annarosa Pesole, José Hernández-Orallo, Enrique Fernández-Macías, and Emilia Gómez. 2020. Does AI Qualify for the Job? A Bidirectional Model Mapping Labour and AI Intensities. In *Proceedings of the 2020 AAAI/ACM Conference on AI, Ethics, and Society (AIES '20), February 7–8, 2020, New York, NY, USA*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/3375627.3375831>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

AIES '20, February 7–8, 2020, New York, NY, USA

© 2020 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-7110-0/20/02...\$15.00

<https://doi.org/10.1145/3375627.3375831>

1 INTRODUCTION

In this paper we present a setting for the analysis and simulation of the *intensity* flows between Artificial Intelligence (AI) research and the labour market. Intensity is understood as the relevance of and effort spent on any undertaking. For instance, in the case of an occupation one can estimate how much time a particular activity requires. In the case of AI, one can estimate how much effort (in terms of activity) is devoted to a certain task in a particular area of research. Without a model, some direct connections can be made, such as the observation that progress in machine translation will have an impact on human translators, or that in order to rationalise the cost in language translation and subtitling of a major video-on-demand company, more progress of AI in this area would be needed. But the connections become more complex when we wonder how much AI research in natural language processing is affecting a lawyer, or what areas in AI should require more activity to alleviate the bottleneck of auditors, or any other profession. A traceable two-way analysis would be a more anticipatory and prescriptive analysis than just predicting what jobs are more suitable of automation, assuming things equal or extrapolating from a predictive model in which we cannot have any intervention. A model mapping labour and AI research that allows for counterfactuals could account for the relation between AI and labour in ways that could better represent different scenarios and guide policies according to them.

Differently from previous approaches that have tried to link directly AI developments with labour-related task characteristics [6], our framework adds an intermediate dimension of cognitive abilities which gives us greater flexibility as well as a broader understanding on the impact of AI on labour tasks. More precisely, on one side, we map 14 generic cognitive abilities taken from the cognitive science literature to 59 generic labour-related tasks from task-based surveys from the workplace. On the other side, we map these 14 generic abilities to a comprehensive list of 328 benchmarks used to promote and measure the progress in different areas of AI.

In this regard, we start with the detailed set of labour-related tasks (and occupations) from [11–13], which are assessed according to the cognitive abilities they typically require. Here we link these cognitive abilities to AI intensity indicators in terms of research activity and interest using AI benchmarks (see Figure 1). We also perform a cluster analysis to see how the AI benchmarks group

together given the underlying structure of their required cognitive abilities in order to further increase the interpretation of the results.

This mapping between tasks and AI benchmarks allows us to accurately assess how the intensity of AI research may affect work-related tasks and corresponding occupations, as well as the other way around: how task and occupation intensity should be translated to AI research. We then use this setting to rank tasks by potential AI impact, and to show which areas of AI research should be intensified to have an impact in particular selected tasks and occupations. The main contributions of this paper are summarised as follows:

- We propose a formal matrix-based bidirectional setting for the analysis of the impact between AI research and the labour market.
- We show how identifying the specific cognitive abilities that can be performed by AI gives a broader understanding on the impact of AI on labour tasks, and vice versa.
- We see the lack of alignment between the intensities coming from the activity in the workplace and the intensities coming from the activity in AI benchmarks.
- We provide a grouped interpretation of the activity in AI research by performing a cluster analysis on AI benchmarks given the underlying structure of their required cognitive abilities.
- We show how our setting allows for the analysis of counterfactual simulated scenarios and the identification of situations where AI research does not match the required abilities in the labour market.
- We develop an online visual approach¹ for showing the intensity flows between AI benchmarks and the labour market tasks and occupations.

2 RELATED WORK

The presented setting builds on the labour economics literature focused on measuring the potential for automation on the labour market [3, 6, 15, 23]. However, we have to draw a clear line between the impact and technological feasibility of AI to modify the workplace and the configuration of tasks and occupations, and a more simplistic view of AI as leading to full automation (substitution through machines) (anonymous). With this paper, we further complement the literature with a formal setting for measuring AI potential in cognitive abilities and, subsequently, in labour-related tasks and occupations. On the AI side, we perform this by relying on AI benchmarks, as used by researchers and industry to encourage and evaluate progress in AI, instead of relying on expert predictions on the future automatibility of occupations, as in [15] and subsequent studies. This is also in contrast to the use of models that quantify the probability of computerisation for different occupations based on their proportion of routine and non-routine tasks [5]. Furthermore, we complement Brynjolfsson et al.'s measure of "suitability for machine learning" for labour-related tasks [6], which draws upon particular technologies in machine learning only. Here, we use a more comprehensive list of AI tasks and benchmarks (which can be further extended and updated to future developments).

