

---

# Decentralized Technologies for AI Hubs

---

**Richard Blythman**  
Algovera  
Ireland  
richard@algovera.ai

**Mohamed Arshath**  
Algovera  
Malaysia  
arshy@algovera.ai

**Salvatore Vivona**  
Algovera  
Canada  
sal@algovera.ai

**Jakub Smékal**  
Algovera  
United Kingdom  
jakub@algovera.ai

**Hithesh Shaji**  
Algovera  
Ireland  
hithesh@algovera.ai

## Abstract

1 AI requires heavy amounts of storage and compute with assets that are commonly  
2 stored in AI Hubs. AI Hubs have contributed significantly to the democratization  
3 of AI. However, existing implementations are associated with certain benefits and  
4 limitations that stem from the underlying infrastructure and governance systems  
5 with which they are built. These limitations include high costs, lack of monetization  
6 and reward, lack of control and difficulty of reproducibility. In the current work,  
7 we explore the potential of decentralized technologies - such as Web3 wallets,  
8 peer-to-peer marketplaces, storage and compute, and DAOs - to address some  
9 of these issues. We suggest that these infrastructural components can be used in  
10 combination in the design and construction of decentralized AI Hubs.

## 11 1 Introduction

12 The field of deep learning is powered by assets such as datasets, models and software, which require  
13 heavy amounts of storage and compute [15]. As a result, data scientists are regular users of AI Hubs  
14 such as GitHub, Kaggle, HuggingFace Hub and ActiveLoop Hub to provide a place for assets to be  
15 stored, shared and further developed. AI Hubs have been a significant factor in democratizing  
16 access to state-of-the-art pretrained models [16] and the contribution of open source to the field of AI  
17 [9].

18 At the same time, existing AI Hubs make certain trade-offs that arise from their underlying technolo-  
19 gies and governance structures [7]. Today's AI Hubs tend to rely exclusively on centralised cloud  
20 services such as AWS, GCP, and Azure, and the high expense of these services is often passed on to  
21 the user. Furthermore, while the assets themselves may be open source, the platform itself is typically  
22 closed source and governed by a centralized entity. The platform ultimately controls the accessibility  
23 of uploaded assets, and monetizes the network effects of user contributions without sharing in the  
24 rewards. Finally, the assets within AI Hubs are often isolated from compute and not well integrated  
25 with workflows, making reproducibility difficult.

26 Decentralized technologies such as peer-to-peer storage, compute and marketplaces, machine learning  
27 frameworks and decentralized autonomous organizations (DAOs) present opportunities for tackling  
28 the above challenges. We explore the benefits offered by these technologies to address some of the  
29 above issues within decentralized AI hubs, which offer an alternative value proposition compared  
30 with existing solutions.

## 2 Limitations of Existing AI Hubs

An AI Hub is a platform that allows data scientists, engineers and other stakeholders to share and discover AI assets such as datasets, models, apps, notebooks, pipelines and other software. AI Hubs require different features and infrastructure such as storage and compute, and may also provide organisational tools on top.

There are a number of existing hubs such as GitHub, Kaggle, HuggingFace (HF) and ActiveLoop with different features as summarized in Table 1. GitHub, owned by Microsoft, is a popular platform for storing software assets. HF achieved success by standardising the code for architectural components and providing a unified API for the popular transformer model for natural language processing (NLP) use cases [16]. As well as storing code, HF Hub provides storage for datasets and models, and compute for demos and inference. ActiveLoop Hub focuses on efficient cloud streaming of datasets for deep learning. Replit is an online integrated development environment (IDE). Unlike GitHub and HuggingFace, where modifying assets requires a separate IDE and command line, Replit users can interact with code and source control for their project through a web-based graphical user interface. Replit provides a shared compute engine that provides collaborative coding similar to Google Docs, where code can be run and displayed to multiple users. However, GPU support has not yet been released. Furthermore, file storage is limited to 0.5 GB for free users and 5 GB for paid users, which is too small for most ML assets. In general, existing AI Hubs are built using centralized infrastructure, which have certain benefits and limitations. Replit has other features such as AI-assisted tools for software development, such as co-pilot and live chat and in-line threads for discussions around code by users.

