# ROADMAP ON AI TECHNOLOGIES & APPLICATIONS FOR THE MEDIA INDUSTRY

## SECTION: "OPEN AI REPOSITORIES AND INTEGRATED INTELLIGENCE: NEXT STEPS TOWARDS AI DEMOCRATISATION"

info@ai4media.eu          www.ai4media.eu

| Authors | Sven Becker (Fraunhofer Institute for Intelligent Analysis and Information Systems - IAIS) |
|---|---|
| | Annerike Flint (Fraunhofer Institute for Intelligent Analysis and Information Systems - IAIS) |
| | Laszlo Friedmann (Fraunhofer Institute for Intelligent Analysis and Information Systems - IAIS) |
| | Andreas Steenpass (Fraunhofer Institute for Intelligent Analysis and Information Systems - IAIS) |

This report is part of the deliverable D2.3 - *"AI technologies and applications in media: State of Play, Foresight, and Research Directions"* of the AI4Media project.

You can site this report as follows:

The information and views set out in this report are those of the author(s) and do not necessarily reflect the official opinion of the European Union. Neither the European Union institutions and bodies nor any person acting on their behalf.

# Open AI repositories and integrated intelligence: next steps towards AI democratisation

The media sector is already applying AI technologies. For example, users consuming media already enjoy live subtitling, automatic translation into other languages, or into plain language. One might argue that these mechanisms have introduced ***integrated intelligence***, i.e. adding a software component that makes a product intelligent, such as live subtitling for a live video stream.

To design an AI system, one needs data, storage and computational resources, algorithms, and AI talent. On all these levels, ***openness*** plays an important role and stretches across multiple levels of AI development.

When it comes to data, thousands of open source datasets are published on different platforms like GitHub (awesome-public-datasets[1]), Kaggle[2] or by Google[3], to name a few. The limitation of storage and computational resources was addressed during the last years by three major AI cloud providers: Amazon Web Services (AWS), Microsoft Azure and Google Compute Platform. In November 2018 GitHub announced that 100 million code repositories are live on GitHub. According to Jason Warner, the repositories stem from 31 million developers from almost all countries in the world, which are collaborating across 1.1 billion contributions.[4] In recent years Microsoft, Google and Amazon have provided development tools to allow practitioners to create AI applications without deep knowledge of machine learning, linear algebra, or statistics. This addresses the issue of limited AI talent and the transfer of AI tools into multiple industry sectors. The provided technologies are designed as ***modular tool boxes***.

In general, the widespread use of open source software has increased in popularity. TensorFlow, a machine learning framework for basic machine learning model development (developed by Google), is in the top 10 of all GitHub projects with over 150,000 star ratings.[5] However, the importance of open source is also reflected in the availability of much more specific AI frameworks for solving tasks for example in the domain of Natural Language Processing (NLP) (Figure 1).

---

[1] Awesome public datasets: https://github.com/awesomedata/awesome-public-datasets
[2] Kaggle datasets: https://www.kaggle.com/datasets
[3] Google public datasets: https://cloud.google.com/bigquery/public-data/
[4] Warner, Jason. "Thank you for 100 million repositories". Last accessed: 16 December 2021.
[5] GitHub Ranking, top 100 stars:https://github.com/EvanLi/Github-Ranking/blob/master/Top100/Top-100-stars.md

