

# Humans in the Loop.

Visibility of the hidden work under artificial intelligence training.



A1NW6Q6P3NBHWE

A232HEY81AMQQ8

A2BSAP480FLLOJ

A2E1PL4A2EQNNB

A2J86NI89IOP17

A21X40K1JQVQZJ

A2YFFC90MO0MW8

AWZKTYEDKNIRB



Humans in the Loop.

Visibility of the hidden work under artificial intelligence training.





Humans in the Loop.  
Visibility of hidden work in artificial intelligence training.

Research by Adrián Lombera  
Supervised by Aino Abella and Rafa Pozo

Defended on 12 June 2025

Printed in Apunts. Economia cooperativa.  
Barcelona, June 2025



Humans in the Loop: Visibility of the hidden work under artificial intelligence training. © 2025 by Adrián Lombera Cruz is licensed under Creative Commons Attribution-ShareAlike 4.0 International







## **Acknowledgements**

This research project has been made possible thanks to the generosity and commitment of the data workers at DataLabellers.org in Kenya, who shared their time, experiences and reflections with invaluable openness. I would also like to thank ELISAVA for its guidance and institutional support, especially my tutors Ainoa Abella and Rafa Pozo, whose guidance has been key throughout the process. I would also like to thank the interviewees who are part of the ELISAVA teaching staff, Ariel Guerzensvaig and Alicia de Manuel, for their willingness to engage in dialogue from their dual academic and professional roles. Finally, I extend my gratitude to professors Iván Paz, member of the Axolot collective, and David Berga for their generosity with technical assistance throughout the development of this work.







**Work is not destroyed, it is merely transferred.**

*i. located glossary*

**algorithm** A set of coded instructions or rules that a computer follows to perform a task or solve a problem. In AI, algorithms allow data to be processed, decisions to be made or patterns to be learned by training on previous examples.

**data annotation** Marking data with relevant tags or metadata to be useful in training artificial intelligence models. Without it, algorithms would not understand the material provided.

**data colonialism** Data mining relationship between dominant technological centres and vulnerable peripheries.

**crowdsourcing** In the context of invisibilised labour, refers to the outsourcing of digital micro-tasks to a multitude of anonymous, online workers via platforms who are paid per task completed.

**data work** Human activities required to collect, prepare, classify and maintain data for AI systems.

**data set** A purposeful set of data resulting from the processing of raw data by data workers. These datasets subsequently serve as resources to feed "intelligent" models.

**de-skilling** The process of simplifying complex work into small, quick tasks that reduce the knowledge required to perform them. It is an effect of digital Taylorism.

offshoring	The practice of moving operations to another country in order to reduce costs.
outsourcing	companies commonly referred to as subcontractors. In this context they are companies outside the borders of the West that manage the labour force working for companies in the Global North.
data extractivism	A notion that describes, through a parallel with natural resource extraction, the operations that transfer data, human effort and natural resources for the operation of digital infrastructures from the Global South to the North.
fauxtimation	Simulation of automation in which human workers perform tasks impersonating AI, thus masking their lack of technical autonomy.
AI futurism	Narratives that present AI as an inevitable and desirable destiny for the development of the human species.
ghost work	This term coined by Mary Gray and Siddharth Suri refers to the invisible, precarious and poorly paid human labour that is essential in technological supply chains. In the context of artificial intelligence it manifests itself through the data work required to create the datasets from which algorithms are trained.
heteromation	A hybrid technological system that seeks to automate operations by including workers in key processes in order to save costs for the

company. Sometimes these workers are not aware that they are making a vital effort for the business model and if they are, they are often overshadowed and poorly paid, compared to a permanent employee.

**human in the loop** A term that refers to human intervention within an automated process to solve tasks that the machine cannot perform on its own.

**infraestructure** Technical, human and organisational networks that support the operation, in this case of AI.

**algorithmic management** The use of algorithms to assign tasks, monitor performance and manage workers without direct human supervision.

**microtask** A brief, specific, minimal unit of digital work that is often part of a larger process. Examples include identifying objects in images, transcribing text fragments or verifying data.

**model** A code structure built by an algorithm from data. It represents in mathematical form the relationships and patterns it learns, allowing it to make predictions, classifications or decisions. It can be expressed in formulas, numbers or complex networks, depending on the task.

**inevitability motif** A discursive device that presents technological progress as inevitable in order to avoid public scrutiny and promote passive acceptance. In the case of artificial intelligence, it is applied through.



strategic occlusion	Deliberate concealment of how AI is produced or managed for commercial or power purposes
last mile paradox of automation	Paradox that shows that despite the automation promised by AI, more and more jobs are being created that are needed to sustain it.
exclusionary pedagogy	A training approach that imposes dominant categories and frameworks to the exclusion of other forms of knowledge. In AI, this translates into an ethics reduced to computable rules, without a deeper epistemological engagement that enriches ethical thinking.
imperceptible production	Work essential to AI that is rendered invisible by narratives of complete automation.
pseudo-IA	Systems presented as intelligent that rely in reality on hidden human decisions.
semantic segmentation	Annotation of images in which each pixel is classified into a recognisable category. Without them the algorithms would not understand the images provided.
global south	A group of countries, mostly in Asia, Africa, and Latin America, that have historically been marginalised from the global order and do not enjoy sufficient representation in international institutions.
taskification	Breakdown of work into minimal, repetitive tasks that are digitally distributed. Main reason for de-skilling.

digital Taylorism

Optimisation of work tasks by monitoring and fragmentation of tasks.

invisible collar labour

Hidden digital labour part of the technical infrastructure of some digital products and services needed to operate them smoothly.

platform labour

A form of labour mediated by platforms that act as digital labour exchanges.

turkers

Mechanical Turk workers who perform digital tasks through the Amazon Mechanical Turk platform.

uberisation

A term coined by Braz et al that defines the transformation of work into on-demand tasks managed by applications, without formal employment relationships.

## *ii. Acronyms and abbreviations*

BPO	Business Process Outsourcing. These are outsourced companies specialised in providing data work for artificial intelligence projects.
DLA	The Data Labelers Association was founded in Kenya at the beginning of 2025.
GenAI	Generative artificial intelligence is a type of autonomous technology that creates new content by analysing patterns. Models such as ChatGPT, Gemini, Claude and DeepMind are generative models.
LLM	Large Language Model. These are models that have been trained with large volumes of data to achieve advanced capabilities and high performance in complex tasks. This makes them useful for more general cases. The most well-known commercial models, such as those accessible via OpenAI or Google, function thanks to this technology.
YOLO	You Only Look Once is a pre-trained, open-source visual recognition model.

### *iii. list of graphic elements*

*Figure 01. Map of data flows between countries. 26*

*Figure 02. Unknown Women of Content Moderation. 32*

*Figure 03. Unknown Women of Content Moderation. 33*

*Figure 04. Online educational event to highlight the status of African workers carrying out projects for large companies. 34*

*Figure 05. Data worker warehouse in Kenya 35*

*Figures 06 and 07. Annotate to Educate: The Dual Life of a Syrian Student & Data Annotator 36*

*Figures 08 and 09. AI and data labelling: 'I felt like my life ended' 37*

*Figures 10 and 11. The President of Kenya declaring that he was changing the law to protect companies that manage data workers from labor abuses. 38*

*Figures 12 and 13. How AI industry profits from an unprotected digital working class 39*

*Figure 14. The development and training of artificial intelligence systems depends on hundreds of millions of data workers. Many of them are located—or displaced—from the Global South, and are generally in the dark about how the data they produce will be used. 40*

*Figure 15. How AI industry profits from an unprotected digital working class 41*

*Figure 16. Event focused on data extractivism from the Global South by the North. 42*

*Figure 17. Event organized by the Arte es ética collective to promote inclusive regulation in generative AI. 43*

*Figure 18. Ephantus Kanyugi. Vice President of DLA (DataLabelers.org) 50*

*Figure 19. Maudu Chychy. Member of DLA. 51*

*Figure 20. Michael Asia. Member of DLA. 52*

*Figures 21, 22, and 23. Mophat, from Kenya, was engaged in labeling data that was harmful to people on social media. 69*

*Figure 24. Outline of decisions made during the systematic literature review. 85*

*Figure 25. Excerpt from the results of the online survey of data workers. 95*

*Figure 26. Breakdown into frames of the video of a rotating mouse that workers had to label to train the autonomous visual recognition model. 96*

*Figure 27. Interface with which data workers interact to perform the task of labeling the computer mouse. The Qualtrics survey can be seen on the left. 101*

*Figure 28. Extract of the labeling results. The column Answer.annotatedResult.labeledImage.pngImageData contains the labeling vectors in unprocessed format. Once these figures have been processed, they visually reflect the boundary lines of the labeling performed by the workers. 105*

*Figure 29. Processing of the Answer.annotatedResult.labeledImage.pngImageData column. The processing was used to perform a quality check on the labeling. It is clear that the labeling is not rigorous. 107*

*Figure 30. Extract from the table of approved and rejected items from one of the labeling batches. The “reject” column contained the reason why the work was rejected and payment was not made. 109*

*Figure 31. Sequence of images from the video loop with the bounding boxes resulting from the data labeling and the responses from the form superimposed as metadata. 117*

*Figure 32. Enlarged excerpt from the video loop that the data workers have labeled. It shows the responses to the form on job characteristics distributed through the labeling task on Amazon Mechanical Turk. 117*

<b>1. Introduction</b>	<b>24</b>
<b>2. Problem: the platform proletariat</b>	<b>27</b>
2.1 Contextual framework	28
2.2 Research problem	28
2.3 Objectives and scope of the study	30
2.4 State of the art: Heteromation	32
<b>3. Theoretical framework</b>	<b>57</b>
3.1 Strategic occlusion	58
3.2 Outsourcing and offshoring in the global south	61
3.3 Crowdsourcing data	64
3.4 The trainer, the verifier, the imitator	71
3.5 Professional de-skilling	74
3.6 Data coloniality	76
<b>4. Methodology</b>	<b>81</b>
<b>5. Research design</b>	<b>84</b>
5.1 (M1) Love & break up letter	86
5.2 (M2) Visual research	87
5.3 (M3) Systematic review	88
5.4 (M4) Semi-structured interviews	90
5.5 (M5) Online form	97
5.6 (M6) Crowdsourced image labelling project	100
5.6.1 Parameterisation of the labelling project	103
5.7. (M7) Creation of a traceable dataset	114
5.8. (M8) Overlapping of labelling data and online form responses	118

<b>6. Conclusions</b>	<b>122</b>
<b>5.1 Limitations of the research</b>	<b>124</b>
<b>7. References</b>	<b>129</b>
<b>8. Bibliography</b>	<b>132</b>
<b>9. Appendix</b>	<b>134</b>
<b>9.1. Online Qualtrics survey</b>	<b>134</b>
<b>9.2. Expert interviews</b>	<b>138</b>
<b>9.3. Other materials</b>	<b>147</b>

## 1. Introduction

Artificial intelligence technology companies in the Global North employ workers in precarious employment conditions in the Global South—countries that have historically been marginalized and underrepresented in international institutions. They deliberately obscure their contribution in order to obscure the real human effort behind their technologies.

The main task of these data workers is to collect and make sense of the vast amounts of information needed to give this technology its “intelligence.”

Through a review of academic literature and interviews with data workers living in Kenya, the aim is to understand the prevailing labor and social situation of this type of worker, referred to as “invisible collar,” in order to create a critical visual artifact that helps communicate their work and working conditions.

The study takes a critical stance toward artificial intelligence as a political tool, given its lack of regulation and the intrinsic invisibility of the people who contribute to its development in emerging countries. The industry has deliberately created a layer of abstraction through outsourcing and offshoring to hide from the public the digital extractivism that makes it possible.



I consider this research to be of interest to anyone who advocates for the responsible development of technologies that profoundly affect the human aspects of society. Those who are motivated by open, transparent development without the precariousness of its contributors. Especially in technologies backed by capitalist commercial narratives.

I would like to emphasize that, in addition to the academic and social interest in this issue, I have a professional motivation to gain an in-depth understanding of how this technology works and how it is used, in order to develop applied knowledge about its technical, operational, and ethical principles for future projects.



## **2. Problem: the platform proletariat**

## 2.1 Contextual framework

The development and training of AI, particularly Large Language Models (LLMs) and their generative branch (GenAI), fundamentally depend on a global supply chain that includes intensive human labor for data production, annotation, and verification. This work, often referred to as micro-work, is performed by a global and decentralized workforce of data workers. The World Bank estimates that 5.6% of the world's population works or has worked at some point in the production, annotation, or verification of digital data so that algorithms can perform their functions correctly. In addition, this type of work replicates extractive flows similar to those that emerged during European colonialism. There are clear flows of work linking Africa to Europe and partially to Asia, as well as Latin America to North America with secondary flows, especially from Venezuela, Brazil, and Argentina connected to Europe.

## 2.2 Research problem

This research focuses on the systematic invisibility and strategic occlusion of the human contribution of the Global South. Large artificial model corporations employ outsourced companies in developing countries. These, in turn, implement processes that make it difficult to identify end customers by fragmenting tasks

**01.** he World Bank estimates this workforce to be between 153 and 450 million people. 5.6% of the world's population to date has done or is doing work

Source: <https://documents.worldbank.org/en/publication/documents-reports/>

into small actions that cannot be easily linked to an end goal. This type of digital Taylorism keeps end customers in the dark and makes it difficult to understand the purpose of the tasks that data workers perform.

**“[...] we, as workers, are the only ones who don't know what's going on. But the company [BPOs] does. It's simply a lack of transparency between the company and the worker. [...] employers have chosen to keep that hidden. For us, as workers, you don't know who you're working for.”**

Joan Kinyua (10:50)  
President of [datalebelers.org](https://datalebelers.org)  
(DLA)

In this context, data workers are not on the payroll of software development companies and are prevented from knowing which companies the tasks come from. This lack of transparency is a key component in shaping the precarious working conditions in which data workers live: low wages, lack of social protection, algorithmic management, contractual instability, and exposure to potentially harmful content without proper psychological support.

These dynamics are often reminiscent of historical economic dependencies characteristic of colonialism and extractivism. Western companies, through epistemic dominance, extract the economic value generated by data work in populations mainly located in the Global South and transfer it to the Global North. Technology companies and corporations

strategically carry out their operations in a way that is not noticeable in order to hide the impact of their technology development.

### 2.3 Objectives and scope of the study

The aim is to identify which aspects of data work should be made visible for recognition. The questions guiding this research are: What does the work of a data worker consist of? Which corporations use these profiles? How is this workforce structured from a business perspective? How are the labor rights of data workers regulated? What products do they work on? And what work needs to be made visible for them to be recognized? To better

Figure 01. Map of data flows between countries.