¹<https://safe-tools.dsic.upv.es/shiny/OTAAI/>

The use of AI benchmarks to analyse the state of the art of AI research has been popularised by the seminal work done by the Electronic Frontier Foundation (EFF) [9], and reports such as the AI Index [24], which also covers jobs briefly. Using the EFF data, [10] make a more explicit connection with the labour market. They measure progress in AI through linear trends in benchmarks across different metrics. However, due to nonlinear performance jumps at certain thresholds of each benchmark, progress in different benchmarks cannot be measured in a comparable manner. We address this issue by translating benchmarks to a measure of AI research activity, and not the incommensurate magnitudes of each benchmark. [10] introduce abilities, but they are specialised for "job task requirements", which limits its independence to the labour connection, and precludes a balanced bidirectional analysis.

In this paper, we integrate several theories of intelligence and cognition in psychology, animal cognition and AI textbooks to give a broader definition of abilities, as a more independent latent layer than human abilities (work-oriented) or AI abilities (technology-oriented). We draw information from a very comprehensive set of AI benchmarks, competitions and tasks (see section 3 for details), ensuring a broad coverage of AI tasks. Unlike many of the previous approaches, we formalise our setting by proposing a unified matrix-based mathematical model for the specification of dynamic intensities for AI and labour tasks. This formalisation allows for the analysis of intensity flows between AI and labour tasks (in both directions) analytically. This makes it possible to study real scenarios as well as simulated ones, using counterfactual or speculative hypotheses varying the intensity levels across tasks or AI benchmarks.

3 DATA

For the two extremes of our mapping, as shown in Figure 1, we need to rely on very different sources of data. We start with a description of labour-related task intensity before moving to a description of research intensity in AI.

3.1 Tasks and occupations

We gather the data about labour-related tasks and occupations from [11–13], comprising a list of tasks and their respective intensity (i.e. relevance and time spent) across occupations.

Concretely, we classify occupations according to the 3-digit International Standard Classification of Occupations (ISCO-3)². Since there is no international data source that covers the full classification, we combine data from three different sources: (1) the European Working Conditions Survey (EWCS)³; (2) the OECD Survey of Adult Skills (PIAAC)⁴; and (3) the database from the Occupational Information Network (O*NET)⁵. While (1) and (2) are surveys that provide data measured at the individual worker level based on replies to questions on what they do at work, (3) is also based on employer job postings, expert research and other sources. O*NET is widely used in the literature on labour markets and technological change [1, 15, 16] and it covers a large share of the task list

²<https://www.ilo.org/public/english/bureau/stat/isco/>

³<https://www.eurofound.europa.eu/surveys/european-working-conditions-surveys>

⁴<https://www.oecd.org/skills/piaac/>

⁵<https://www.onetonline.org/>

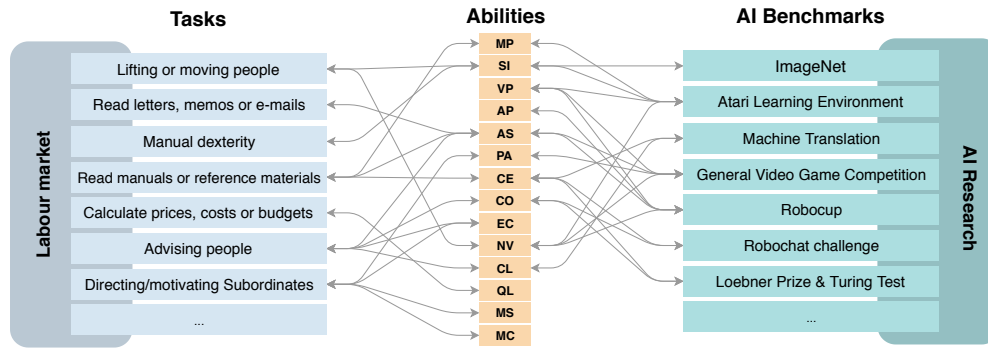


Figure 1: Bidirectional and indirect mapping between job market and Artificial Intelligence. The notation we use is t for the tasks, a for the abilities and b for the benchmarks. The arrows are represented by correspondence matrices W (task-ability correspondence) and R (ability-benchmark correspondence).

that we use in our analysis. However, the occupational level of the data precludes a further analysis of the variation in task content within occupations. Moreover, much like the EWCS for Europe, the O*NET is based on US data only. Therefore, even if there are likely differences in the task content of occupations across countries (due to institutional as well as socio-economic differences) we cannot analyse these differences in the present analysis.