**Top 10 Open AI repositories on GitHub**

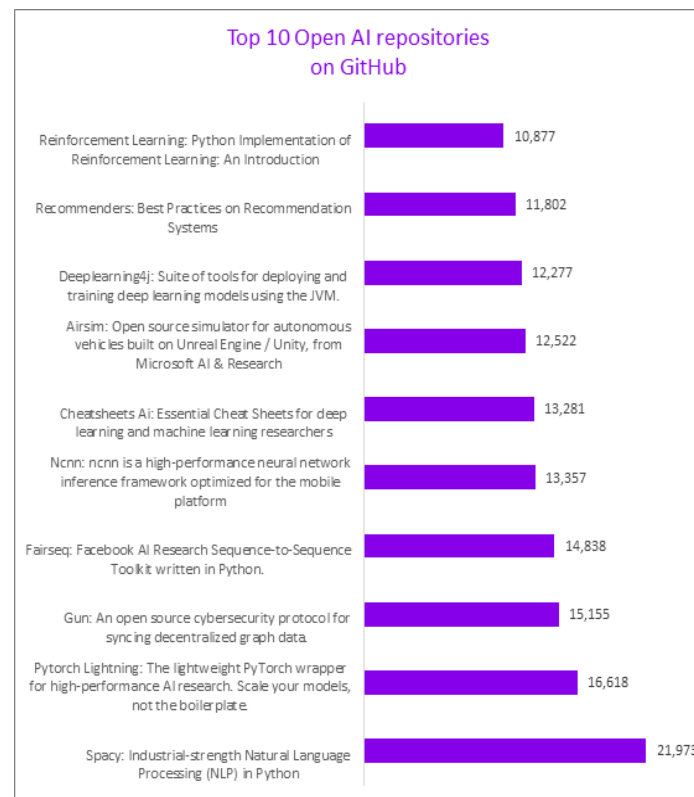| Repository | Stars |
|---|---|
| Reinforcement Learning: Python Implementation of Reinforcement Learning: An Introduction | 10,877 |
| Recommenders: Best Practices on Recommendation Systems | 11,802 |
| Deeplearning4j: Suite of tools for deploying and training deep learning models using the JVM. | 12,277 |
| Airsim: Open source simulator for autonomous vehicles built on Unreal Engine / Unity, from Microsoft AI & Research | 12,522 |
| Cheatsheets Ai: Essential Cheat Sheets for deep learning and machine learning researchers | 13,281 |
| Ncnn: ncnn is a high-performance neural network inference framework optimized for the mobile platform | 13,357 |
| Fairseq: Facebook AI Research Sequence-to-Sequence Toolkit written in Python. | 14,838 |
| Gun: An open source cybersecurity protocol for syncing decentralized graph data. | 15,155 |
| Pytorch Lightning: The lightweight PyTorch wrapper for high-performance AI research. Scale your models, not the boiler plate. | 16,618 |
| Spacy: Industrial-strength Natural Language Processing (NLP) in Python | 21,973 |

*Figure 1: List of the top 10 AI repositories on GitHub, according to Awesome Open Source[6].*

**Openness** in information technology (IT) is referred to as a feature of an IT system. Such a system is characterised by interoperability, portability, and extensibility. These can be implemented using interfaces, standards, and IT architectures. Next to these technical aspects, openness is also based on partnership between the involved partners (IT customers, IT vendors and/or IT service providers).[7] Schlagwein et al. define openness as "[...] often deeply embedded in information technology (IT) and [as] both a driver for and a result of innovative IT."[8]

**Open source** is source code which is freely available, modifiable and redistributable. The distribution terms of open-source software are handled by licenses (e.g., Apache License, 2.0 or MIT License). **Open access (OA)** refers to open and free-of-charge access to scientific publications for all entities worldwide. OA also means that research results are ideally flexibly reusable. OA publications are always electronic or published online, which may appear in printed form later on.[9]

A **repository** is a storage location for software. It usually includes a table of contents and code documentation. It is typically under version control and serves as a working directory for software development, often related to a specific project or software module.

---

[6] Awesome Open Source: https://awesomeopensource.com
[7] Steven, Vettermann. "Code of Openness". CPO. ProSTEP iViP. Retrieved 10 January 2017.
[8] Schlagwein, D., Conboy, K., Feller, J. et al. "Openness" with and without Information Technology: a framework and a brief history. J Inf Technol 32, 297–305 (2017). https://doi.org/10.1057/s41265-017-0049-3
[9] P. Suber. "Open Access Overview". Retrieved 29 November 2014.

**Democratisation** is the action of making something available to everyone. This includes introducing democratic systems and principles. Recently, many major tech companies, e.g. Google or Microsoft, have made the democratisation of AI one of their major goals. At this point it is still unclear how this claim differs from e.g., the open access movement.[10] However, it is clear that the democratisation of AI must champion for broader participation. This concern stems from today's reality: the ability to engineer and make use of AI technologies rests in the hands of a privileged few.[11] To disrupt the workings of today, the barriers to participating in AI development need to be reduced.

## Research challenges

Today, AI technology is largely only accessible to advanced users – users that know how to code, how to access repositories, and have the knowledge to apply code to their specific use cases. On the one hand, low tech companies wanting to apply AI might lack the knowledge of what the technology can do for them or lack the AI talent or financial means to employ them. On the other hand, AI savvy companies fear losing their intellectual property and restrict the accessibility to code, models, and data. Bridging the gap and championing democratisation and participation will be one of the major tasks to advance AI in a way that it becomes a common good. More specifically, this task will have to address the following challenges:

*Suitable general APIs*. First, the AI tools in any modular toolbox must implement suitable APIs. These APIs are often designed with a specific application in mind. On the other hand, the tools should be reasonably versatile, that is, it should be possible to reuse them for other, and future, applications. For this purpose, the implemented APIs must be sufficiently general. What, for example, should be the format of a speech-to-text transcription? The transcription can surely be represented as a string, but for some applications, it might also be interesting to provide confidence scores for the recognised words. This also leads to the question whether formats which are widely used for data exchange between *microservices* today such as JSON and XML are sufficient for representing the input and output of AI tools. How can the standardisation process which is needed to synchronise the data formats for a given modular toolbox be organised?