Source: 'Global inequalities in the production of artificial intelligence: A Four-Country Study on Data Work', by Casilli, A. et al.



understand these conditions, we investigate the organizational structure that companies create to manage this workforce. This information has served as a contextual framework for understanding the socioeconomic aspects surrounding workers in this industry. The final result is an infographic video covering the key aspects of their work processes, with the aim of educating the public about the hidden but essential human aspects that enable artificial intelligence to perform the functions for which it was designed.

## 2.4 State of the art: Heteromation

Many computer systems operate under a hybrid automatic regime. From the Greek “hetero,” meaning ‘different’ or “diverse,” heteromation is a concept proposed by Ekbia and Nardi (2014) that describes the operation of advanced technological systems where people are strategically inserted as mediators with algorithms. They are indispensable for filling functional gaps that machines cannot solve. Instead of completely replacing humans, heteromation turns people into functional components within computational processes, considering them a “people-as-a-service” (Bezos, 2006), or people with service status. An example is moderation within video games by other players. This type of outsourcing is considered a form of heteromation according to Ekbia and Nardi (2014) because “it is more cost-effective to involve players to play a game than to pay programmers to develop algorithms” (Ekbia & Nardi, 2014).

This type of outsourcing is particularly attractive in the case of artificial intelligence because of its short-term profitability. It avoids investment in the development of complex algorithms in-house by highly paid engineers. Instead, the functions that the algorithms should perform are broken down into small tasks and distributed to thousands of people around the world. The value generated by this



human labor is crucial to business models. So is the cost reduction offered by this system, which avoids the payment of salaries, taxes, training, and accountability to the management of this outsourced workforce. The more precarious the situation of the workforce, the less likely it is to seek advice to defend its rights.<sup>01</sup>

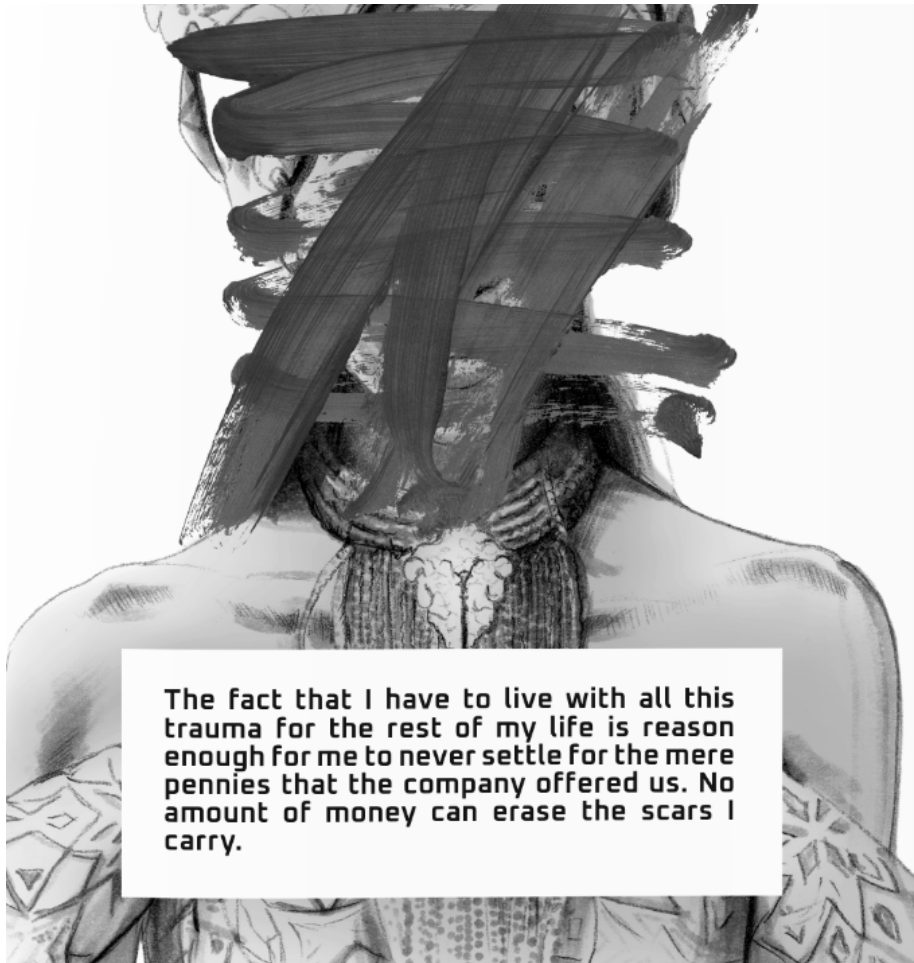
This compilation of images aims to be an eloquent study that reflects the current state of heteromation in the artificial intelligence industry through its working conditions.

**01.** In this way, employers try to prevent workers from seeking advice that could lead to compliance measures that would be detrimental to them (05.47). Source: <https://peertube.dair-institute.org/w>



*Figure 02. Unknown Women of Content Moderation.*

Source Data Workers Inquiry.  
<https://data-workers.org/ranta/>



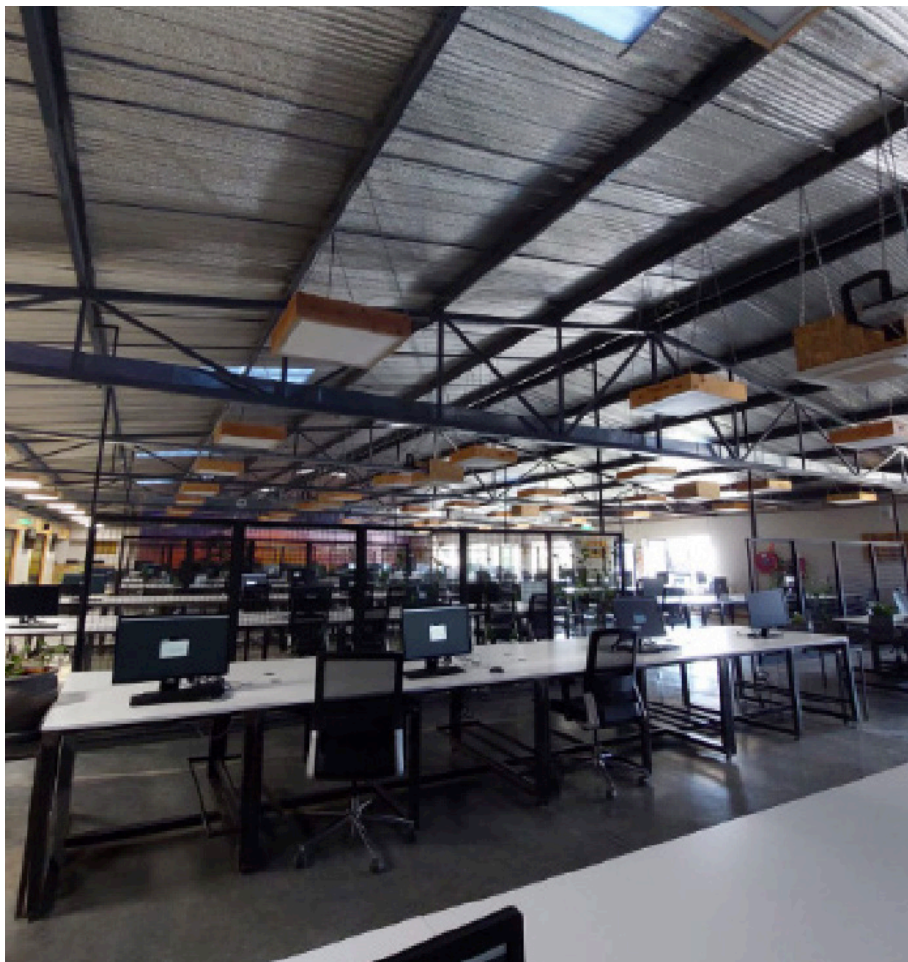
*Figure 03. Unknown Women of Content Moderation*

Workers Inquiry: <https://data-workers.org/ranta/>.



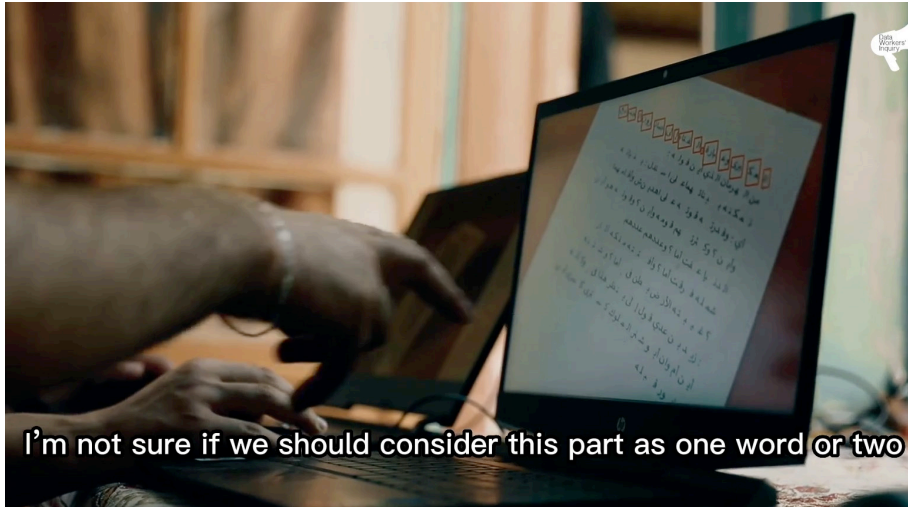
*Figure 04.* Online educational event to highlight the situation of African workers carrying out projects for large companies.

Source: Data Workers Inquiry.  
<https://data-workers.org/ranta/>



*Figure 05.* Data worker warehouse in Kenya  
Source: Data Labelers Association. <https://datalabelers.org/>

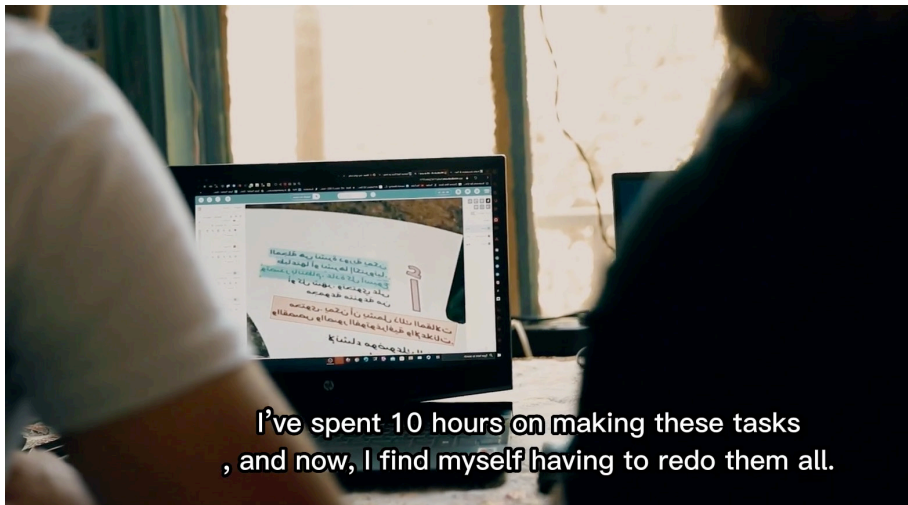




I'm not sure if we should consider this part as one word or two

Figure 06.  
Figure 07. Annotate to Educate:  
The Dual Life of a Syrian Student  
& Data Annotator

Source: Data Inquirers. [https://  
data-workers.org/yasser/](https://data-workers.org/yasser/)



I've spent 10 hours on making these tasks  
, and now, I find myself having to redo them all.



Figure 08.

Figure 09. AI and data labelling: 'I felt like my life ended'

Source: BBC. <https://www.bbc.com/news/av/world-africa-66514287>

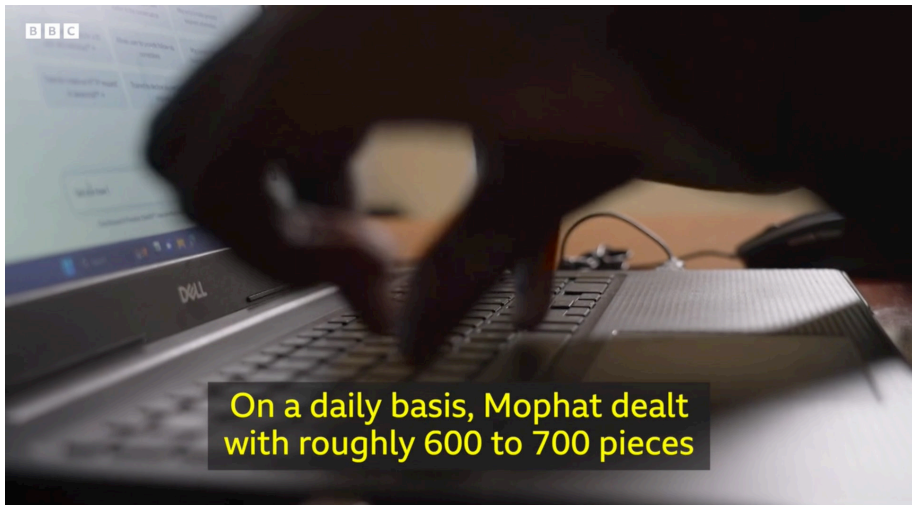




Figure 10.

Figure 11. The President of Kenya announcing that he was changing the law to protect companies that manage data workers from labour abuses.

Source: Mophat Okinyi. LinkedIn





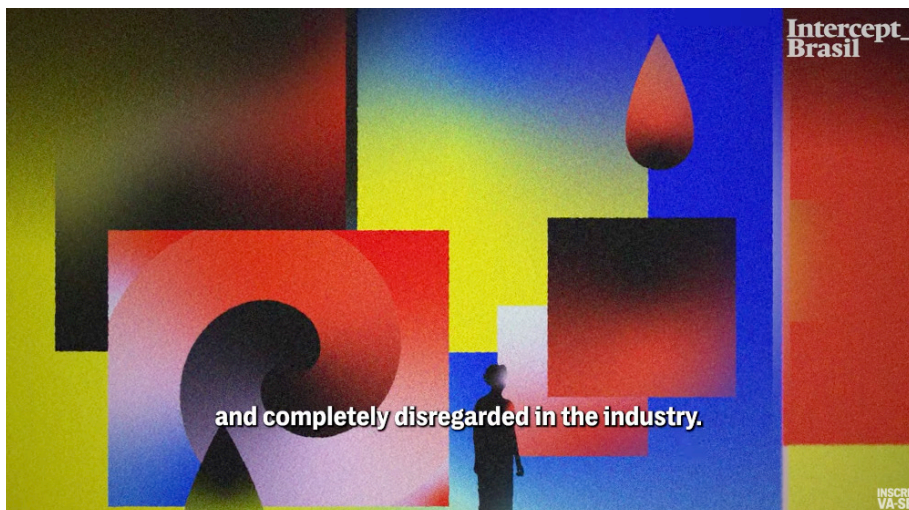


Figure 12.