In these sources, task intensity for different occupations is derived either as a measure of time spent on specific tasks (e.g., the intensity for the task “Lifting or moving people” is obtained from survey question “Does your main paid job involve lifting or moving people?” and the corresponding 7-point scale answers ranging from “All of the time” to “Never”), or curated by occupational experts and provided on a standardised occupational level (e.g., the extent to which the task is required to perform a job). Due to the varying nature of survey data, we need to be aware of issues such as measurement error, high variation in responses across individuals and biased responses. Consistency in the measurement of task intensity across the different data sources is measured with Cronbach’s alpha, which is calculated from the pairwise correlation between items that measure similar concepts. All tests yield high correlations and Cronbach’s Alpha values of between 0.8 and 0.9.

Finally, in order to make the measures of task intensity comparable across all three data sources, we equalise scales and levels of all variables. For this purpose, we rescale the variables to a $[0, 1]$ scale with 0 representing the lowest possible intensity and 1 representing the highest possible intensity of each variable. Moreover, we average scores measured on an individual level (i.e., all variables from PIAAC and EWCS) to the unified level of standardised 3-digit occupation classifications. The final database contains the intensity of 59 tasks across 119 different occupations.

3.2 AI benchmarks

We consider a comprehensive set of AI benchmarks for our setting based on our own previous analysis and annotation of AI papers [17, 19–21] as well as open resources such as *Papers With Code*⁶ (the largest, up-to-date, free and open repository of machine learning papers). It includes data from multiple (verified) sources, including

⁶<https://paperswithcode.com/>

academic literature, review articles and code platforms focused on machine learning and AI.

From the aforementioned sources we track the reported evaluation results on different metrics of AI performance across separate AI benchmarks (e.g., tasks, datasets, competitions, awards, etc.) from a number of AI domains. We cover computer vision, speech recognition, music analysis, machine translation, text summarisation, information retrieval, robotic navigation and interaction, automated vehicles, game playing, prediction, estimation, planning, automated deduction, among others. This ensures a broad coverage of AI tasks, well beyond perception, such as the ability to plan and perform actions on such plans. Specifically, our framework uses data from 328 different AI benchmarks, after selecting those with enough information available for different evaluation metrics.

When aiming at evaluating the progress in specific AI areas, we need to pay attention to the set of criteria about how a system is to be evaluated. Even if the metrics that are used in each benchmark improve, it would be misleading to consider that the progress in AI should be analysed by aggregating these values. First, these magnitudes are incommensurate, so aggregating the score in a video game with the result of translation task is meaningless. Second, the results are obtained by specific systems solving particular tasks. There is no understanding on how to build systems that can solve all these tasks at the same time.

Therefore, instead of using the rate of progress with particular performance metrics, we analyse the activity level or *intensity* for a benchmark, measured in terms of the production (e.g., outputs such as research publications, news, blog-entries, etc) from the AI community. Benchmarks that have increasing trends in their production rates –not their performance metrics– indicate that more AI researchers and practitioners are working on them (i.e., there is a clear research effort and intensity). Note that this is not an indication of progress, although, presumably, effort may lead to some progress eventually.

In order to derive the activity level or *intensity*, we use some proxies. In particular, we performed a quantitative analysis using data obtained from *AI topics*⁷, an archive kept by the *Association*

⁷<https://aitopics.org>

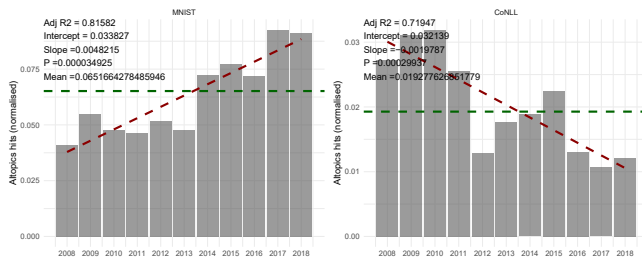


Figure 2: Average rate of activity level or intensity (green dashed line) for a couple of illustrative AI benchmarks over the last decade (2008-2018).

for the Advancement of Artificial Intelligence (AAAI)⁸. This platform contains a myriad of AI-related documents (e.g. news, blog entries, conferences, journals and other repositories from 1905 to 2019) that are collected automatically with NewsFinder [7]. In this regard, in order to calculate the intensity in each particular benchmark, we average the number of hits (e.g., documents) obtained from *AI topics* per benchmark and year over a specific period of time. Note that the number of hits are normalised to sum up to 100% per year. Figure 2 shows the activity trends for two different benchmarks. Our measure of intensity is the average over the period 2008-2018.