Table 1: Existing AI Hubs

Existing AI Hub	GitHub	Kaggle	Huggingface	ActiveLoop	Replit
Launch	2008	2010	2016	2018	2016
Users	SWEs	Data Scientists	Data Scientists	Data Scientists	SWEs
Monetization	No	Prizes	No	No	No
Storage/Asset	Code	Code, Datasets, Notebooks	Code, Datasets, Models, Apps	Datasets	Code
Compute/Hosting	No	GPU (Notebooks)	GPU (Inference, Apps)	No	CPU
Cloud Infrastructure	Centralized	Centralized	Centralized	Centralized	Centralized
Governance	Centralized	Centralized	Centralized	Centralized	Centralized

### 2.1 High Storage and Compute Costs

The computational cost of AI research is increasing exponentially, creating to higher barriers to entry for participants [15]. As a result, cloud services, such as storage and compute, are a significant expense for AI startups. Currently, three companies make up approximately two thirds of the market share of cloud service [8]. More than half of Amazon’s profits has come from Amazon Web Services, and 20% of AWS customers deliver 80% of revenue with the widest margins come from small and medium-sized customers [12]. Popular AI Hubs like GitHub, HuggingFace, ActiveLoop and Replit rely exclusively on centralised cloud platforms.

### 2.2 Lack of Monetization and Reward

There are few online platforms where data scientists can perform paid work independently [7]. Within existing AI Hubs, money only travels from the user to the platform itself. While AI Hubs like HuggingFace do offer free services and contribute to open source development, they also charge users for premium services that are not open source. In contrast, all contributions by users must be open source, with no ability to offer paid services.

Open source tools and libraries are widely used by commercial platforms and products within software development and AI [9], although the contributions are not typically rewarded. Platforms invite assets to be uploaded by users, but do not share any generated revenue or platform ownership with users, even when directly monetizing their contributions. For example, GitHub Copilot is a commercial product for code generation that uses a model trained on user-contributed code. HuggingFace’s paid inference API can be used to accelerate the deployment of user-contributed models.

## 72 **2.3 Lack of Control**

73 Generally, software developers and data scientists do not have full control and autonomy with their  
74 creations on centralized platforms. In one case, GitHub reverted malicious changes (and suspended  
75 the account) of a developer to their own popular open source library, raising questions around the  
76 rights of developers to do what they wish with their code [14]. In the field of AI, there has been  
77 an ongoing discussion on whether open sourcing disruptive models should be commonplace, since  
78 there is the potential for harm and bias. For example, AlphaFold can be used for discovery of novel  
79 toxic molecules. Language models can be trained on abusive content and used by online bots. Large  
80 models that are trained on the corpus of internet data reproduce bias within generated text and images.  
81 As a result, platforms like HuggingFace have come under pressure to gate or remove access to models.  
82 On the other hand, it can be argued that open sourcing the model puts the technology in the hands  
83 of more people that can study and solve issues around safety and bias. In other words, there is an  
84 orthogonal risk involved with centralization of AI in the hands of a few. Keeping models closed  
85 source effectively turns large tech companies into gatekeepers, who may not always be relied upon to  
86 adjudicate on disputes in an unbiased manner.

87 Finally, it is difficult for owners to manage fine-grained access to assets. Traditional access tokens  
88 like OAuth 2.0 [5] and API keys for datasets and models can be widely shared, and licenses for  
89 datasets and software are often broken. While possible to keep repositories private, this is often a  
90 paid feature and the encryption key is held by the platform rather than the user.

## 91 **2.4 Difficult to Reproduce**

92 The limitations of existing hubs such as GitHub for AI can make reproducibility more difficult. For  
93 example, academic papers often contain links to AI Hubs for the purpose of reproducibility. This  
94 may include code on GitHub, and datasets and model weights stored on the cloud. Nonetheless,  
95 reproducing experiments is difficult and can require many steps such as downloading datasets, running  
96 processing scripts and installing environments. This issue results from a variety of factors such as the  
97 lack of standardisation and interoperability of in the format of assets (such as dataset and code), and  
98 the decoupling of assets from compute environments and infrastructure needed to operate on them.  
99 Some of these issues can be resolved by using containers and notebooks to replicate environments  
100 and bring compute to code. HuggingFace Hub uses Gradio and Streamlit apps. Replit integrates code  
101 repositories with compute environments, but has limited storage for assets such as datasets and model  
102 weights.