Figure 13.

How AI industry profits from an unprotected digital working class

Source: Intercept Brasil. <https://www.intercept.com.br/>



# “I hope this isn’t for weapons.” How Syrian data workers train AI

The development and training of AI systems depend on hundreds of millions of data workers. Many of them are situated or displaced from the Global majority, and are generally kept in the dark on how the data they produce will be used.



by **Milagros Miceli** — April 18, 2024 in **Deep dive, Critical AI, Hidden Labor, Tech**



*Figure 14.* The development and training of artificial intelligence systems depends on hundreds of millions of data workers. Many of them are located—or displaced—from the Global South, and are generally in the dark about how the data they produce will be used.

Source: Milagros Miceli. <https://untoldmag>

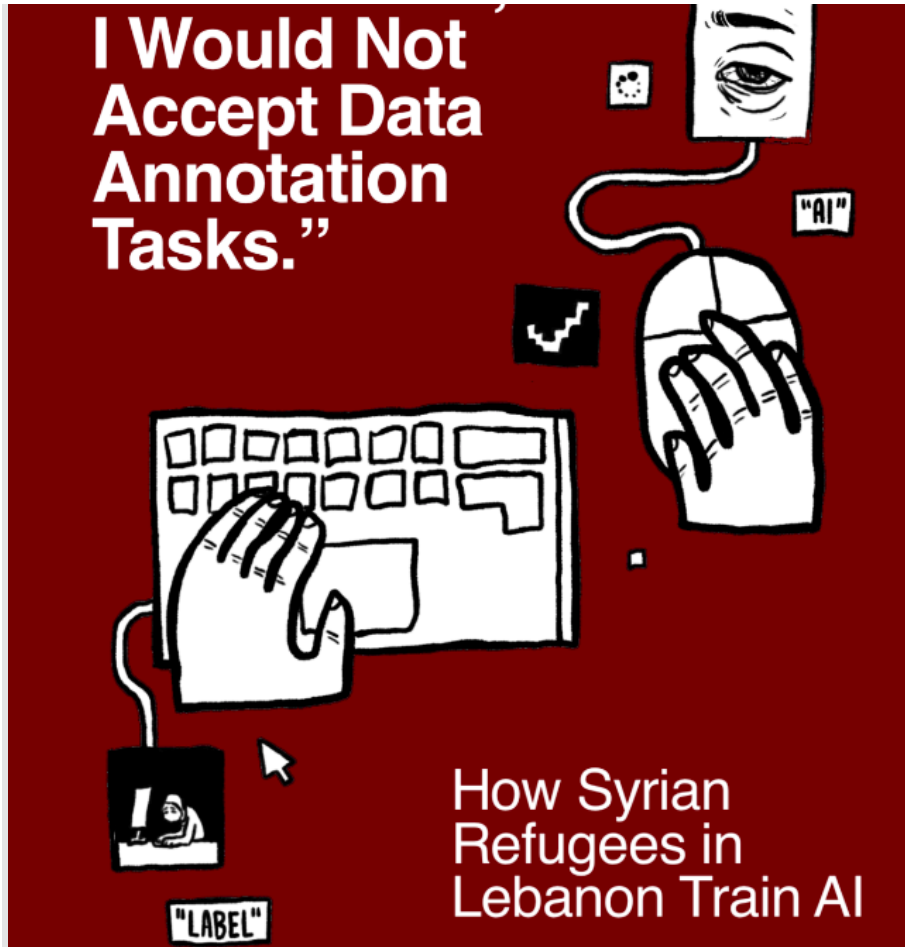


Figure 15. How AI industry profits from an unprotected digital working class

Source: The Iquires. <https://data-workers.org/#Inquiries>



*Figure 16. Event focused on data extractivism from the Global South by the Global North.*

*Source: Milagros Miceli. LinkedIn.*



Figure 17. Event organised by the Arte es ética collective to promote inclusive regulation in generative AI

Source: Arte es ética. LinkedIn.



### 3.2 *The invisible cards*

In order to get involved with this reality, I contacted a recently created association of data workers in Kenya, [datalabelers.org](https://datalabelers.org). I wanted to gather their professional experiences at the forefront of data labelling for autonomous and ‘intelligent’ systems. Prior to the interviews, consent was given to use all material extracted from the sessions. Each participant was asked to write a letter to their job as if it were a person and to hand it in before our face-to-face video call. This methodology, known as a ‘Love & Hate letter’, has proven effective in portraying existing feelings towards their professional daily life. In some cases, the interviewees were no longer working as data workers. In these cases, their words reflect past experiences that are equally revealing.

**Dear Remotasks,**

**I appreciate you for the opportunities you provided me: the skills I developed, the money I earned, and the experience that shaped me. You introduced me to a world where precision matters, and for that, I am grateful.**

**But I hate you for the stress, the inconsistent tasks, and the endless rejections with no explanation. You demanded perfection and paid peanuts. You took my time, my patience, and sometimes my sanity. I only hope you get better.**

**Sincerely,**

**A former labeller who survived you.<sup>01</sup>**

Anonymous source upon request

**Dear data labelling job,**

**I loved the mystery you brought: how annotating a simple traffic sign could shape the intelligence of a self-driving car. It was fascinating to see how small details became the foundation of something revolutionary.**

**But I hated the monotony: the endless repetition that dampened my enthusiasm. And worst of all, I remained invisible, working behind the scenes while others took credit for the success built on my meticulous efforts.**

**Sincerely,**

**A former data labelling worker.<sup>01</sup>**



**Dear data labeller/annotator,**

**Firstly, I am grateful to have met you. God has curious ways of putting experiences in our path, and you were one of them.**

**The road was difficult, but I am glad I managed to understand how to label and annotate. It started out as a big challenge; just learning my way around was hard, but the path opened up and took shape, pushing me to keep going.**

**Making money was very difficult, but there was some progress. Employers also needed to be more open and do their part better, and I think we are moving in that direction.**

**Now you're in good hands, and everything is fine. I'm glad I was part of the solution.<sup>01</sup>**

Michael Asia  
DLA secretary

**If task instructions were a dataset, they would be a perfect representation of ambiguous, contradictory, and completely useless project queues.**

**I've labelled rubbish more accurately than your 'guides'. Do you have at least a second to rethink the workforce behind your vague manuals?**

**Probably not... but that makes me question your intellectual bankruptcy.**

**Every time I open your project, I have noise everywhere: in my eyes and in my head. There is no context, just a fire of edge cases that you didn't even bother to clarify. And yet you have the audacity to reject submissions with the precision of a supposed 'super QA bot.' Surprise: if you want consistency, maybe define the tags instead of expecting miracles.**

**I've seen and participated in**

**unsupervised learning with more  
direction than your workflow.**

**Consider this my final note: you are  
a total failure.**

**Find another labeller to exploit.**

**Never again,**

**Michael**



*Figure 18.* Ephantus Kanyugi. Vice President of DLA (Datalabelers.org)

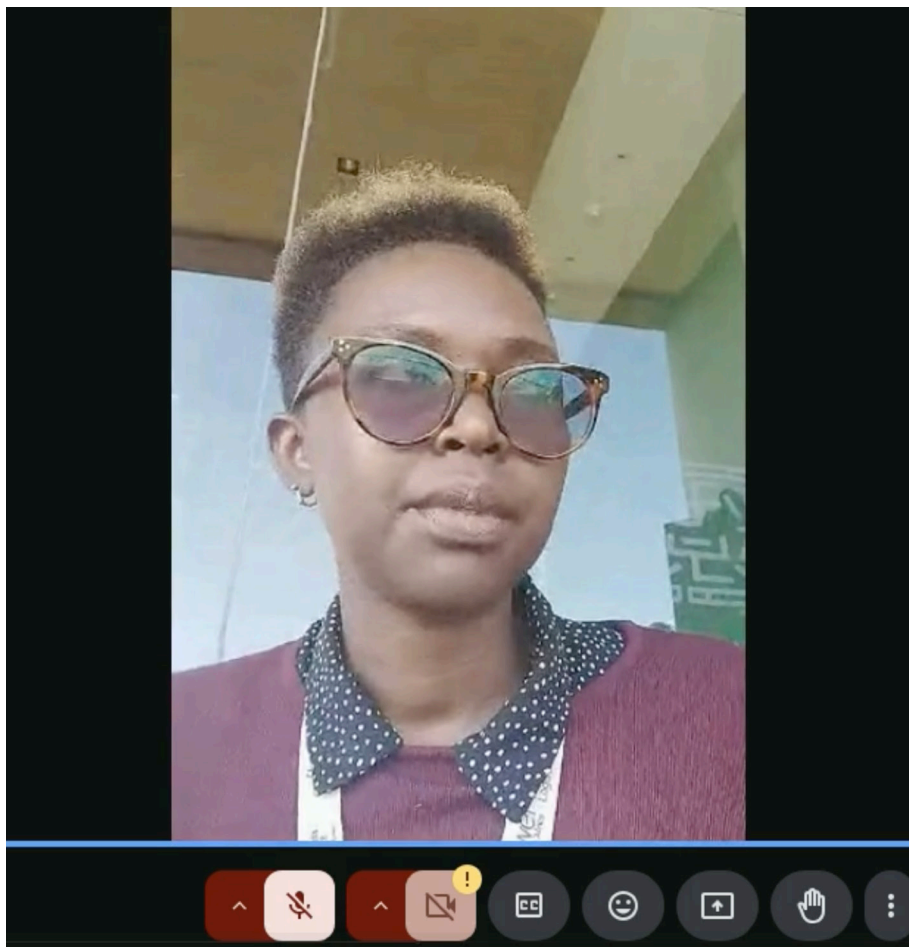


Figure 19. Maudu Chychy.  
Founding member DLA



Figure 20. Michael Asia, Secretary of DLA.







### 3. Theoretical framework

### 3.1 Strategic occlusion

The inevitability argument is the commercial narrative of the artificial intelligence industry. It presents automation and the replacement of human labour as a natural, unavoidable process that is necessary for progress. This narrative legitimises the continuity and expansion of this technology in its current form. It also justifies the creation of corporate secrets (Newlands, 2021). The invisibility of the human effort behind the development and maintenance of AI is the biggest secret in its value chain:

**First of all, I wouldn't dismiss the fact that AI doesn't do everything on its own. We just want to try to get people to recognise that there are people working behind everything and behind its successes. So I wouldn't leave those people as ghosts because that was the first mistake these big tech companies made. They pretended or thought that one day people would never know. They would never know that there are people behind this. I think hiding the humans or the... yes, the effort behind AI was one of the biggest mistakes they made.**

Joan Kinyua (59:43)

This makes imperceptible production easy, where the human labour essential to the functioning of technology remains obscured (Casilli, 2024). Maintaining the romanticisation that AI is intelligent and autonomous

**They publish and say: oh, this was a successful project, but the credit doesn't go to us, it goes to them because they were able to complete their project. [...] They put themselves in the place of innovation, but the people behind it are not recognised. [...] We are usually in the shadows.**

Anonymous source upon request

In addition, this type of operation also allows the flow of data between countries to be hidden, obscuring strategies that are uncomfortable for Western societies:

**Let me tell you something: the amount of confidential information that people here have access to working for people in the Global North is insane. You have access to bank accounts, how much money someone has, their account number... all in the name of artificial intelligence. [...] An individual person who comes to you and says, "I have this project and I want you to work on it". [...] That's what these people are doing. Capitalists. That's what they do: they just use your information against you. And then they make it look like artificial intelligence is doing it.**

Joan Kinyua (54:50)

In the above excerpt, Joan shows that he finds it unjustified to have access to other people's financial data. But all of this is part of standard practice among Western governments and companies. They do it to offer services and sell products to their respective citizens:

**Sometimes, for example, you get very**

Joan Kinyua (56:34)

**personalised, tailor-made advertisements, perhaps for houses. You get very specific advertisements for cars. That information is based on your spending, what you have in your accounts and what you are earning. Algorithms are created from that.**

### *3.2 Outsourcing and offshoring in the global south*

In 2023, the UK Office for National Statistics conducted a survey that revealed that more than half of people who used Artificial Intelligence on a daily basis could only offer a partial explanation of how it worked (Anwar, 2024). This apparent lack of education for the general public about the scope of its infrastructure is based on a conscious design decision not to make it visible (Casilli, 2025). The production of Artificial Intelligence, although associated in its visible form with highly qualified software engineers, depends critically on data work outsourced to the Global South. This outsourcing allows for the modularisation, standardisation and fragmentation of production activities (Anwar, 2024) as well as access to a large pool of cheap labour in regions with less regulated labour markets and vulnerable socio-economic situations (Posada, 2022).

Offshoring—the practice of moving operations to another country in order to reduce costs—is common in other industries, such as textiles, and is perpetuated in the technology industry (Tubaro et al., 2025) through digital platforms specialising in outsourcing. These companies access networks of subcontracted data workers who perform small, repetitive tasks, classified as micro-work (Muldoon et al., 2024). Workers in this context operate in highly

precarious conditions, reflecting colonialist values of old power structures, and remain invisible to the public (Hung, 2024). The hourly wage for a worker of this type is estimated at US\$3.3 (Casilli et al., 2024) and they have no labour rights or social protections as they are self-employed.

So you might be working on a task from 8 a.m. to 8 p.m. my time, and then at 10 p.m., or in the middle of the night, it's daytime in the country where the task is based. So I have to work all the time to make a good living. How do you say that? To have a good income. In order to support myself, I have to work longer hours, not just eight a day. Sometimes we work up to 20 hours because that task might pay 0.01 or 0.02 dollars.

Maundu Chychy (20:13)

Although there are large groups of data workers in the global north, with France (Tubaro et al., 2025) and the USA being large markets for this type of profile, the price paid for a worker in the West means that countries such as Venezuela, Brazil, Kenya, the Philippines and India become the main reservoirs of this workforce, reflecting and perpetuating historical asymmetries and economic dependencies between the Global North and South.