4 MODEL

In the following subsections we describe the main components of our model, as originally illustrated in Figure 1. We use the following notation:

- \mathbf{t} : (labour) task intensity vector.
- \mathbf{W} : task-ability correspondence matrix.
- \mathbf{a} : ability vector.
- \mathbf{R} : ability-benchmark correspondence matrix.
- \mathbf{b} : benchmark intensity vector.

We define them in more detail below.

4.1 Intensity vectors

Vector \mathbf{t} denotes task intensities. In section “Tasks and occupations” we described the data we use, meaning that \mathbf{t} have dimension (59×1) , on a $[0, 1]$ scale with 0 and 1 representing the lowest and highest possible intensity respectively. This vector reflects the occupational task intensity in the abilities assigned to the tasks in each occupation (note that each occupation has a different \mathbf{t} vector).

On the other hand, \mathbf{b} denotes a benchmark intensity vector (328×1) , with relative values in $[0, 1]$. This vector shows the average (normalised) number of documents obtained from *AI topics* per benchmark and year over a specific period of time as explained in section “AI benchmarks”.

4.2 Cognitive abilities

In previous works [2, 4], labour-related tasks and those that are usually set in AI as capacities are usually matched directly, even if the elements on the left list in Figure 1 are very different from the elements on the right. However, tasks and benchmarks can be

⁸<https://www.aaai.org/>

mapped through an intermediate layer of latent factors, what we refer to as ‘cognitive abilities’, also at a level of aggregation that is more insightful. For this characterisation of abilities we look for an intermediate level of detail, excluding very specific abilities and skills (e.g., music skills, mathematical skills, hand dexterity, driving, etc.) but also excluding very general abilities or traits that would influence all the others (general intelligence, creativity, etc.). As we just cover cognitive abilities, we also exclude personality traits (e.g., the big five [14]). Although we consider the latter essential for humans, their ranges can be simulated in machines by changing goals and objective functions.

For our purposes we use 14 categories as the result of the integration of several tables and figures from [18], originally collected from psychometrics, comparative psychology, cognitive science and artificial intelligence (see Figure 1). The 14 categories are defined as follow: *Memory processes* (MP), *Sensorimotor interaction* (SI), *Visual processing* (VP), *Auditory processing* (AP), *Attention and search* (AS), *Planning and sequential decision-making and acting* (PA), *Comprehension and compositional expression* (CE), *Communication* (CO), *Emotion and self-control* (EC), *Navigation* (NV), *Conceptualisation, learning and abstraction* (CL), *Quantitative and logical reasoning* (QL), *Mind modelling and social interaction* (MS), and *Metacognition and confidence assessment* (MC). The hierarchical theories of intelligence in psychology, animal cognition and the textbooks in AI are generally consistent (at least partially) with this list of abilities, or in more general and simple terms, with this way of organising the vast space of cognition. The definition of the cognitive abilities can be found in [27].

4.3 Mapping

To generate the mapping between labour-related tasks and cognitive abilities, a multidisciplinary group of researchers conducted an annotation exercise for each item of the task database. More precisely, in a cross-tabulation of the vector of tasks \mathbf{t} of length $p = |\mathbf{t}| = 59$ and cognitive abilities \mathbf{a} of length $m = |\mathbf{a}| = 14$, each annotator was asked to put a 1 in a task-ability correspondence matrix \mathbf{W} (59×14) if an ability is inherently required, i.e. absolutely necessary to perform the respective task (see the rubric in the Appendix [22, Section A]). In order to increase robustness in the annotations, we followed a *Delphi Method* approach [8], repeating this process in order to increase agreement among annotators, and finally obtaining the share in percentage terms for each combination of task and ability. Similarly, we also linked the cognitive abilities with our list of AI benchmarks (which will be also described in detail in the following sections). Specifically, a group of AI-specialised researchers was asked to consider how each AI benchmark is related to each cognitive ability: in a cross-tabulation of the vector of benchmarks \mathbf{b} of length $n = |\mathbf{b}| = 328$ and cognitive abilities \mathbf{a} of length $m = |\mathbf{a}| = 14$, we put a 1 in the ability-benchmark correspondence matrix \mathbf{R} (14×328) if an ability is inherently required, i.e. absolutely necessary to solve the respective benchmark. Full information about this mapping procedure can be found in [25, 26]

4.4 Two-way interpretation

We can then translate the benchmark intensity vector \mathbf{b} to cognitive abilities as a matrix-vector multiplication $\mathbf{Rb} \rightarrow \mathbf{a}$ thus obtaining

an ability intensity vector \mathbf{a} (14×1). We can also analyse task intensity, by weighting the task-ability mapping matrix by the ability intensity vector \mathbf{a} as a matrix-vector multiplication $\mathbf{W}\mathbf{a} \rightarrow \mathbf{t}$ thus obtaining a new task intensity vector \mathbf{t} (59×1).