## 103 **3 Decentralized Technologies for AI Hubs**

104 Decentralized technologies - such as Web3 payments, wallets, marketplaces, storage and compute,  
105 learning frameworks and DAOs - have the potential to alleviate some of the limitations of existing  
106 AI Hubs discussed above. Examples of projects working on these individual projects are shown in  
107 Figure 1.

### 108 **3.1 Payments**

109 There are few options for AI workers to monetize their creations and rewards generated by their  
110 contributions are often not shared, as discussed in Section 2.2. We believe that building in mechanisms  
111 for monetization and ownership by users would create a healthier and more sustainable ecosystem  
112 and economy. This can be achieved using cryptocurrencies (such as Bitcoin, Ethereum, Polygon,  
113 Ocean and Filecoin) and stablecoins (such as DAI or USDC), which can be used for micro- and  
114 streaming payments to stakeholders such as data scientists, data providers and compute providers  
115 with low transaction fees. Thus, integrated payments offer many opportunities for use with machine  
116 learning frameworks such as active learning and data crowdsourcing.

### 117 **3.2 Web3 Wallets**

118 As discussed in Section 2.3, data scientists typically do not have control of what they create online.  
119 Even if a repository containing assets is private, the platform holds the private keys. A Web3 wallet  
120 can be used to put the user in control of their private keys. The word wallet tends to have financial

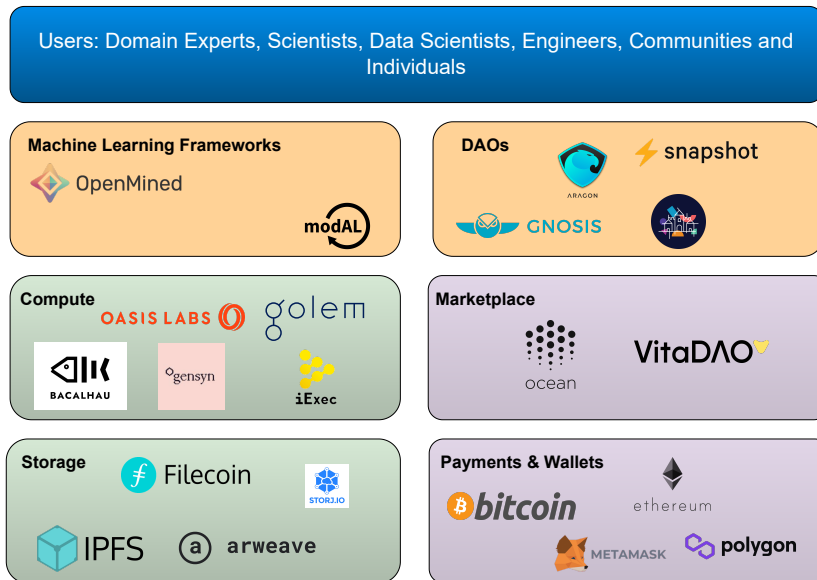


Figure 1: The decentralized AI stack (or Web3 AI stack), consisting of decentralized technologies that offer opportunities for decentralized AI Hubs. The users of decentralized AI Hubs are the many stakeholders required for undertake successful projects.

121 connotations. However, wallets are often used in the real world as a place where you hold ownership  
 122 and identity documents (such as a driver's license). Similarly, Web3 wallets can be used for ownership  
 123 and identity in the digital world. Wallets are interoperable in the sense that you can use the same  
 124 wallet to signify ownership of assets across many different protocols. Web3 wallets include software  
 125 wallets (such as MetaMask<sup>1</sup>) and hardware wallets (such as Trezor<sup>2</sup>).

### 126 3.3 Marketplace

127 Traditional AI Hubs and marketplaces are typically operated by a centralized entity serving as a  
 128 middle man. In the ideal scenario, the operator provides services in exchange for transaction fees and  
 129 acts as a mediator for conflict resolution between users. However, centralized hubs and marketplaces  
 130 also have the power to capture an outsized proportion of the value generated in a market economy as  
 131 network effects grow. This contributes to the issues discussed in Sections 2.2 and 2.3.