*Simplified and configurable modular toolboxes*. Second, setting up an actual application by combining AI tools from a modular toolbox is often not just straightforward. Instead, it requires specific knowledge and skills regarding the used tools, the toolbox or, in other words, the platform for connecting these tools, and the computational environment such as the server or the cluster where the application is to be deployed. These requirements can be a high burden for inexperienced individuals and teams. Therefore, the following questions need to be answered: How can the deployment of AI applications arising from modular toolboxes be simplified? On the other hand, how can very simple "one-click" solutions be designed in a way

---

[10] Garvey C. A framework for evaluating barriers to the democratization of artificial intelligence. InThirty-Second AAAI Conference on Artificial Intelligence 2018 Apr 29.

[11] Wolf CT. Democratizing AI? experience and accessibility in the age of artificial intelligence. XRDS: Crossroads, The ACM Magazine for Students. 2020 Jul 9;26(4):12-5.

that they remain configurable for more experienced users, for example with respect to the scalability and the allocated resources of the deployed services?

***Legal and ethical concerns***. Finally, the usage of modular toolboxes also raises legal and ethical questions. When applications are designed by combining AI tools from modular toolboxes, how can the contributions of all involved developers and researchers be acknowledged? What are the requirements with respect to licensing and data protection? How can the transparency and the security of such an application be ensured? How can its trustworthiness be assessed? Who will be responsible for infringements and damages caused by such an application? Although many approaches have been proposed and tested, these questions have not yet been fully answered.

## Societal and media industry drivers

### Vignette 1: Streaming pipeline for news live feeds based on open source AI modules

It is the year 2030, Catarina is an editor at a news outlet. Her daily operational task is overseeing the websites' live feeds, e.g., weather report, traffic information, or live tickers for real-time news events. These feeds are automatically created by an AI system. It has been created from an open source repository and it has learned the semantics from millions of weather reports written in English. A streaming pipeline has been established that automatically feeds new input from other news outlets into the model to improve its output. The AI system makes use of social media accounts and, if applicable, includes videos or images from social media users on the live feed, e.g., photos from extreme rain, wind etc. Additionally, all content that is published online has an attached digital ledger that regulates ownership. Therefore, each digital object has an attached license, defining the possibility to (re)use the asset for specified purposes. The AI system might include an open licensed photo of a sunny beach in today's weather feed and is therefore able to clearly identify and notify the owner. If necessary, remuneration is also easily managed through being able to clearly identify how many times a content piece has been used. If content is duplicated without consent by the owner, digital watermarks are included.

### Vignette 2: Supporting cinematographic creativity by open AI tools

After two and a half years of studying at the Madrid Film School, Ricardo is about to become a director. For his graduation, he is working on the realisation of a movie which he had in his mind for some years already. The story line has been updated a few times and is perfect now. Ricardo owns a decent camera and a powerful computer for the postproduction of the footage. For reducing the costs, he has borrowed a professional microphone and he has also casted some of his friends. However, the budget is very low, and he cannot afford to pay musicians for composing and recording music for his movie. Furthermore, professional postproduction software nowadays includes AI tools for adjusting the lightning of a scene and even the mimic expressions of the actors and can therefore be very expensive. Luckily Ricardo has learned as part of his curriculum that there is a European open repository for AI tools where he can find what he needs free of charge. He downloads a tool for automatic music composition and an AI tool for special effects. These tools are easy to integrate into his software setup and he has already used them for other projects. As a further benefit, he knows that they come from

trustworthy sources and are compliant with the legal requirements such as licensing and privacy protection. Relieved that he does not have to use unlicensed demo versions anymore, Ricardo finishes the postproduction within a few days because the next creative project is already on his mind.

**Vignette 3: Using Open Source NLP for news media delivery in any language**

Sophia is a developer at a news outlet and is responsible for multiple website features. Users of the website should be able to choose any language to consume any content. For years, the news outlet management was concerned to provide these features due to bias in the AI models. Since last year, Sophia can easily incorporate models from an international cooperation for Open Source NLP to ensure that language is not a barrier for digital content. Sophia interfaces her software to large, already trained AI translation models. Those models were trained on thousands of openly curated language data sources and are easily accessible in the cloud. The cooperation ensures large, diverse, and well curated data sets for models to be trained on. The respective open repositories for code, data and models have governance mechanisms in place that regulate the shared resources to meet the trustworthiness guidelines, e.g. model interpretability, high model performance, minimised bias and transparency to users through open documentation and system architecture design. Hereby, she follows her company's media guidelines for trustworthy content.