This gives us a leftward interpretation of Figure 1 as:

$$\mathbf{R}\mathbf{b} \rightarrow \mathbf{a} \text{ and } \mathbf{W}\mathbf{a} \rightarrow \mathbf{t}$$

which together makes $\mathbf{W}\mathbf{R}\mathbf{b} \rightarrow \mathbf{t}$. This is interpreted as “benchmarks require abilities, which are required for tasks”.

By using this framework we can analyse flows in both directions mathematically. Therefore, we can also give the rightward interpretation as:

$$\mathbf{t}^T\mathbf{W} \rightarrow \mathbf{a}^T \text{ and } \mathbf{a}^T\mathbf{R} \rightarrow \mathbf{b}^T$$

which together makes $\mathbf{t}^T\mathbf{W}\mathbf{R} \rightarrow \mathbf{b}^T$. This is interpreted as “tasks require abilities, which are required for benchmarks”.

Note that since both \mathbf{W} and \mathbf{R} mean “requires” (in the direction of abilities), it makes sense to distribute the values when a task or a benchmark requires many abilities. So, assuming that more abilities require more effort, we normalise both \mathbf{W} and \mathbf{R} through abilities. This means that in \mathbf{W} rows are normalised to sum up 1, and in \mathbf{R} columns are normalised to sum up 1, and values are thus in $[0, 1]$.

5 ANALYSIS AND RESULTS

We analysed the correspondence between the two edges of our model. By comparing the values of \mathbf{b} as propagated rightwards from \mathbf{t} ($\mathbf{t}^T\mathbf{W}\mathbf{R} \rightarrow \mathbf{b}^T$) against the values of \mathbf{b} that originate directly from the benchmark intensities, we see very low correlations between these vectors. Figure 5 in the Appendix [22, Section E] shows some discrepancy scatterplots illustrating this. This picture is general, and we can conclude that the intensities do not match: the focus on AI benchmarks today does not correspond with the labour activities having highest intensity according to our data. Could this be different? In order to answer this question, in what follows we analyse the results bidirectionally, exploring several hypotheses and professional profiles.

5.1 From AI to labour

As an illustrative example of how the model can be used in a single direction, we can obtain the task intensity vector \mathbf{t} from the original benchmark intensity vector \mathbf{b} . This illustrates the leftward interpretation of Figure 1.

While in Figure 3 in the Appendix [22, Section C] we show how our model works when specific AI benchmarks are selected, Figure 3 shows a sorted list of labour tasks according to the computed values in \mathbf{t} from the analysis of *AI topics*. Those with the highest values consist almost entirely of information gathering and processing tasks (e.g., read letters or manuals, articles, bills, etc.), as well as performing tasks without using explicit instructions, relying on patterns and inference instead (e.g., learning, solving unforeseen problems, learning-by-doing, etc.). On the other hand, the lowest-scoring tasks are largely non-cognitive tasks that require a high degree of physical effort and dexterity (e.g., steadiness, manual/finger dexterity, etc.). This probably reflects a limited coverage of robotic benchmarks, which usually involve more propriosensory perception and manipulation. At the same time, there are also plenty

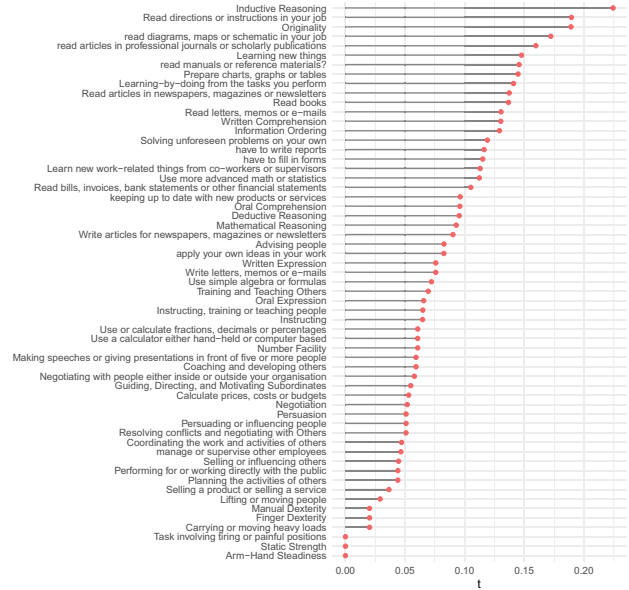


Figure 3: Labour-related tasks ranked in descending order based on by their intensity vector \mathbf{t} .

of interpersonal tasks that include a human component. These are considered non-routine tasks (e.g., persuasion, supervision, communication or people management, etc.), all of which generally require social and emotional skills.