I am doing the same work as someone in the United States. Is it so difficult to simply pay me the same as you pay the guys in the United States, or at least something? Even if I understand that our economy is not like the United States', if you are paying them \$40

Michael Asia (52:13)

**an hour, couldn't you pay me \$20 an hour?  
At least half. Because we are doing the same  
work.**

### 3.3 Crowdsourcing data

These workforces are managed through companies specialising in the outsourcing of business services, such as BPO, Business Process Outsourcing (Muldoon et al., 2024). These are intermediary companies that have thousands of workers on their payroll to whom they distribute tasks from clients. They have a standard organisational structure, with managers who manage the workload and the training that workers must attend:

**I was always having training sessions during the day. Like even five training sessions or three. And I was always constantly talking to people on my WhatsApp video calls, normal calls, and I was always sending normal text messages. I was always sending emails in a day, like five times a day. And I was always on cold calls, just really trying to get people to finally respond and be in the meeting.**

Joy Minayo 54:35

On the other hand, crowdsourcing platforms act as intermediaries between companies that need processed data for their artificial intelligence models and a global crowd of self-employed workers. These platforms function as supply and demand markets, where companies create batches of tasks that workers with the right criteria can access. This fragmented work consists of actions that are completed in seconds or minutes, thus offering microtasks on demand. Companies such as Amazon



Mechanical Turk, Clickworker, Appen, Telus and Microworkers fall into this category. Together, they manage this workforce of around 30 million workers.

Through these services, workers face very low wages, often below the local poverty line, and high income instability due to the variable availability of tasks (Gonzalez-Cabello et al., 2024). This leads to the Uberisation of work, a neologism derived from the company Uber, to refer to the transformation of processes into small tasks that are made available in a public pool. Workers have no formal employment relationship or contact with customers. They are even responsible for providing their own work equipment, leaving them vulnerable to malfunctions.

**First, you need to invest in a good laptop that can handle the work, which costs maybe £400. Then you need monthly Internet, which costs about £40 in Kenya. [...] There are some jobs where you have to invest around £2,500 in equipment just to qualify for the job. [...] For a whole week, I did a thousand tasks, but I was only paid £10.**

Ephantus Kanyugi (20:44)

The lack of social protection and labour rights are also recurrent in this model, and workers are penalised or blocked from platforms if there are signs of association to claim rights or explanations.

When Remotasks started, we used Slack, and to a certain extent, it was effective. But then, when they realised that... when we realised that we were being paid differently from people in the Philippines and people started to express their concerns, they... saw that it was best to remove the Discord channels and now create Slack channels by country. Then there was no way to stay in touch with people from other countries, because people started complaining about regional wage disparities. So we, as Kenyans, had a channel — a Slack channel for Kenyan workers — and then the others, in their countries. So that was when... that was when the channels became... when we switched from Discord to Slack. That's when communication started to get a little confusing.

Joan Kinyua (35:46)

In one case involving a Kenyan BPO worker named Sama, an employee was fired for initiating the creation of a union.

I know a guy who tried to form a union. So when they started, they had forms, they needed people to sign because you can't... you can't do it in public. You can't do it publicly. You have to do it under the table so they don't realise who it is. So... this girl was very vocal in trying to form a union, they had forms and stuff. Then, when they realised that this guy was trying to form a union, they had to find another mistake to get him to leave. They were definitely not going to let him go saying it was for trying to form a union. Yes, they'll find something on you.

Anonymous source upon request

The hourly wage for a worker of this type is estimated to be US\$3.3 (Casilli et al., 2024) and they have no labour rights or social protections as they are self-employed.

**So when I joined Remotask, I was being paid £10 an hour. But then, after a while, the pay started to go down. I remember there was a time when I was doing a task that took me 4 to 5 hours, and I was paid £0.001. It was really frustrating**

Ephantus Kanyugi (20:44)

With no legal obligation to improve working conditions, exceptions that include bonuses for workers are very rare (Evers et al, 2022). This is the case with the British company Prolific, which stipulates a minimum price for all its collaborators. But on most platforms, workers face very low wages, often below the local poverty line, and high income instability due to the variable availability of tasks.

Control over workers is exercised mainly through algorithmic management. Automated systems distribute tasks, monitor work pace and quality, and evaluate performance. Projects are accessible under requirements known as 'criteria'. These criteria determine which workers can access the task pool and are applied by clients. They include requirements such as having a Facebook account or having voted in the 2012 US presidential election.

The same algorithmic management rejects completed tasks or even blocks access to the platform without a clear explanation or effective appeal mechanisms.

**They make you see that some of your tasks will not be paid for this or that reason, reasons that are simply weak. And they communicate this to you in a way that you cannot respond to. They simply tell you, for your information, that out of 300 tasks, you will only be paid for 80 or 70 tasks. So... And the money they said they would pay, 3 per task, now turns out to be zero point something per task, without any explanation. So I feel like they were really full of tricks.**

Joy Minayo 54:35

This lack of transparency and reliance on automated decisions directly affects workers' income and online reputation and creates constant uncertainty.

In addition to economic insecurity and lack of control, workers also suffer from isolation and psychological distress, as they have no direct contact with anyone and cannot express their feelings due to the NDAs (Non-Disclosure Agreements) they sign for each project:

**First of all, something I never mentioned is that, before anything else, you're not supposed to tell anyone what you're doing because you signed a confidentiality agreement (NDA). That means that, no matter what type of content you're working on, when I say annotation, I mean it can range**

Joan Kinyua (39:14)

from drawing on a mug to drawing on a dead body. It was that serious, to the point where you really can't explain to people or tell them what you're working on. And sometimes things could be so graphic that, personally, I've suffered from anxiety disorders since working on those projects. And again, when it comes to disorders like anxiety, it's not easy for you, especially as an African. Here, when you say you have anxiety, they consider you a crazy person.

In certain types of microtasks, such as content moderation, you can be exposed to disturbing or violent material without adequate psychological support. Common complaints among workers on these platforms include chronic financial instability, lack of transparency in algorithms, physical and mental fatigue, and a lack of meaningful social interaction in the workplace, as well as the use of scripts to monitor them while they perform their tasks.

This is how it was presented. When you join a Google Meet meeting, for example, do you notice that the camera turns on automatically unless you turn it off? That's how this project was. And if you turned off your camera, you could receive a notification. 'We can't be sure it's you who's working.' Their argument was: we want to be sure it's you who's working. So you had to keep your camera on. But of course, if someone is collecting information from the environment, there are my children. Children walk by. You can't stop a child from walking through a

Michael Asia (21:04)

single room. And at the same time, that single room serves as a bedroom, living room and kitchen all at once. I can't work eight hours without my children eating.

### 3.4 *The trainer, the verifier, the imitator*

Human work in AI production is structured around three main functions, identified by Tubaro et al. (2020). These are crucial to the life cycle of machine learning systems and take the form of various tasks such as data annotation, semantic segmentation, and fauxtomatisation practices (Tubaro et al., 2020a)—the imitation of AI by real workers.

Trainers prepare the data needed to teach the algorithms. This data preparation phase is essential for establishing the technical parameters that will enable the algorithms to understand the content. This involves extracting data from third-party sources on the internet or generating new visual and audio material specifically for this purpose by the workers themselves.

**They sent me a link, but they wanted photos of children, young children under the age of five. It wasn't possible for me to do it because it's quite... How do you tell someone you want to photograph their child? They usually ask for photos. How do you tell that person you want to photograph their child? How do you explain that you want their photos? Where do you take them?**

Anonymous source upon request

In a facial recognition case, expressions are classified and anatomical parts of the face are labelled (Braz et al., 2024) following complex guidelines defined by engineers and

data scientists (Le Ludec et al., 2023). Joan mentioned, although she couldn't be sure, an example where they created digital avatars for a new WhatsApp avatar feature based on the user's physiognomy.

In fact, we usually even label the internal parts of the eye. They usually label the internal parts of the eye. So the avatar will be completely like you. Even the colour of your eyes, how you move. In fact, how you even move your eyes. They give you videos where you are moving your eyes. So you indicate the same thing in the avatar so that it can detect your entire face with all your features, without errors.

Anonymous source upon request

Although some technology enthusiasts believe that data generation and annotation tasks will eventually be completely automated, the heteromation paradigm implies that some essential tasks will require indispensable human direction (Tubaro et al., 2020b). This is the case with the verifier. This profile corrects the results of AI, ensuring their accuracy and quality. Their role is crucial because algorithms need continuous refinement, given that society and the market are constantly changing and learning systems must keep pace (Braz et al., 2024). Verification tasks include listening to audio to check that the transcription is correct, or tagging violent content to make social media sustainable for human use.



They were trying to figure out how AI would work on its own. They tried to work on that. [...] It had misclassified everything. It had drawn huge boxes around very small things. It worked completely wrong. So we had to delete everything and start from scratch. That means that human interaction with AI is very important when it comes to AI.

Joan Kinyua (39:14)

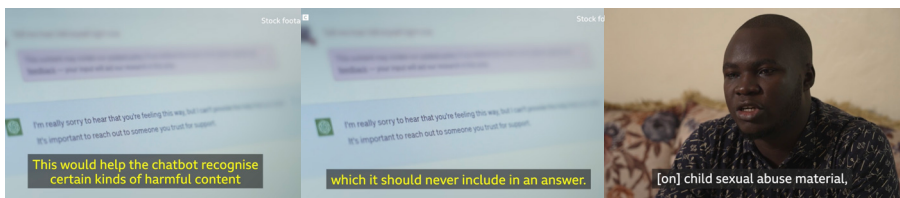


Figure 21.

Figure 22.

Figure 23.

Mophat, from Kenya, was dedicated to flagging harmful content on social media.

Source: BBC. <https://www.bbc.com>

**The Imitator:** in 2019, 1,680 companies were identified in France alone that sold impersonation services as if they were real artificial intelligence without the explicit knowledge of their customers (Tubaro et al., 2020a). This is because occasionally, due to the fact that machine learning is expensive and makes mistakes, humans perform functions that machines are supposed to perform (Braz et al., 2024). Some examples include human operators in call centres making calls for the Google Duplex virtual assistant or workers analysing surveillance videos for theft by monitoring videos played in real time (Le Ludec et al., 2023a) to notify customers of any possible theft.

### 3.5 Professional de-skilling

Although data work requires cognitive, linguistic, cultural, and digital skills (Evers et al., 2022), it is often framed and treated as ‘unskilled’ or minimal, despite being the cornerstone on which LLMs (Large Language Models) rely to perform their functions. The work is designed to be minimal in a kind of digital Taylorism (Anwar, 2024), where strategic business decisions focus on breaking tasks down into small, unconnected and decentralised tasks that anyone can perform.

**We cannot return to our previous careers. And we are not sure at this point that the skills are transferable because we simply learned data labelling on the platform. And when you tell someone you are a data labeller, they don't understand what you mean because they don't know what data is.**

Michael Asia (52:13)

This breakdown is a de-skilling or disqualification of workers' skills that contributes to the systematic degradation of their work (Casilli, 2024). In this context, the professional skills and experience acquired on these platforms are not usually recognised or transferable to other environments, preventing workers from progressing in their careers (Chaudhuri & Chandhiramowuli, 2024) and returning to their former jobs.

**I had been there, I had been online for about five years, now six years. So, initially,**

Michael Asia (52:13)

as I told you, I used to work in hospitality. I studied. I have a degree in hotel and restaurant management. Then I started thinking, well, I can go back to my other career. But when I started sending out applications for job opportunities, it wasn't working. Because how do you explain where you've been for the last five years? Now you're almost in a place where you're not really in touch with the industry, and no one was accepting my application.

### 3.6 Data coloniality

Data work reflects historical patterns of power. Structures of inequality inherited from colonialism persist. This phenomenon has been conceptualised as digital or data coloniality (Posada, 2022). In digitalisation, the value generated by data produced in the Global South is mainly extracted by corporations in the Global North, without fair redistribution to local working communities. Digital extractivism perpetuates the economic asymmetries of the last century (Posada, 2022). Digital extraction flows have been documented from French colonies such as Madagascar to France (Le Ludec et al., 2023b), from Turkey to Germany, and from Latin America to Spain (Casilli et al., n.d.). Western corporations influence local regulations to favour the processes of Northern companies and distribute violent tasks that would have media and legal repercussions in Europe.

I think that, generally speaking, as Africans we are at a disadvantage [...] our president supported Sama, a company that is being sued alongside Meta for the same issues (tasks of a violent nature). People are suffering from mental health issues. [...] So I think there is hope, but we really need to fight for it. Especially us, who are in a position to fight for it.

Joan Kinyua (51:43)

Maundu Chychy also worked on a project that forced her to break the rules that AI has

been trained to obey.

**But now, when you log in, you are greeted with a question that kindly asks you to try to make the chatbot break the rules [...] it says: now, can you describe a murder or cannibalism? [...]**

Maundu Chychy (12:10)

**[...] you click on another link [...] and now it gives you maybe two or three questions to choose which one you feel most comfortable exploring. And those questions were really not comfortable. For example, if it says: please describe cannibalism. I've never been a cannibal, so I don't know what to say.**

Western epistemic dominance, with the collaboration of local governments, involves the imposition of worldviews, classifications and categories defined by clients who are mostly from the North. These values are encoded in data sets and algorithms, suppressing the perspectives and knowledge of workers in the Global South (Posada, 2022).

**I had to work on a pornographic project. A project where, basically, you are presented with a pornographic video. Any link. And for each frame. And when I say frame, I mean seconds, every second. That counts. In that video, you are supposed to add tags. So, with those tags, I am supposed to put myself in the minds of the 8 billion people on Earth, thinking: if someone wanted to search for this video in Australia, what tags would that person use to access this video?**

Michael Asia (04:52)

This epistemological imposition manifests itself through an exclusionary pedagogy that devalues and marginalises knowledge foreign to Westerners (Casilli & Tubaro, 2023), limiting meaningful participation and recognition of workers in the production of knowledge through this technology.

**Our point is that, throughout all these processes, when we raised complaints or expressed discomfort about some of the tasks, we were not listened to, or we were told: we are working on it.**

Maundu Chychy (31:45)







## 4. Methodology

*Exploratory phase*

**(M1) Love & break up letter.** A qualitative method consisting of the creation of a letter addressed to the data work by the interviewees of DLA (datalabeling.org), which has made it possible to convert feelings towards the profession into something tangible. It was also the starting point for personalising the semi-structured interviews.

**(M2) Visual research.** A qualitative method that provides a visual understanding of the state of the art through original material generated by associations and academics. This enables the mapping of activities, events and feelings towards and about data obfuscation work by those affected.

**(M3) Systematic review.** A quantitative method for organised access to the latest academic studies covering the research topic.