## Future trends for the media sector

***Modular tool boxes*** are the future of open AI repositories. The core tasks are always the same, for the media sector e.g., keyword extraction, named entity recognition, face recognition, object detection. These tools can then be connected in series and adapted to the intended use. A difficulty lies in the definition of interfaces between the modules. On the one hand, they should be specific enough to be able to use the functionality of the respective tool, and on the other hand, they should be general enough so that the tool can be deployed in a multitude of workflows.

In connection with the modular tool boxes, the trend is moving in the direction of ***open source software***, and even more specifically in the direction of software without copyleft and without restrictions on commercial use. A recent study from the European Commission highlights that procuring open source software instead of proprietary software could increase the digital autonomy of the public sector. This could reduce total cost and avoid vendor lock-in effects. Moreover, the study predicts that an additional 0.4% - 0.6% of GDP would be generated, if contributions to open source code are increased by 10%.

The interplay of open repositories and open source software will lead to the **democratisation of AI technology** that is accessible for everybody no matter the technological skills or tools.

The European strategy for data foresees the creation of **interoperable data spaces** in strategic sectors[12], including the media sector. The common European data space for media will ensure

---

[12] Common European Data Spaces: http://dataspaces.info/common-european-data-spaces/

interoperable and easy access to key datasets. It will provide an ecosystem for the creation of solutions, tools and models for the creation, curation and distribution of media content[13].

## Goals for next 10 or 20 years

Following the arguments from the previous subsections, the main goal for the years to come is the establishment of an *open repository for AI tools* which can be used by the media, but also by other sectors. The tools within this repository should cover a *wide range of AI applications* and they should also continuously make the most recent research available.

One example for a trailblazing open repository is the AI4EU Experiments[14] platform from the AI4EU project[15] (see Figure 2), funded through the European Commission's Horizon 2020 program. As of today, the platform has more than 200 onboarded AI models and provides functionality for matching models. When creating pipelines, the software automatically provides info for matching ingoing and outgoing pipelines, making it easy for users to create AI pipelines from scratch.

The platform is the cornerstone of a future ecosystem for AI tools and will be growing with contributions from all European-funded AI projects. It is also open to all market participants who wish to contribute technologies. Regarding the provided AI tools, the focus should be on open source software in order to ensure transparency, availability, and sustainability. However, commercial tools should not be excluded from the repository.
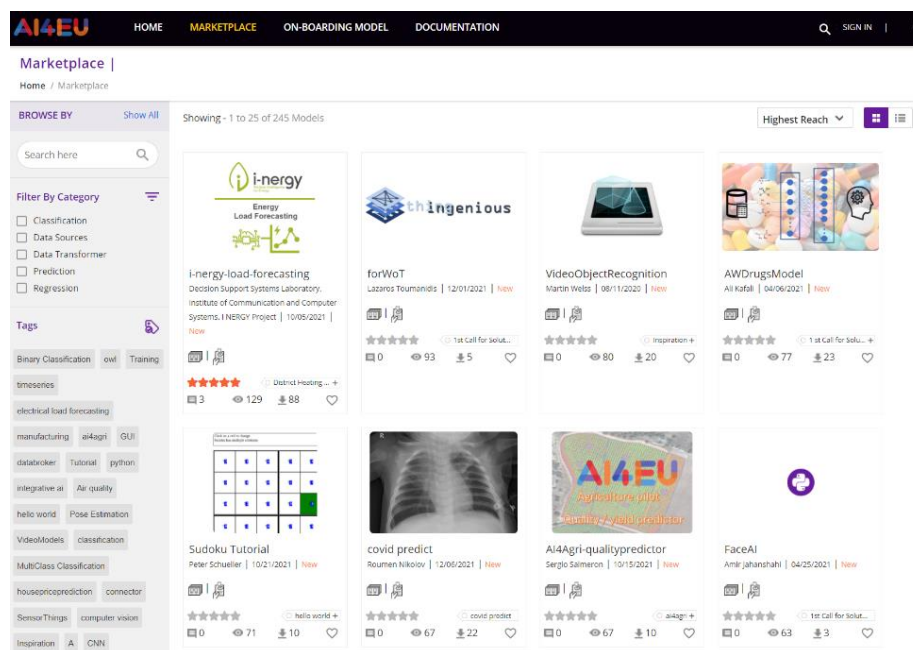


*Figure 2: Screenshot of the AI4EU Experiments platform.*

---

[13] European Commission, Staff working document on data spaces (23 Feb. 2022): https://digital-strategy.ec.europa.eu/en/library/staff-working-document-data-spaces (pp. 36-37)
[14] AI4EU Experiments Platform: https://aiexp.ai4europe.eu/
[15] AI4EU - The European AI on Demand Platform (funded by H2020 under grant agreement no 825619): https://www.ai4europe.eu

# AI4media

ARTIFICIAL INTELLIGENCE FOR
THE MEDIA AND SOCIETY

info@ai4media.eu          www.ai4media.eu