Note that the above considers the current activity (as extracted from the *AAAI AI topics* data) and the tasks that would be affected if this activity would be transformed into progress in the areas the benchmarks represent and assuming that different abilities can be combined seamlessly.

5.2 From labour to AI

Following the rightward interpretation of our setting, we can also analyse, given a particular (set of) occupation(s) and their corresponding set of tasks, which sort of AI benchmarks should attract more interest or require more effort from the AI research community in order to have a potential impact in the selected occupation(s).

We can do (1) one specific labour-related task or (2) a combination of tasks conforming particular occupations. Figure 4 in the Appendix [22, Section D] shows some illustrative examples of (1). Regarding (2), we can also compute the AI benchmarks intensity scores by selecting relevant occupations from the ISCO-3 specifications. In this sense, we focus on nine illustrative occupations: (a) general office clerks; (b) shop salespersons; (c) agricultural, forestry and fishery labourers; (d) medical doctors; (e) mining and construction labourers; (f) sales, marketing and public relations professionals; (g) mobile plant operators; (h) waiters and bartenders; (i) market gardeners and crop growers.

Because of the large number of AI benchmarks (328), we have clustered these benchmarks into six groups to make the interpretation of results easier (details in the Appendix [22, Section B]). Figure

4 depicts benchmark intensity scores for the nine selected occupations mentioned above. For instance, in order for AI developments to have an effect on general office clerks, AI research should focus on those benchmarks related to inspection and data extraction as well as on those focused on the development of narratives, question answering and social interaction.

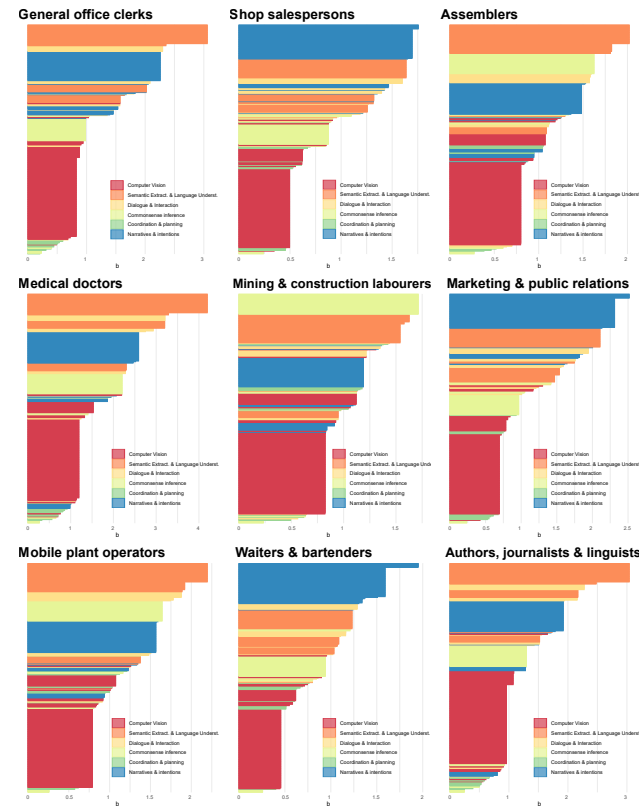


Figure 4: AI benchmarks ranked in descending order based on by their intensity vector b given their task intensity vectors t from six different occupations. Benchmarks coloured according to the cluster they belong to.

If we pay attention to those benchmarks where more progress is apparently taking place in AI (visual and auditory perception using deep learning and sensorimotor interaction, through (deep) reinforcement learning), we see that these cognitive abilities are generally at the bottom for the nine selected occupations. This means either that (1) some of these skills are taken from granted (e.g., recognising objects and moving around in the workplace) or (2) many tasks in the workplace require skills for which there is not a high AI research activity at the moment. About (1), in our annotations, we included abilities when ‘absolutely necessary’. Consequently, we considered that many of the tasks used in the workplace do not inherently require that a robot or a human visually recognises static or moving elements, as other capabilities could be used instead (e.g., blind people may “read manuals or reference manuals” using Braille).