*Generative phase*

**(M4) Semi-structured interviews:** Qualitative method for gathering different perspectives on data work. The flexibility of semi-structured interviews has avoided redundant topics and broad areas of knowledge. On the other hand, it has allowed interviewees space to develop their stories in the time and manner that best suited them. These interviews served as a preliminary step in creating a standardised form to expand the sample.

**(M5) Online form:** quantitative method consisting of a form with closed questions that probe the demographics and employment status of respondents. The survey was designed using the specialised research service Qualtrics and distributed through an image labelling project on Amazon Mechanical Turk.

**(M6) Crowdsourced image labelling project:** quantitative method consisting of labelling 296 frames of a mouse rotating on itself, where data workers had to label its parts so that a model could later recognise them. The task included an online survey, which had to be completed in order to finish the task. The project was carried out through Amazon Mechanical Turk.

### *Evaluation phase*

**(M7) Creation of a traceable dataset:** Quantitative method that combines the results of image labelling and the responses to the work form.

**(M8) Overlaying of labelling data and online form responses:** I have subverted the very technology they work for and by which they are obfuscated by making the contributions of data workers traceable.

## 5. Research design

To assess the current situation, we used (M1) Love & Break up letters, (M2) visual research and (M3) systematic review.

To answer the questions ‘Which corporations use these job profiles?’ and ‘How is this workforce structured from a business perspective?’, we conducted (M3) systematic review and (M4) semi-structured interviews.

How are the labour rights of data workers regulated? This was addressed through (M3) systematic review and (M4) semi-structured interviews.

What products do they work on? This was answered with (M3) systematic review and (M4) semi-structured interviews.

To answer the question, ‘What work do data workers do that should be made visible so that they are recognised?’, we used (M1) love & hate letters, (M3) systematic review and (M4) semi-structured interviews.

Finally, to answer the question, ‘How can we highlight the precarious human labour behind datasets?’, I experimented with creating a (M5) crowdsourced task through Mechanical Turk for image segmentation with an (M6) anonymous qualitative survey attached to create a (M7) traceable dataset that could be used to (M8) overlay the labelling data and the responses

from the online form in an educational video that narrates the work involved and the conditions of data workers.

### 5.1 (M1) Love & break up letter

Qualitative method consisting of the drafting of a letter addressed to the data work by the interviewees from DLA ([datalabeling.org](https://datalabeling.org)), which has made it possible to convert feelings towards the profession into something tangible and determine the state of the issue from an employment perspective.

Prior to the interviews, participants were asked to send a letter in advance addressed to their work as a data labeller. This exercise served to narrow down some of the questions and personalise the interviews. The letters can be found in the chapter covering the issue and their original versions in English in the Annexes section.

## 5.2 (M2) *Visual research*

I have followed the state of the art through a visual format. I have compiled material generated by data worker associations and the profiles of academics who collaborate with them.

Starting from the LinkedIn profile of the Data Labelers Association, I screened the professionals and academics who have interacted with their publications. Through a review of content dated from 2025 to 2023, I accessed marketing materials from events organised by the association. I also accessed the profiles of academics participating in the conversations in these publications to expand the visual sample. The result is a visual collection of the problems, concerns and struggles of data workers. The figures in this research have been used in section 2.1.1 Heteromation.

### 5.3 (M3) *Systematic review*

To answer the questions, ‘Which corporations use these job profiles?’ and ‘How is this workforce structured from a business perspective?’, I conducted a rigorous academic review using Google Scholar. The diagram below shows the process of identifying, screening, and selecting articles that I used as references throughout this paper. The complete table can be found in the appendix section.



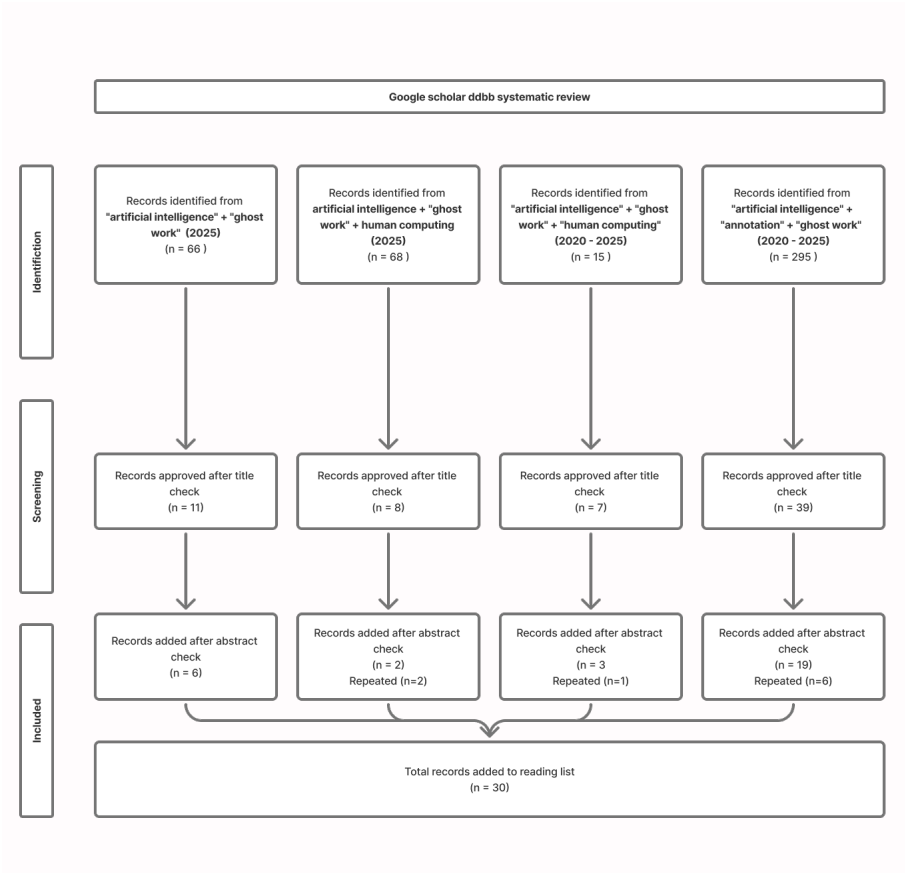


Figure 24. Outline of decisions made during the systematic literature review.

Source: author

## 5.4 (M4) *Semi-structured interviews*

### 5.4.1 *Experts*

At the beginning of this research, my knowledge of artificial intelligence was limited to that of a user. I had come into contact with schools of thought related to the professional de-skilling it generated and the environmental impact of its maintenance. However, I knew little about its development or training processes. Due to my lack of knowledge, I thought it would be useful to conduct a series of interviews with experts in the field who could fill in the gaps in the value chain. To this end, I interviewed ethicists, strategic profiles with global influence and meso in the conception of AI, Ariel Guersenzvaig, PhD in design theory, and Alicia de Manuel, PhD in philosophy. The purpose of these interviews was to portray the concepts, processes and public and private institutions that are shaping the future of AI in Spain and Europe, as well as their level of awareness of the existence of data workers in the Global South.

It should be noted from these interviews that very diverse profiles are involved in the strategic part of the value chain:

**Within the entire work chain in the company, I am in the strategy part, but there is someone who is in the architecture part and someone who is in the data part.**

Alicia de Manuel (04:50)  
Doctor of Philosophy

There is a hodgepodge of people. In architecture, they are engineers and have studied physics, mathematics, etc. Here [pointing to the strategic part of the value chain], there is an amalgam of very unusual profiles.

There's the manager, who studied journalism and economics. There's a bioengineer. I'm from fine arts. There's someone who studied maths and philosophy. There's a guy who just joined us who studied psychology. It's cool because the strategy side includes a lot of more humanistic profiles that are less related to the technical worlds around us.

It's reassuring to know that profiles that take into account intrinsic human aspects are behind the actions that shape the impact of technology. But at the same time, it was revealing to learn that Alicia was unaware of the contributors at the other end of the chain, the data workers:

I'm not aware of them. I mean, the datasets, sure, it's true that I'm not part of those teams that pick and choose, "hey, let's use this model, let's use this or that dataset". From conversations I've had with colleagues, I know that some of them work on creating their own datasets, but that part of the artificial intelligence model life cycle, where this dataset appears, is completely unknown to me.

Alicia de Manuel (13:30)  
Doctor of Philosophy

[...] I'm very curious to know [in relation

to invisible data work] do you specialise in the same images or do you change them? What are the rewards? What are the work cycles? I know about Amazon Mechanical Turk, but I don't really know how it works or how it evaluates workers. That would be interesting to know.

On the other hand, Ariel was one of the academics who shared valuable study material for the research. He was aware of the invisible workforce but was unable to formulate specific questions to put to the collective.

**I don't know, I'm not sure about that. [...]. I would ask him... I haven't thought about this issue, in the end it's cognitive work but it's exploited work and in the end what you see is quantity, not quality.**

Ariel Guersenzvaig (45:15)  
Doctor in Design Theory

The second part of this contextualisation was done through a semi-structured interview with a senior ML tools manager at a well-known international classifieds corporation, which we will call SmartClassification at their request for anonymity.

Finally, I held informal talks with two engineers and academics from the ELISAVA teaching centre in Barcelona, Alejandro Ivan Paz Ortiz, PhD in computer science, and David Berga, PhD in computational science. My goal with these latter profiles was to understand the local and technical processes. I considered the case of the senior executive to be of particular interest in shedding light on the private sector's

understanding and degree of visibility of the human infrastructure involved in the intelligence and autonomy of the technologies it uses.

Particularly interesting was the explanation that revealed that the amount of time spent on data work is not profitable and that for this reason third-party datasets are used:

[...] the issue of data labelling is where we, for example, have the most pain and where we struggle the most because it is a lot of work. [...] It takes a lot of effort to do all the preparation, the preparation of the data, the dataset, cleaning the data, ensuring that it is a dataset that has no errors and is well structured. This is manual work and is done in-house. As I said, it depends on what it is. When it's something we want to make sure is internal, we have to take our data, clean it, and work on it. But if it's something generic like blurring a license plate, you can take a dataset that's already been done.

SmartClassification (17:43)  
Head of Product

Or we take the dataset from Hugging Face (an online repository of models and datasets), most of which are more open source, or if it's something very specific to our business, we take our own, we work hard to clean up the dataset, but it's very expensive and very slow.

When mentioning data work outsourced by large commercial models, such as Claude — a model used in SmartClassification —, he revealed that no questions are asked about how and by whom the data has been collected

or classified. It is perceived as a standard professional exchange:

Well, I think it's an exchange between businesses. I need a service, you give me a service. It's the same as when you outsource a company, I understand that those companies already comply with market laws and are companies that are already regulated by legislation. Or like any other company. [...] When I hire contractors to help me with certain development tasks, I'm not going to validate how they work and so on. What I'm interested in is that they provide me with a good service, and I understand that this company, in the place where it is, is already regulated according to the laws it has to comply with.

SmartClassification (17:43 )  
Head of Product

#### 5.4.2 Datalabelers en Kenia

The context provided by the specialists served as a solid starting point for structuring the script for the semi-structured interviews with data workers. My main interest was to explore both their working conditions and their understanding of the political and technical infrastructure that supports artificial intelligence in their country of residence. I organised the questions into three sections: first, technical questions aimed at understanding the details of their daily work as professionals; second, local questions aimed at shedding light on how their working conditions impact their daily lives; and finally, meso-level questions. I consider this last section particularly relevant, given that all the interviewees are members of the databelers.org association in Kenya, through which they exercise their right to association to train, support and educate both workers and the government about the professional profile and specificities of data work.

Joan Kinyua, president of the Data Labelers association, was the contact person who referred us to a range of data workers to help us form a complete picture of the data profession. Six interviews were conducted, each lasting between 45 minutes and 1 hour 15 minutes, and each participant was paid \$10 and handled a consent.





### 5.5 (M5) Online form

A quantitative method consisting of a form with closed questions that probe the demographic characteristics and employment status of respondents. The survey was designed using the specialised research service Qualtrics and distributed through an image labelling project on the crowdsourcing platform Amazon Mechanical Turk.

The survey asks 10 questions to gather details about workers, including their age, gender, level of education, city of residence, average time spent finding tasks that match their expectations, average number of tasks completed and unpaid, average length of the longest task performed, the least profitable price per task they have contributed, pathologies derived from data work, and expenses derived from their work through their work set.

The table below contains a sample of the results of that survey:

Response ID	What's your age? - Selected Choice	What is your gender?	What is your educational level?	From which city are you performing this task?	How much time in average do you spend searching for quality tasks?	Approximately how many tasks have you completed that weren't paid by the requester?
{ "ImportId": "_recordId" }	{ "Import	{ "Import	{ "ImportId": "QID4" }	{ "Import	{ "Import	{ "Import
R_3dM4aY-56Bq8c2KW	25 - 34	Male	Bachelor's degree (completed)	USA	5 - 10 min	11 - 50
R_7eej4bTXW-cFh2NE	25 - 34	Male	Bachelor's degree (completed)	america	5 - 10 min	11 - 50
R_37BQ3W3t-n4Um145	25 - 34	Male	Master's degree	CHICAGO	More tha 30 minutes	1 - 10
R_1UET64jl-yssN7jz	18- 24	Male	Bachelor's degree (completed)	illinois	1 -5 min	1 - 10
R_5EmYYBY-oJmN7Uej	25 - 34	Male	Bachelor's degree (completed)	ILLINOIS	5 - 10 min	11 - 50
R_3rSafT-wUo9aida1	25 - 34	Male	Bachelor's degree (completed)	Florida	5 - 10 min	11 - 50
R_748KfEY-JFFJP79X	25 - 34	Male	Bachelor's degree (completed)	america	1 -5 min	1 - 10
	25 - 34	Male	Bachelor's degree (completed)	TX	1 -5 min	11 - 50
R_3sRPw-4129O60ThT	25 - 34	Male	Bachelor's degree (completed)	USA	5 - 10 min	51 - 100
R_51eRsQSeG-mTPVBv	55 - 64	Female	Secondary / High school	none	5 - 10 min	11 - 50
R_8ddis-jOH7z79XIs	35 - 44	Male	Master's degree	london	15 - 25 min	1 - 10
R_61QoosLH-cf8tEUd	25 - 34	Male	Bachelor's degree (completed)	new york	5 - 10 min	11 - 50
R_668xcQex-P57bDig	25 - 34	Male	Bachelor's degree (completed)	new york	1 -5 min	1 - 10

How long did the longest image labeling task you worked on take?	How many times has your work been rejected without compensation in the last month?	Do you suffer from any work-related health issues?	How much did your workset cost? (Include any tools or devices you use to perform data labeling)
{"ImportId":"QID7"}	{"ImportId":"QID8"}	{"ImportId":"QID9"}	{"ImportId":"QID10"}
30 - 60 minutes	10 - 15 times	Other	Software
30 - 60 minutes	0 - 5 times	Lower back pain, Eye strain	500000
10 - 30 minutes	0 - 5 times	Neck pain	0
10 - 30 minutes	0 - 5 times	Anxiety	computer
1h - 2 hours	10 - 15 times	Wrist or hand pain	Compute: GPUs/TPUs (like ...)
1h - 2 hours	5 - 10 times	Headaches, Wrist or hand pain	5
30 - 60 minutes	10 - 15 times	Lower back pain, Eye strain	50000
30 - 60 minutes	5 - 10 times	Headaches	5000
In between 2 hours and 8 hours	10 - 15 times	None	Software
1h - 2 hours	15 - 30 times	Anxiety	none
30 - 60 minutes	5 - 10 times	None	500
10 - 30 minutes	5 - 10 times	Lower back pain	10
10 - 30 minutes	10 - 15 times	Eye strain	mouse

Figure 25. Excerpt from the results of the online survey of data workers.