132 Using decentralized marketplaces protocols for tracking publication, ownership of (and access to)  
 133 assets has the potential to mitigate these risks. All operations are stored on an immutable public  
 134 distributed ledger such that provenance can be tracked. For AI use cases, assets can include datasets,  
 135 models, algorithms, apps, notebooks and manuscripts. Examples of decentralized marketplaces  
 136 include Ocean Protocol [13] and VitaDAO [4]. These protocols use non-fungible tokens (NFTs) to  
 137 represent ownership of the underlying intellectual property (IP), and fungible tokens to represent  
 138 access rights to assets under different types of licenses. The details of published assets (and associated  
 139 metadata) are encrypted and stored on-chain, along with access control parameters. A decentralized  
 140 identifier (DID) is issued to represent the asset's decentralized digital identity, and a DID Document  
 141 (DDO) is used to include additional information relevant to the asset.

142 Access gated by tokens on a blockchain has advantages compared to traditional access token like  
 143 OAuth 2.0, by solving the "double spend problem". They act as access tokens that can only be used  
 144 by one individual or for a period of time. If a user receives a token on a blockchain, the user can still  
 145 share it with someone else but this means the original user will no longer have access. This facilitates  
 146 more fine-grained access control by owners.

<sup>1</sup><https://metamask.io/>

<sup>2</sup><https://trezor.io/>

### 147 3.4 Storage

148 While details about the assets are stored on-chain with decentralized marketplaces, the data associated  
149 with the asset are often too large to store on chain. As discussed in Section 2.1, storage on centralized  
150 cloud providers is expensive. Furthermore, these services are less robust and more prone to censorship  
151 (see Section 2.3). Popular dataset and model hubs like HuggingFace and ActiveLoop Hub rely on  
152 centralised cloud platforms.

153 Decentralised protocols for storage have the potential to vastly reduce the costs incurred by data  
154 scientists for storing raw and processed versions of datasets and model weights. This makes it  
155 possible to download files from multiple locations that aren't managed by a single organization. The  
156 interplanetary file system (IPFS) [1] is a peer-to-peer protocol for storing and accessing data in a  
157 permissionless and censorship-resistant way. IPFS clusters enable data orchestration across swarms  
158 of IPFS peers by allocating, replicating, and tracking assets. Another important feature that IPFS  
159 offers is the ability to verify the validity of assets using Content Addressable Identifiers (CIDs), based  
160 on the content's cryptographic hash.

### 161 3.5 Compute

162 Access to compute is a necessity for AI projects, and the provision of services by a handful of central-  
163 ized companies has resulted in inflated costs (see Section 2.1). At the same time, the experiments and  
164 results of AI studies are often difficult to reproduce, as discussed in Section 2.4. While less mature  
165 than peer-to-peer storage solutions, decentralized protocols for providing compute resources aim to  
166 reduce the barrier-to-entry for compute providers and remove the centralised overheads on scaling  
167 [3]. This provides more options for end consumers, resulting in reduced cost. Compute can be run  
168 where the data is stored - called Compute over Data (CoD) by the Bacalhau project<sup>3</sup>, or Compute to  
169 Data (C2D) by Ocean Protocol - rather than transporting data to the location of the compute which  
170 is expensive. In this setting, decentralized compute infrastructure presents many opportunities for  
171 integration with privacy-preserving machine learning. Running compute jobs in a trustless setting  
172 requires verification that it was carried out correctly, which can be difficult for non-deterministic  
173 compute such as deep learning. Gensyn<sup>4</sup> have developed a novel system for providing proof under  
174 this condition.

### 175 3.6 Machine Learning Frameworks

176 Decentralizing infrastructure for storage and compute, and integrating payments has the potential to  
177 open up new use cases of AI. This require advancements in decentralized frameworks for machine  
178 learning. For example, privacy-preserving machine learning (PPML) - through libraries such as  
179 Openmined<sup>5</sup> - has the potential to unlock learning on private data such as health records and user data.  
180 Integrated payments can be used with active learning frameworks with libraries (such as modAL [2])  
181 and tools for crowdsourcing human intelligence (such as Turkit [10]).