6 CONCLUSIONS

We have developed a setting for the analysis of the relationship between Artificial Intelligence and the labour market in both directions. The setting combines occupations and tasks from the labour market with AI research benchmarks through an intermediate layer of cognitive abilities. The identification of the specific cognitive abilities that can be performed by AI gives a broader understanding on the impact of AI, as the inner layer is more independent of particular occupations, tasks or AI benchmarks. Although not included in the paper, we can also generate simulations outwards, setting a particular combination of ability intensities and propagate how tasks and occupations would be affected and what benchmarks would be more relevant. This analysis could also be done inwards.

In the paper we have seen examples where we can assess, in a very detailed way, how technological intensity of AI research may affect work-related tasks and corresponding occupations, as well as the other way round: how task and occupation intensity should be translated into AI research. We have seen the discrepancy between AI intensity and labour intensity and have used this setting to rank tasks by potential AI impact. In the end, we can determine which areas of AI research should be intensified if we sought to have a technological impact in particular selected task and occupations.

Despite its popularity in AI, using AI benchmarks to pulse the progress of AI research is fraught with caveats and criticisms, especially if performance metrics are used as an indication of progress. Instead, our model is based on intensities: we analyse whether some located activity on one edge translates on some located activity on the other edge. We use proxies for activities (such as time spent in a particular labour-related task or the research activity as per Figure 2). The use of activity versus progress makes this setting adoptable for the governance and assessment of AI R&D in academia and industry. In future work this analysis can be refined as more data becomes available on the relevance of specific work-related tasks as well as new AI benchmarks are introduced. Overall, we already present a powerful and flexible open tool⁹ to map AI research and the impact on labour bidirectionally. The major merit of our model is not being predictive, but being prescriptive: we can decide priorities and make AI research interventions accordingly, to procure that AI does qualify for the job.

ACKNOWLEDGMENTS

This material is based upon work supported by the EU (FEDER), and the Spanish MINECO under grant RTI2018-094403-B-C3, the Generalitat Valenciana PROMETEO/2019/098. F. Martínez-Plumed was also supported by INCIBE (Ayudas para la excelencia de los equipos de investigación avanzada en ciberseguridad), the European Commission (JRC) HUMAINT project (CT-EX2018D335821-101), and UPV (PAID-06-18). J. H-Orallo is also funded by an FLI grant RFP2-152.

⁹The presented setting and posterior analysis is flexible by updating data about benchmarks and professions, as well as the computed rates of intensity in AI benchmarks as measured using *AI topics*. Further details about the complete set of occupations, tasks, benchmarks and the associated intensity rates based on the results from *AI topics* or work surveys can be found in *link anonymised*.