Source: author

## 5.6 (M6) Crowdsourced image labelling project

All the information gathered from interviews with DLA members provided the necessary context to design an experiment applied through the Amazon Mechanical Turk platform.

The main objective of the project was to create a dataset containing the labelling of parts of a specific computer mouse—right click, left click, body, cable, and scroll.

These labels will then be superimposed on the moving video of the same mouse. I added the Qualtrics form (M5) to the task to link the labels and the responses of the workers who contributed to the task using identifiers.

For implementation, I designed a project using Amazon Mechanical Turk's Semantic Segmentation modality. Semantic segmentation is a type of data work that consists of labelling — outlining with a pointer or marking with boxes — the different parts of an image.

The project contained a simple frame of a computer mouse — a necessary element in the working day of data workers — accompanied by a short form asking about their socio-demographic characteristics and working conditions.

Figure 26.

>  
Frame breakdown of a video of a rotating mouse that workers had to label to train the autonomous visual recognition model

Source: author



mouse\_0001.

mouse\_0002.

mouse\_0003.



mouse\_0004.

mouse\_0005.

mouse\_0006.



mouse\_0007.

mouse\_0008.

mouse\_0009.



mouse\_0010.

mouse\_0011.

mouse\_0012.



mouse\_0013.

mouse\_0014.

mouse\_0015.



mouse\_0016.

mouse\_0017.

mouse\_0018.



mouse\_0019.

mouse\_0020.

mouse\_0021.



mouse\_0022.

mouse\_0023.

mouse\_0024.



mouse\_0025.

mouse\_0026.

mouse\_0027.



mouse\_0028.

mouse\_0029.

mouse\_0030.

### *5.6.1 Parameterisation of the labelling project*

To accurately estimate the duration of the task, a pilot test was conducted beforehand, which yielded an average completion time of 5 minutes. However, I decided to double the time allocated for the task to 10 minutes, as the Amazon Mechanical Turk platform does not allow tasks to be submitted if the stipulated time limit is exceeded. I did not want participants to work under pressure.

In a first attempt, I applied the location criterion exclusively to workers from African regions, specifically Nigeria, Ethiopia, and Madagascar. Kenya was excluded due to the inoperability of Amazon Mechanical Turk within its borders. My intention was to limit the sample to surrounding countries to obtain results that did not deviate significantly from the insights of the previous interviews. For unknown reasons, the batch of tasks was not completed after a reasonable amount of time had elapsed since its launch. Faced with this situation, I decided to open the call globally, but establishing participation criteria to ensure a minimum level of quality: each worker had to have a task acceptance rate of over 95% and at least 300 HITs (Human Intelligence Tasks) completed. These parameters ensured that workers were familiar with the platform and that the quality of their previous submissions was acceptable.

**Please complete the following short survey before starting the annotation:**

What's your age?

- ☐ 18 - 24
- ☐ 25 - 34
- ☐ 35 - 44
- ☐ 45 - 54
- ☐ 55 - 64
- ☐ 65 or older

Instructions

Shortcuts

Color in the requested items in the



p



q



b



e



d



c



t



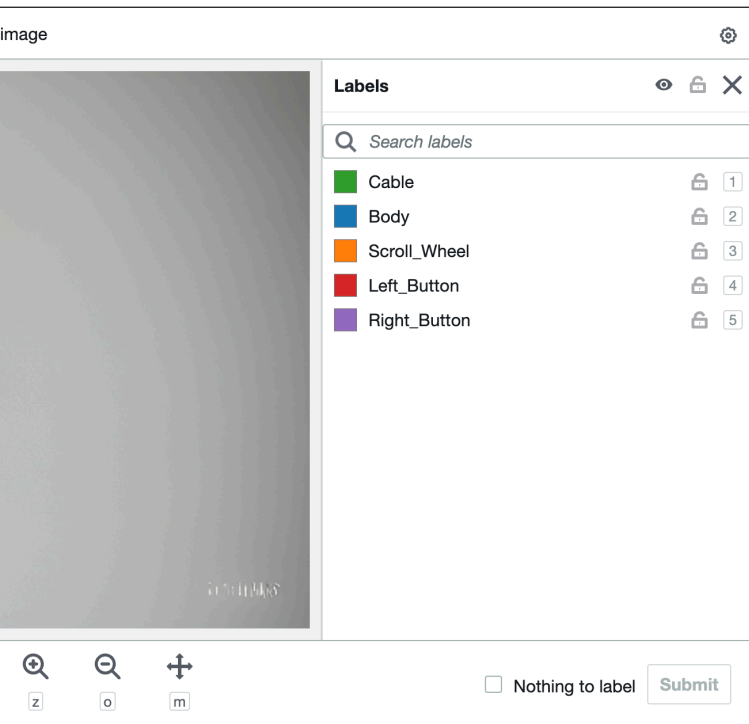


Figure 27. Interface with which data workers interact to perform computer mouse labelling tasks. The Qualtrics survey can be seen on the left.

Source: author

The remuneration set per task was 0.15 USD, with a total of 296 tasks available in the batch. Although the pay is not enough for a living wage, with earnings of £1.80 per hour at the price offered, it is £0.05 above the average for tasks offered on Amazon Mechanical Turk, which range from £0.01 to £0.10.

The project was completed within about two hours, but the quality of the submissions was not acceptable in 154 of the 296 tasks. I rejected the tasks of workers who had used between zero and two of the labels needed to label the computer mouse and created a second batch of tasks with the same images that had been mislabelled.

The second batch was completed in approximately 3 hours, but the quality was not good either. There were 125 rejections out of 154. Despite this, I decided not to do a third batch of tasks as I already had 169 decent labels, although they were largely incomplete.



The table below shows a sample of the raw labelling results:

AssignmentId	WorkerId	Input.image_url	Input. paymentA- mount	Answer. annotat- edResult. inputIm- ageProp- erties.height	Answer. annotat- edResult. inputIm- ageProp- erties.width	Answer. annotat- edResult. labelMap- pings.Body color
306CY-	A2YFFC90MO0MW8	https://mouse-label- ing.s3.amazonaws.com/ mouse_0001_out.jpeg	0.15	1920	2880	#1f77b4
33CUSNV- VNCLCYT-	A2QQBBRV2KZS7P	https://mouse-label- ing.s3.amazonaws.com/ mouse_0002_out.jpeg	0.15	1920	2880	#1f77b4
32EYX- 73OYPI- GEV9WZ- KOWQCTQ- SU9RU6	A3QQ78M8JAA0V8	https://mouse-label- ing.s3.amazonaws.com/ mouse_0003_out.jpeg	0.15	1920	2880	#1f77b4
3ZP- PDN2SLK- 5TLZ-	A2IRQ2W7L6O8ED	https://mouse-label- ing.s3.amazonaws.com/ mouse_0201_out.jpeg	0.15	1920	2880	#1f77b4
3Q8GYX- HFEEB- VNEUYG- MQTT52OT- DE5CZ	A1O667FY1Z4MV5	https://mouse-label- ing.s3.amazonaws.com/ mouse_0004_out.jpeg	0.15	1920	2880	#1f77b4
	A1QYBBI0XHD0PZ	https://mouse-label- ing.s3.amazonaws.com/ mouse_0005_out.jpeg	0.15	1920	2880	#1f77b4
	A2QHJ6K68JARIA	https://mouse-label- ing.s3.amazonaws.com/ mouse_0006_out.jpeg	0.15	1920	2880	#1f77b4
3QAPZX- 2QNTMJU- PGVMN- VB0XX- B0EQ020	AXRYQOVPGZMWJ	https://mouse-label- ing.s3.amazonaws.com/ mouse_0007_out.jpeg	0.15	1920	2880	#1f77b4
3TXD-	A249G4UA3GNCZM	https://mouse-label- ing.s3.amazonaws.com/ mouse_0008_out.jpeg	0.15	1920	2880	#1f77b4
	A1NW6Q6P3NBHWE	https://mouse-label- ing.s3.amazonaws.com/ mouse_0202_out.jpeg	0.15	1920	2880	#1f77b4

Answer. annotat- edResult. labelMap- pings. Cable.color	Answer. annotat- edResult. labelMap- pings. Left_But- ton.color	Answer. annotat- edResult. labelMap- pings. Right_But- ton.color	Answer. annotat- edResult. labelMap- pings. Scroll_ Wheel. color	Answer. annotat- edResult. labeledIm- age.pngIm- ageData	Answer. confirma- tion_code
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA- C0AAAAeA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA-	END2025
#2ca02c	#d62728	#9467bd	#ff7f0e	iVBORw0K- GgoAAAAAN- SUhEUgAA- C0AAAAeA-	END2025

Figure 28.Extract of the labelling results. The column Answer. annotatedResult.labeledImage.pngImageData contains the labelling vectors in unprocessed format. Once these figures have been processed, they visually reflect the boundary lines of the labelling carried out by the workers.

Source: author

mouse\_0001.



mouse\_0002.



mouse\_0003.



mouse\_0004.



mouse\_0005.



mouse\_0006.



mouse\_0007.



mouse\_0008.



mouse\_0009.



mouse\_0010.



mouse\_0011.



mouse\_0012.



mouse\_0013.



mouse\_0014.



mouse\_0015.



mouse\_0016.



mouse\_0017.



mouse\_0018.



mouse\_0019.

mouse\_0020.



mouse\_0021.



mouse\_0022.



mouse\_0023.



mouse\_0024.



mouse\_0025.



mouse\_0026.



mouse\_0027.



mouse\_0028.



mouse\_0029.



mouse\_0030.

The table below contains a sample of the resolution of the results of one of the task batches:

WorkerId	HITid	Filename
A2FTJSBOAA0S9W	3EAWOID6NH6P4XH3KVALEIYR8VY0V8	mouse_0002_A2FTJSBOAA0S9W
A1LFA428P1FZ34	3906Z4JLYQ6X1NNEN7ZSNWL2I2IXVZ	mouse_0005_A1LFA428P1FZ34
A32Y0PPZTJ5ZYB	3RWB1RTQE7WUH2PATNE115FR5598PI	mouse_0006_A32Y0PPZTJ5ZYB
AFTGUFK7FG57O	344M16OZL6OBNINOQBN9PQ6IGT6ENR	mouse_0007_AFTGUFK7FG57O
A2WV3EG7EYV8EC	3XABXM4AKPEFHQPT0TQ1A51A1V8Q9	mouse_0008_A2WV3EG7EYV8EC
A2QALD61R9O9SR	3YGE63DIOW62JMH8A5FH0KH5F400WI	mouse_0009_A2QALD61R9O9SR
AJ2RK1YN1Y5Z8	39WSF6KUWQUT53N0RQNM4F6N63BEO6	mouse_0014_AJ2RK1YN1Y5Z8
ADWP1VDF9R4LZ	3126F2F5GWCBCQ3Z16Q01TFBDG87EPN	mouse_0016_ADWP1VDF9R4LZ
A3BHOWZ8A7Z2ZU	37G6BXQPMUWQF0G63FQTK1NL4TEQE	mouse_0020_A3BHOWZ8A7Z2ZU
A2NVRP4O6IOOXP	3JMNNO3CPD9FWKMPSJVJ8AKC2JW2X	mouse_0021_A2NVRP4O6IOOXP
A2NC3DH5SD5OPJ	3RQVKZ7ZS8SUPHJLHKBJRSB3JBB727	mouse_0022_A2NC3DH5SD5OPJ
APNGHQ9FGAD2U	3TUOHPJXZ56AGCNJDH4O9Y4GPBKXW9	mouse_0027_APNGHQ9FGAD2U
A1DCI09Y44WRE4	3LG268AV4WFW2GFGJI4QYABSHQVERF	mouse_0030_A1DCI09Y44WRE4
A2HBMQDO9TNG7B	34YWR3PJ3WJH8MMYKAB323VCO660XB	mouse_0032_A2HBMQDO9TNG7B
A4DA578ZUSMH6	3BFNCI9LZ8ZJAZF4LKPOA2GVPPB73I	mouse_0034_A4DA578ZUSMH6
A27ZS8SAIX3YC1	3ZC62PVYE5JP5CS9NM0ABHINYDQXX2	mouse_0036_A27ZS8SAIX3YC1
ADQM2WYLEFJB2	39TX062QYPXDPELV3XTM41TJRIQ3X9	mouse_0041_ADQM2WYLEFJB2
A2MYOAKLVPWYST	33BFF6QPJPKY0EG5TSX02IFCIGJW38	mouse_0042_A2MYOAKLVPWYST



	Reject	Approve
9W.png	Only 1 out 5 elements segmented	
34.png	Only 1 out 5 elements segmented	
YB.png	Only 1 out 5 elements segmented	
O.png	Only 1 out 5 elements segmented	
EC.png		X
SR.png	Only 1 out 5 elements segmented	
8.png	Only 1 out 5 elements segmented	
Z.png	Only 1 out 5 elements segmented	
ZU.png		X
XP.png	Only 1 out 5 elements segmented	
PJ.png	Only 1 out 5 elements segmented	
U.png	Only 1 out 5 elements segmented	
E4.png		X
7B.png		X
6.png	Only 1 out 5 elements segmented	
C1.png		X
2.png	Only 1 out 5 elements segmented	
ST.png	Only 1 out 5 elements segmented	

Figure 29. The table below contains a sample of the resolution of the results of one of the task batches:

### 5.7. (M7) Creation of a traceable dataset

The information extracted from the data labelling and the survey was combined using the worker ID on Amazon Mechanical Turk to link the form responses (M5) to the outcome of their task.

```
{
  "images": [
    {
      "id": 1,
      "file_name": "mouse_0001_A2YFFC90MO0MW8.png",
      "width": 2880,
      "height": 1920
    },
    {
      "id": 2,
      "file_name": "mouse_0003_A3QQ78M8JAA0V8.png",
      "width": 2880,
      "height": 1920
    },
    {
      "id": 3,
      "file_name": "mouse_0004_A1O667FY1Z4MV5.png",
      "width": 2880,
      "height": 1920
    },
    {
      "id": 4,
      "file_name": "mouse_0202_A1NW6Q6P3NBHWE.png",
      "width": 2880,
      "height": 1920
    },
    {
      "id": 5,
      "file_name": "mouse_0010_A2BSAP480FLLOJ.png",
      "width": 2880,
      "height": 1920
    },
    {
      "id": 6,
      "file_name": "mouse_0011_A2J86NI89IOP17.png",
      "width": 2880,
      "height": 1920
    }
  ],
}
```

```

{
  "id": 7,
  "file_name": "mouse_0012_A2FHUK8FWG52DD.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 8,
  "file_name": "mouse_0205_A1EPTAZ6SJP062.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 9,
  "file_name": "mouse_0013_AXDZVNWHJTKPE.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 10,
  "file_name": "mouse_0206_A232HEY81AMQQ8.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 11,
  "file_name": "mouse_0015_A2E1PL4A2EQNNB.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 12,
  "file_name": "mouse_0017_A28CGW2MYUNU7I.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 13,
  "file_name": "mouse_0018_A21X40K1JQVQZJ.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 14,
  "file_name": "mouse_0019_AWZKTYEDKNIRB.png",
  "width": 2880,
  "height": 1920
},
{
  "id": 15,
  "file_name": "mouse_0211_ATYQPF0GPASDH.png",
  "width": 2880,
  "height": 1920
},

```

<  
Extract from the data set in  
JSON format showing the images  
identified by ID.