### 182 3.7 DAOs

183 Decentralized autonomous organizations (DAOs) are systems that allow communities to coordinate  
184 and take part in self-governance, as determined by a set of self-executing rules on a blockchain  
185 [6]. DAOs have previously been suggested as governance structures for digital data trusts [11]. We  
186 suggest that DAOs can be used to (i) govern assets within AI Hubs, and (ii) to create decentralized  
187 AI Hubs that are governed by communities rather than single entities. Tools for governing assets  
188 within DAOs include multisignature wallets (such as Gnosis<sup>6</sup>) and profit-sharing mechanisms (such  
189 as Superfluid<sup>7</sup>). Multisignature wallets provide functionality for sharing ownership and control of  
190 assets with multiple individuals in teams in a trustless manner, while profit-sharing mechanisms can

---

<sup>3</sup><https://github.com/filecoin-project/bacalhau>

<sup>4</sup><https://www.gensyn.ai/>

<sup>5</sup><https://www.openmined.org/>

<sup>6</sup><https://gnosis-safe.io/>

<sup>7</sup><https://www.superfluid.finance/>

191 be used to distribute the revenue generated by assets. Tools for governing the infrastructure of AI  
192 Hubs include decentralized voting systems (such as Snapshot<sup>8</sup>).

## 193 4 Conclusion

194 In this work, we reviewed the trade-offs made by existing AI Hubs, and explored the ability of  
195 a collection of decentralized technologies to mitigate some of their limitations. Decentralized AI  
196 Hubs have the potential to reduce the barrier-to-entry for cloud infrastructure, increase monetization  
197 opportunities for independent AI teams, put ownership and control in the hands of creators, and  
198 improve reproducibility of research.

## 199 References

- 200 [1] Juan Benet. IPFS - Content addressed, versioned, p2p file system, 2014.
- 201 [2] Tivadar Danka and Peter Horvath. modAL: A modular active learning framework for Python. 2018.
- 202 [3] Ben Fielding and Harry Grieve. Gensyn: The hyperscale, cost-efficient compute protocol for the world’s  
203 deep learning models. 2022.
- 204 [4] Tyler Golato and Paul Kohlhaas. VitaDAO. 2021.
- 205 [5] Dick Hardt. The OAuth 2.0 authorization framework. Technical report, 2012.
- 206 [6] Samer Hassan and Primavera De Filippi. Decentralized autonomous organization. *Internet Policy Review*,  
207 10(2):1–10, 2021.
- 208 [7] Abhishek Kumar, Benjamin Finley, Tristan Braud, Sasu Tarkoma, and Pan Hui. Marketplace for AI models,  
209 2020.
- 210 [8] Rahul Kumar. Cloud market share 2022: An overview of growing ecosphere. <https://www.wpoven.com/blog/cloud-market-share/>, 2022. Accessed: 2022-08-30.
- 211
- 212 [9] Max Langenkamp and Daniel N Yue. How open source machine learning software shapes AI. In  
213 *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, pages 385–395, 2022.
- 214 [10] Greg Little, Lydia B Chilton, Max Goldman, and Robert C Miller. TurKit: Human computation algorithms  
215 on Mechanical Turk. In *Proceedings of the 23rd annual ACM symposium on User interface software and*  
216 *technology*, pages 57–66, 2010.
- 217 [11] Kelsie Nabben. Decentralised autonomous organisations (DAOs) as data trusts: A general-purpose data  
218 governance framework for decentralised data ownership, storage, and utilisation. *Available at SSRN*, 2021.
- 219 [12] Jordan Novet. How Amazon’s cloud business generates billions in profit. [https://www.cnbc.com/2021/](https://www.cnbc.com/2021/09/05/how-amazon-web-services-makes-money-estimated-margins-by-service.html)  
220 [09/05/how-amazon-web-services-makes-money-estimated-margins-by-service.html](https://www.cnbc.com/2021/09/05/how-amazon-web-services-makes-money-estimated-margins-by-service.html),  
221 2021. Accessed: 2022-08-30.
- 222 [13] Ocean Protocol. Tools for the Web3 data economy. 2018.
- 223 [14] Emma Roth. Open source developer corrupts widely-used libraries. [https://www.theverge.com/2022/](https://www.theverge.com/2022/1/9/22874949/developer-corrupts-open-source-libraries-projects-affected)  
224 [1/9/22874949/developer-corrupts-open-source-libraries-projects-affected](https://www.theverge.com/2022/1/9/22874949/developer-corrupts-open-source-libraries-projects-affected), 2022. Ac-  
225 cessed: 2022-09-30.
- 226 [15] Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. Green AI. *Communications of the ACM*,  
227 63(12):54–63, 2020.
- 228 [16] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric  
229 Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. HuggingFace’s transformers: State-of-the-art  
230 natural language processing. 2019.

---

<sup>8</sup><https://snapshot.org/>