REFERENCES

- [1] Daron Acemoglu and David Autor. 2011. Skills, tasks and technologies: Implications for employment and earnings. In *Handbook of labor economics*. Vol. 4. Elsevier, 1043–1171.
- [2] Daron Acemoglu, David Dorn, Gordon H Hanson, Brendan Price, et al. 2014. Return of the Solow paradox? IT, productivity, and employment in US manufacturing. *American Economic Review* 104, 5 (2014), 394–99.
- [3] Melanie Arntz, Terry Gregory, and Ulrich Zierahn. 2016. The Risk of Automation for Jobs in OECD Countries: A Comparative Analysis. *OECD Social, Employment and Migration Working Papers* 2, 189 (2016), 47–54. <https://doi.org/10.1787/5jlz9h56dvq7-en>
- [4] David Autor. 2013. *The “task approach” to labor markets: an overview*. Technical Report. National Bureau of Economic Research.
- [5] David H Autor, Frank Levy, and Richard J Murnane. 2003. The skill content of recent technological change: An empirical exploration. *The Quarterly journal of economics* 118, 4 (2003), 1279–1333.
- [6] Erik Brynjolfsson, Tom Mitchell, and Daniel Rock. 2018. What Can Machines Learn, and What Does It Mean for Occupations and the Economy?. In *AEA Papers and Proceedings*, Vol. 108. 43–47.
- [7] Bruce G Buchanan, Joshua Eckroth, and Reid Smith. 2013. A Virtual Archive for the History of AI. *AI Magazine* 34, 2 (2013), 86–86.
- [8] Norman Dalkey and Olaf Helmer. 1963. An experimental application of the Delphi method to the use of experts. *Management science* 9, 3 (1963), 458–467.
- [9] Peter Eckersley, Yomna Nasser, et al. 2017. EFF AI progress measurement project.
- [10] Edward W Felten, Manav Raj, and Robert Seamans. 2018. A Method to Link Advances in Artificial Intelligence to Occupational Abilities. In *AEA Papers and Proceedings*, Vol. 108. 54–57.
- [11] Enrique Fernández-Macias and Martina Bisello. 2017. *Measuring The Content and Methods of Work: a Comprehensive Task Framework*. Technical Report. European Foundation for the Improvement of Living and Working Conditions.
- [12] Enrique Fernández-Macias, Martina Bisello, Sudipa Sarkar, and Sergio Torrejón. 2016. *Methodology of the construction of task indices for the European Jobs Monitor*. Technical Report. European Foundation for the Improvement of Living and Working Conditions.
- [13] Enrique Fernández-Macias, Emilia Gómez, José Hernández-Orallo, Bao Sheng Loe, Bertin Martens, Fernando Martínez-Plumed, and Songül Tolan. 2018. A multi-disciplinary task-based perspective for evaluating the impact of AI autonomy and generality on the future of work. *CoRR abs/1807.02416* (2018). [arXiv:1807.02416](http://arxiv.org/abs/1807.02416)
- [14] Donald W Fiske. 1949. Consistency of the factorial structures of personality ratings from different sources. *The Journal of Abnormal and Social Psychology* 44, 3 (1949), 329.
- [15] Carl Benedikt Frey and Michael A Osborne. 2017. The future of employment: how susceptible are jobs to computerisation? *Technological Forecasting and Social Change* 114 (2017), 254–280.
- [16] Maarten Goos, Alan Manning, and Anna Salomons. 2009. Job polarization in Europe. *American economic review* 99, 2 (2009), 58–63.
- [17] José Hernández-Orallo. 2017. Evaluation in artificial intelligence: from task-oriented to ability-oriented measurement. *Artificial Intelligence Review* 48, 3 (01 Oct 2017), 397–447. <https://doi.org/10.1007/s10462-016-9505-7>
- [18] José Hernández-Orallo. 2017. *The measure of all minds: evaluating natural and artificial intelligence*. Cambridge University Press.
- [19] Fernando Martínez-Plumed, Shahar Avin, Miles Brundage, Allan Dafoe, Sean Ó hÉigeartaigh, and José Hernández-Orallo. 2018. Accounting for the neglected dimensions of AI progress. *arXiv preprint arXiv:1806.00610* (2018).
- [20] Fernando Martínez-Plumed and José Hernández-Orallo. 2018. Analysing Results from AI Benchmarks: Key Indicators and How to Obtain Them. *arXiv preprint arXiv:1811.08186* (2018).
- [21] F. Martínez-Plumed and J. Hernandez-Orallo. 2018. Dual Indicators to Analyse AI Benchmarks: Difficulty, Discrimination, Ability and Generality. *IEEE Transactions on Games* (2018), 1–1. <https://doi.org/10.1109/TG.2018.2883773>
- [22] Fernando Martínez-Plumed, Songül Tolan, Annarosa Pesole, José Hernández-Orallo, Enrique Fernández-Macias, and Emilia Gómez. 2020. Does AI Qualify for the Job? A Bidirectional Model Mapping Labour and AI Intensities (Appendix). <http://hdl.handle.net/10251/133314> (2020).
- [23] Ljubica Nedelkoska and Glenda Quintini. 2018. *OECD Social, Employment and Migration Working Papers No. 38*. Technical Report 38. 1–124 pages.
- [24] Yoav Shoham, Raymond Perrault, Eric Brynjolfsson, Jack Clark, James Manyika, Juan Carlos Niebles, Terah Lyons, John Etchemendy, and Z Bauer. 2018. The AI Index 2018 Annual Report.
- [25] Songül Tolan, Annarosa Pesole, Fernando Martínez-Plumed, Enrique Fernández-Macias, José Hernández-Orallo, and Emilia Gómez. 2019. Artificial Intelligence and Jobs: From Tasks to Cognitive Abilities. *RENIR Workshop on the impact of automation and artificial intelligence on regional economies, May 27-28* (2019).
- [26] Songül Tolan, Annarosa Pesole, Fernando Martínez-Plumed, Enrique Fernández-Macias, José Hernández-Orallo, and Emilia Gómez. 2020. Measuring the occupational impact of AI beyond automation: tasks, cognitive abilities and AI benchmarks. *Submitted for publication* (2020).
- [27] Karina Vold and Jose Hernandez-Orallo. 2019. AI Extenders: The Ethical and Societal Implications of Humans Cognitively Extended by AI. In *AAAI / ACM Conference on Artificial Intelligence, Ethics, and Society*.