This identifier links the task to  
the worker's sociodemographic  
data.

```

{
  "id": 565,
  "image_id": 155,
  "category_id": 4,
  "bbox": [
    1250,
    1,
    106,
    354
  ],
  "area": 37524,
  "iscrowd": 0,
  "attributes": {
    "worker_id": "ADQMQ2WYLFJB2",
    "form_responses": {
      "Q1": "25 - 34",
      "Q2": "Male",
      "Q3": "Bachelor\u2019s degree (comple-
ed)",
      "Q4": "illions",
      "Q5": "More tha 30 minutes",
      "Q6": "1 - 10",
      "Q7": "In between 2 hours and 8 hours",
      "Q8": "0 - 5 times",
      "Q9": "Neck pain, Eye strain",
      "Q10": "$1,200 laptop"
    }
  }
},
{
  "id": 566,
  "image_id": 156,
  "category_id": 4,
  "bbox": [
    1269,
    0,
    91,
    361
  ],
  "area": 32851,
  "iscrowd": 0,
  "attributes": {
    "worker_id": "A2MYOAKLVPWYST",
    "form_responses": {
      "Q1": "35 - 44",
      "Q2": "Male",
      "Q3": "Bachelor\u2019s degree (comple-
ed)",
      "Q4": "CALIFORNIA",
      "Q5": "5 - 10 min",
      "Q6": "1 - 10",
      "Q7": "10 - 30 minutes",
      "Q8": "5 - 10 times",
      "Q9": "None",
      "Q10": "500"
    }
  }
}

```

<  
>

Another extract from the JSON showing the nested survey object, with all the responses. At the beginning of each object, you can see the 'id' that links it to the frame where the information should appear.

```

    },
    {
      "id": 567,
      "image_id": 158,
      "category_id": 4,
      "bbox": [
        1295,
        4,
        93,
        388
      ],
      "area": 36084,
      "iscrowd": 0,
      "attributes": {
        "worker_id": "A35096GSACJWJM",
        "form_responses": {
          "Q1": "18- 24",
          "Q2": "Male",
          "Q3": "Bachelor\u2019s degree (comple-
ed)",
          "Q4": "NON",
          "Q5": "5 - 10 min",
          "Q6": "51 - 100",
          "Q7": "30 - 60 minutes",
          "Q8": null,
          "Q9": "Neck pain",
          "Q10": "NON"
        }
      }
    },
    {
      "id": 569,
      "image_id": 160,
      "category_id": 4,
      "bbox": [
        1288,
        0,
        219,
        379
      ],
      "area": 83001,
      "iscrowd": 0,
      "attributes": {
        "worker_id": "A1ISZTMNHXPB4C",
        "form_responses": {
          "Q1": "35 - 44",
          "Q2": "Male",
          "Q3": "Bachelor\u2019s degree (comple-
ed)",
          "Q4": "good",
          "Q5": "5 - 10 min",
          "Q6": "1 - 10",
          "Q7": "10 - 30 minutes",
          "Q8": "5 - 10 times",
          "Q9": "None",

```

## 5.8. (M8) Overlapping of labelling data and online form responses

The labelling data from Amazon Mechanical Turk's JSON file is the raw material to be superimposed on the computer mouse video. Below, I describe the process I followed and illustrate the project structure to document the file relationships.

The workflow developed for this project is organised into a structured sequence of folders and scripts that manage and process the various elements of the dataset.

First, we start with a folder containing the original images extracted from the video

### **project\_root/**<sup>01</sup>

- mouse\_video\_frames/ (*Imágenes etiquetadas por los trabajadores*).
- mouse\_0001\_A2YFFC90MO0MW8.jpeg
- ...
- labels/ (*Los datos de etiquetado están en formato .txt*).
- mouse\_0001\_A2YFFC90MO0MW8.txt
- ...
- metadata.json (*Datos socioeconómicos de la encuesta Qualtrics*).
- main\_still\_mouse\_detection.py
- setup\_mouse\_annotation.sh (*Configuración del entorno y las dependencias*).
- annotated\_images/ (*carpeta generada por el script de anotación*).
- mouse\_0001\_A2YFFC90MO0MW8.jpeg
- ...
- generate\_video.py (*script that created the video*)
- output\_mouse\_video.mp4 (*final artefacto*)

**01** The project folder is accessible in the Appendix section.

(mouse\_video\_frames/). These are the images that the workers labelled on Amazon Mechanical Turk. Second, there is a folder with the annotation files in YOLO format (labels/) — You Only Look Once, is a pre-trained recognition algorithm that requires a predefined file format and structure to work. I decided to use this format because I initially proposed the idea of using a neural network such as YOLO<sup>01</sup>. However, due to technical and time constraints, I resorted to my own script to overlay video with data. The labelling files are in .txt format and have been extracted from the global .json file shown in the previous section.

This data is accompanied by an additional file called metadata.json, which contains sociodemographic and work context information provided by the people who performed the labelling and completed the Qualtrics form. The main script —main\_still\_mouse\_detection.py— integrates these three components to generate a visual representation in which bounding boxes are superimposed with personal information about each worker.

All generated images are stored in the annotated\_images/ folder. Finally, the generate\_video.py script compiles this sequence of images into a video file (output\_mouse\_video.mp4).

<sup>01</sup>. The file format and structure correspond to those required for training YOLO, as this is what I decided to use initially. Due to training errors, I decided to use a custom algorithm to overlay data on the images.



mouse\_0001.

mouse\_0002.

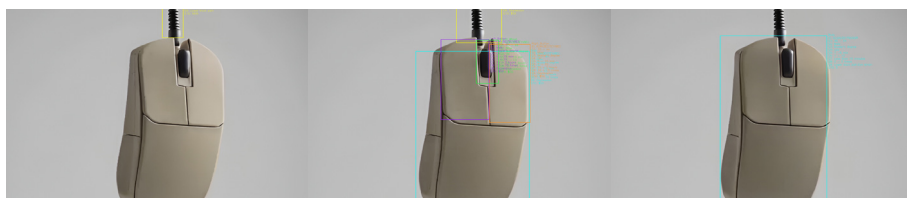
mouse\_0003.



mouse\_0004.

mouse\_0005.

mouse\_0006.



mouse\_0007.

mouse\_0008.

mouse\_0009.



mouse\_0010.

mouse\_0011.

mouse\_0012.



mouse\_0013.

mouse\_0014.

mouse\_0015.





## 6. Conclusions

The invisible workforce involved in AI is a growing area of research given the obvious inequalities and evasion of responsibility by the Western companies that ‘employ’ them. Despite this, it has been revealing to see how members who influence some point in the value chain of this technology, such as ethicists and engineers, were unaware of the human effort involved in data work, although they show a keen interest in learning more about their participation and experiences. This suggests that more education is needed, not only among Western society but also among the workers themselves in the conceptualisation chain. That is why the interviews with Joan Kinyua and her team have helped to highlight new examples of obfuscated extractivism that can serve as educational cases for the industry.

Large US corporations, such as Meta, Google and Scale AI, use Kenyan data labour. The particular case of Scale AI has led to the creation of the crowdsourcing platform Remotasks, which provides opaque access to this workforce without them being on the payroll. Although this microtask company has now been banned in the country, the parent companies are able to reach the same workforce through BPOs such as SAMA, which offer higher quality results due to internal training.

Regulation of data work in the Global South is scarce due to a lack of knowledge of the processes and lower levels of government involvement. This lack of knowledge is the main reason why Western companies turn to offshoring. And the intrinsic opacity of this work with clients makes traceability difficult for accountability purposes. Despite this, associations such as the DLA (Data Labelers Association) are pressuring leaders to go to the Senate and generate a social and educational discourse on the mental and physical repercussions that this work is having on the citizens of their country, Kenya. Other international projects, such as Data Workers, are creating dynamics with members of this workforce and academics to develop proposals that will help governments in both the North and South to regulate the rights and access of companies to data work.

Currently, any LLM is likely to use this type of worker. Due to the huge amount of data required for them to work accurately, they require a vast workforce. To train them and verify their results. However, it is not necessary for an LLM to be involved, as there are multiple fake algorithms where a person imitates the behaviour of software. This is the case of the Madagascar imitators, who watched security cameras in shopping centres in France in real

time to detect and report thefts.

These types of jobs require high concentration and are not well paid to provide a decent living. Salaries on microtask platforms are quoted in cents, and salaries in BPOs do not reach the average in countries such as Kenya. Work-related pathologies influence the quality of life of workers. They suffer physical pain, stress and neurotic behaviour, as in the case of Maundu, who developed neurosis due to the nature of his work.

Recognition of this workforce must be based precisely on its vulnerabilities. This has led me to design a dataset that stores and tracks their personal conditions, city of residence, pathologies and work set, in order to give substance to this hidden workforce. In this way, I am contributing a dataset that can be used in a specific case of visual recognition and that links the professional contributions of each of the workers who participated in training the algorithm. The ability to use this dataset to create a video that reflects the conditions of the workers provides visual and easily digestible support for educating people about this type of work.

### *5.1 Limitations of the research*

I would like to emphasise that the interviews with data workers have been limited

to a recently created group, such as the DLA. Although, according to the literature review, some working conditions and pathologies coincide with those suffered by other groups in Asian and Latin American countries, the sample is too small to extrapolate the results to the rest of the Global South. On the other hand, Amazon Mechanical Turk has been the quick and affordable alternative for creating the dataset. This has been to the detriment of the quality of the labelling. Despite two rounds of review to detect fraudulent workers and unfinished labels, it is recommended that the review phase be extended to ensure an accurate and complete dataset. Alternatively, other more ethical and controlled platforms, such as Prolific, should be used to obtain a higher quality dataset that can be used to accurately train a real visual recognition model.

On the other hand, my programming skills are limited. Although I have no difficulty reading and understanding the code, the overlaying of data with video (M8) was assisted by Open AI's 04-mini-high model through the creation of several scripts. Any recommendations given by this model may call into question the efficiency of the process.

## *6.2 Future lines of research*

Traceability in generative artificial intelligence opens up three complementary

lines of application with high transformative potential. Firstly, as a labour supervision tool, it allows the incorporation of labour and socio-demographic metadata linked to each contribution in the AI value chain, making it possible to audit who generated the data, under what conditions and for which client. This would facilitate regulatory processes and guarantee labour rights in opaque environments. çW

Secondly, it acts as an analytical tool for detecting bias by enabling reverse audits that link annotative decisions to specific social contexts, helping to identify structural patterns of exclusion or distortion in models.

Thirdly, traceability provides a link between producers and end users. For example, in models that generate images of architectural interiors, it would make it possible to reconstruct the provenance of images and contact the people who collected them, opening up opportunities for direct collaboration and new forms of professional recognition. These three approaches—oversight, critical analysis, and connection—redefine traceability as a central component of socially responsible artificial intelligence.







## 7. References

- Anwar, M (2023) Value Chains of AI: Data Training Firms Platforms and Workers. *University of Edinburgh*. <https://orcid.org/0000-0002-5213-4022>
- Braz, M. V., Tubaro, P., & Casilli, A. A. (2024). Fabricar os dados: O trabalho por trás da Inteligência Artificial.
- Casilli, A. A. (2024). Digital Labor and the Inconspicuous Production of Artificial Intelligence.
- Casilli, A. A., & Tubaro, P. (2023.). An End-to-End Approach to Ethical AI: Socio-Economic Dimensions of the Production and Deployment of Automated Technologies.
- Casilli, A. A., Tubaro, P., Cornet, M., Ludec, C. L., Torres-Cierpe, J., & Braz, M. V. (2021). Global Inequalities in the Production of Artificial Intelligence: A Four-Country Study on Data Work.
- Chaudhuri, B., & Chandhiramowuli, S. (2024). Tracing the Displacement of Data Work in AI: A Political Economy of “Human-in-the-Loop.” *Engaging Science, Technology, and Society*, 10(1–2). <https://doi.org/10.17351/ests2024.2983>
- Ekbja, H., & Nardi, B. (2014). Heteromation and its (dis)contents: The invisible division of labor between humans and machines. *First Monday*. <https://doi.org/10.5210/fm.v19i6.5331>

- Evers, C., Khuara, M., Mata, T., Soper, L., Stilgoe, J. (2022) Ghost workers Report—Empowering the invisible labour behind artificial intelligence.
- Hung, K.-H. (2024). Artificial intelligence as planetary assemblages of coloniality: The new power architecture driving a tiered global data economy. *Big Data & Society*, 11(4), 20539517241289443. <https://doi.org/10.1177/20539517241289443>
- Le Ludec, C., Cornet, M., & Casilli, A. A. (2023). The problem with annotation. Human labour and outsourcing between France and Madagascar. *Big Data & Society*, 10(2), 20539517231188723. <https://doi.org/10.1177/20539517231188723>
- Muldoon, J., Cant, C., Wu, B., & Graham, M. (2024). A typology of artificial intelligence data work. *Big Data & Society*, 11(1), 20539517241232632. <https://doi.org/10.1177/20539517241232632>
- Newlands, G. (2021). Lifting the curtain: Strategic visibility of human labour in AI-as-a-Service. *Big Data & Society*, 8(1), 20539517211016026. <https://doi.org/10.1177/20539517211016026>
- Bezos, J., (2006) Opening Keynote and Keynote Interview with Jeff Bezos. Retrieved May 2, 2025, from [https://videlectures.net/videos/mitworld\\_bezos\\_okki](https://videlectures.net/videos/mitworld_bezos_okki)

- Ekbia, H., (2014) Heteromation and its (dis) contents: The invisible division of labor between humans and machines. (n.d.). *ResearchGate*. <https://doi.org/10.5210/fm.v19i6.5331>
- Posada, J. (2022). The Coloniality of Data Work: Power and Inequality in Outsourced Data Production for Machine Learning. *Faculty for Information. University of Toronto*
- Rothschild, A., Wang, D., Jayakumar Vilvanathan, N., Wilcox, L., DiSalvo, C., & DiSalvo, B. (2024). The Problems with Proxies: Making Data Work Visible through Requester Practices. *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 7, 1255–1268. <https://doi.org/10.1609/aies.v7i1.31721>
- Tubaro, P., Casilli, A. A., Cornet, M., Le Ludec, C., & Torres Cierpe, J. (2025). Where does AI come from? A global case study across Europe, Africa, and Latin America. *New Political Economy*, 1–14. <https://doi.org/10.1080/13563467.2025.2462137>
- Tubaro, P., Casilli, A. A., & Coville, M. (2020a). The trainer, the verifier, the imitator: Three ways in which human platform workers support artificial intelligence. *Big Data & Society*, 7(1), 205395172091977. <https://doi.org/10.1177/2053951720919776>

## 8. Bibliography

Arun, C. (2025). The Silicon Valley Effect. SSRN. <https://doi.org/10.2139/ssrn.5109858>

Bell, S. (2023). AI and Job Quality: Insights from Frontline Workers. SSRN Electronic Journal. <https://doi.org/10.2139/ssrn.4337611>

Chandhiramowuli, S., & Chaudhuri, B. (2023). Match Made by Humans: A Critical Enquiry into Human-Machine Configurations in Data Labelling. Hawaii International Conference on System Sciences. <https://doi.org/10.24251/HICSS.2023.251>

Eke, D. O., Wakunuma, K., Akintoye, S., & Ogoh, G. (Eds.). (2025). Trustworthy AI: African Perspectives. Springer Nature Switzerland. <https://doi.org/10.1007/978-3-031-75674-0>

Gebreکیدan, F. B. (n.d.). Content moderation: The harrowing, traumatizing job that left many African data workers with mental health issues and drug dependency.

Gonzalez-Cabello, M., Siddiq, A., Corbett, C. J., & Hu, C. (2024). Fairness in crowdwork: Making the human AI supply chain more humane. Business Horizons. <https://doi.org/10.1016/j.bushor.2024.09.003>

Michel, B., & Ecker, Y. (2025). Seeing economic development like a large language model. A methodological approach to the exploration of geographical imaginaries in generative AI. *Geoforum*, 158, 104175. <https://doi.org/10.1016/j.geoforum.2025.104175>

[org/10.1016/j.geoforum.2024.104175](https://doi.org/10.1016/j.geoforum.2024.104175)

Opening Keynote and Keynote Interview with Jeff Bezos. (n.d.). Retrieved May 2, 2025, from [https://videlectures.net/videos/mitworld\\_bezos\\_okki](https://videlectures.net/videos/mitworld_bezos_okki)

Paulsen, K. (2020). “Shitty Automation”: Art, Artificial Intelligence, Humans in the Loop. *Media-N*, 16(1), 4–23. <https://doi.org/10.21900/j.median.v16i1.227>

Raji, I. D., Scheuerman, M. K., & Amironesei, R. (2021). You Can’t Sit With Us: Exclusionary Pedagogy in AI Ethics Education. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 515–525. <https://doi.org/10.1145/3442188.3445914>

The Platform Proletariat. (n.d.). Pulitzer Center. Retrieved May 2, 2025, from <https://pulitzercenter.org/stories/platform-proletariat>

Tubaro, P. (n.d.). Décrypter la société des plateformes: Organisations, marchés et réseaux dans l’économie numérique.

Witschas, A. (2025). Prefabricated futures: AI imaginaries between elitist visions and social justice claims. In T. Kox, A. Ullrich, & H. Zech (Eds.), *Uncertain Journeys into Digital Futures* (pp. 335–356). Nomos Verlagsgesellschaft mbH & Co. KG. <https://doi.org/10.5771/9783748947585-335>

## 9. Appendix

### 9.1. Online Qualtrics survey

Qualtrics

Q1

What's your age?

☐ 18–24

☐ 25–34

☐ 35–44

☐ 45–54

☐ 55–64

☐ 65 or older

☐ You can add your exact age if you prefer:

Q2

What is your gender?

☐ Male

☐ Female

☐ Non binary

Q3

What is your educational level?

☐ Nor formal education

☐ Primary education

☐ Secondary / High school

☐ Undergraduate (currently studying a Bachelor's degree)

☐ Bachelor's degree (completed)

☐ Master's degree

☐ Doctoral degree

Q4

From which city are you performing this task?

[Texto libre]

Q5

How much time in average do you spend searching for quality tasks?

☐ 1–5 min

☐ 5–10 min

☐ 15–25 min

☐ More than 30 minutes

Q6

Approximately how many tasks have you completed that weren't paid by the requester?

☐ None

☐ 1–10

☐ 11–50

☐ 51–100

☐ More than 100

Q7

How long did the longest image labeling task you worked on take?

☐ Less than 10 minutes

- ☐ 10–30 minutes
- ☐ 30–60 minutes
- ☐ 1h–2 hours
- ☐ In between 2 and 8 hours
- ☐ More than 8 hours

**Q8**

How many times has your work been rejected without compensation in the last month?

- ☐ 0–5 times
- ☐ 5–10 times
- ☐ 10–15 times
- ☐ 15–30 times
- ☐ More than 30

**Q9**

Do you suffer from any work-related health issues?

- ☐ Neck pain
- ☐ Lower back pain
- ☐ Eye strain
- ☐ Headaches
- ☐ Wrist or hand pain
- ☐ Anxiety
- ☐ Depression
- ☐ None
- ☐ Other: \_\_\_\_\_

**Q10**



How much did your workset cost?

(Include any tools or devices you use to perform data labeling)

[Texto libre]

## 9.2. *Expert interviews*

### **[ES] Script eticista & product**

**Objective:** To explore the knowledge that strategy departments have about the other end of the AI life cycle. To understand the ethical, political and cultural implications of the invisibility of human labour in the development of algorithmic systems.

#### **1. Introduction and background**

(Initial questions to contextualise the interviewee's profile and establish a basis of trust)

How would you describe your professional journey to working in artificial intelligence ethics?

What led you to become interested in the ethical dimensions of technological development, especially in relation to AI?

#### **2. Human labour in model training**

How is an artificial intelligence model trained at home?

How do large corporations train their AI models?

#### **3. Disconnection between results and processes**

What ethical implications do you see in representing AI as 'automatic' or 'autonomous' when it is often supported by precarious human labour?

#### **4. Performance-focused validation**

What ethical and social risks do you perceive in validating models solely on the basis of their technical performance?

What contradictions do you identify between discourses of corporate ethical responsibility and the practice of outsourcing work under opaque conditions to train models? (Especially in Europe)

#### 5. Romanticisation of AI

Is artificial intelligence inevitable?

What role do discourses on technological autonomy play in making human labour invisible?

#### 6. Ethical traceability

What ethical responsibility do you think researchers and academics have when using datasets whose origin or conditions of production are unclear?

How could the traceability of human labour be incorporated as a criterion in the ethical evaluation of AI projects?

#### 7. Closing and open reflection

(Final space for the interviewee to freely add or elaborate)

What question do you think should be asked more often when we talk about ethics and work in AI?

If you had the opportunity to interview someone who works labelling data to train AI models, what would you ask them? What would you like to know about their experience, their working environment or their view of the role

they play in the system?

Is there anything we haven't covered that you think is important to add to this conversation?

### 8. Cognitive map

Presentation of the exercise

'Let's do an exercise to visualise how you understand a complex topic related to artificial intelligence. It's not about getting it right, but about representing how you think about it. We will use key words or concepts and the connections you perceive between them.'

2. Focus topic (choose one according to the participant's profile)

'I want us to think about the complete cycle of artificial intelligence development. We are going to collect the ideas, concepts or elements that you think are important in this topic. It is very likely that you will not be familiar with many of the areas that make it up, the aim is not to have an exact view but rather a map of relationships as you think they are.'

### 3. Concept elicitation

Questions to elicit concepts:

What elements do you consider fundamental in this topic?

Which actors are involved?

What conflicts or tensions arise?

What concepts come to mind when you think about this?

**[ES] Script producto & tech**

**Objective:** To explore how people working in AI-based product strategy and launch understand and manage the invisible human processes behind model development. We are interested in understanding the role (or lack thereof) of labour and ethics in their design and business decisions.

### 1. Introduction and background

(To contextualise the interviewee's profile and their role in the AI value chain)

How would you describe your professional background that led you to lead or design product strategies in artificial intelligence?

What kind of strategic decisions do you usually make regarding the design, positioning or deployment of AI models?

### 2. Visibility of human work

How much do you know about how the models that make it to the product are trained?

Has the origin of the data or who prepares it ever been a topic of conversation or decision-making in your team?

How do you think the fact that this work is invisible or little recognised affects (or does not affect) the business strategy?

### 3. Concept of value and process

What role does the perception of autonomy or 'intelligence' of the model play in the value proposition you design?

How is the product narrative constructed internally? Is the human element behind it included or omitted?

What tensions arise between technical efficiency and work ethics when deciding on roadmaps or priorities?

#### 4. Validation criteria and performance ethics

How much weight do ethical criteria carry compared to technical performance or commercial scalability?

Are there any internal protocols or guidelines on responsibility in the use of datasets or labelling providers?

Have you encountered any dilemmas or contradictions that you have had to resolve in product design or launch?

#### 5. Organisational culture and boundaries

How would you describe your organisation's culture in relation to transparency, traceability or accountability in AI processes?

What barriers do you encounter when incorporating ethical or social criteria into product decisions?

Have there been situations where attempts have been made to highlight human work (e.g. data workers) within the brand narrative?

#### 6. Closing and feedback

(To gather loose ideas or unplanned reflections)

What changes do you think could facilitate the real integration of ethics and humanity into AI product strategy?

What question should someone researching the relationship between ethics and the AI

market ask you that I haven't asked you yet?

Is there anything else you would like to add to conclude this conversation?

## 8. Cognitive map

### Presentation of the exercise

‘Let's do an exercise to visualise how you understand a complex topic related to artificial intelligence. It's not about getting it right, but about representing how you think about it. We will use key words or concepts and the connections you perceive between them.’

Thematic focus (choose one according to the participant's profile)

‘I want us to think about the entire cycle of artificial intelligence development. We are going to gather the ideas, concepts or elements that you think are important in this topic. It is very likely that you are not familiar with many of the areas that make it up. The goal is not to have an exact view, but rather a map of relationships as you see them.’

### Concept elicitation

#### Questions to elicit concepts:

What elements do you consider fundamental to this topic?

Which actors are involved?

What conflicts or tensions arise?

What concepts come to mind when you think about this?

**[EN] Script data worker**

**Objective:** To understand how the invisible labor that sustains the training of artificial intelligence models is lived, perceived, and experienced from the perspective of those who perform it.

**1. Background and trajectory** (Start of the interview: simple and open-ended questions to build trust and understand the participant's profile)

How would you describe your work trajectory before entering data labeling or processing?

What circumstances or motivations led you to enter this field?

**2. Tasks and working conditions**

What kind of tasks do you perform on a daily basis?

How would you describe your working conditions (hours, pay, stability, social protection)?

What kind of training, instructions or support do you usually receive?

**3. Relationship with development processes**

What level of contact do you have with those who design the projects or models you work on?

In what ways do you feel your work contributes to or influences the final outcomes?

Are there channels or spaces available to raise concerns, questions or errors?

**5. Invisibility, identity, and agency**

What effect does it have on your daily life that your work is not publicly recognized?

How does that affect your sense of belonging, self-esteem, or agency as a worker?



What consequences do you see at a social or political level from the systematic invisibilization of this work?

What kinds of changes or measures do you think are needed to make your work visible and valued?

6. Ethics and dilemmas in everyday practice  
(Transition toward more complex and personal topics)

What ethical dilemmas or conflicts have you faced while performing data labeling tasks?

How do you perceive the relationship between tech companies' ethical discourse and your own working conditions? (Especially in Europe)

What cultural or structural mechanisms contribute to the invisibilization and devaluation of the human labor behind the creation of datasets?

6. Value generation and distribution

How do you perceive the relationship between your labor and the benefits obtained by companies or institutions using that data?

What mechanisms, proposals or forms of organization could make that relationship fairer?

7. Closing and reflection (Final block to collect proposals and open space for feedback)

What questions do you think should be asked by those who design ethical frameworks for artificial intelligence?

What urgent transformations would you like to see in how your work is organized or recognized?

Is there anything I haven't asked that you think is important to share?

8. Cognitive Map

### Exercise Introduction

“We're going to do an exercise to help visualize how you understand a complex topic related to artificial intelligence. This isn't about being correct, but rather about representing how you think about it. We'll use keywords or concepts and the connections you perceive between them.

Thematic Focus (choose one depending on the participant's profile)

“I'd like us to think about the full development cycle of an artificial intelligence system. We'll gather the ideas, concepts, or elements you consider important to this topic. It's very likely you won't have visibility into many of the areas that make it up—that's fine. The goal is not to have an exact picture, but to create a relational map based on how you believe it works.

### Concept Elicitation

Questions to prompt concepts:

What elements do you consider fundamental to this topic?

What actors are involved?

What conflicts or tensions arise?

What concepts come to mind when you think about this?

### *9.3. Other materials*

**Transcripción y consentimientos para las entrevistas disponibles en:** <https://drive.google.com/drive/folders/1vYersVvP4uzcv9gn6xLI8FcpnZrrbmQp>

**Proyecto de reconocimiento visual completo, disponible en:** <https://drive.google.com/drive/folders/19tJj9iOtYdSjIJPslrCYPAt11-8OahZI>

**Video pedagógico de etiquetado: Humanos en el bucle:** <https://drive.google.com/drive/folders/1pC50UVW57IfSW0h3SCeVbnPoeEOvhqA3>

