



# International Journal of Artificial Intelligence and Machine Learning

Publisher's Home Page: <https://www.svedbergopen.com/>

Research Paper

Open Access

## Artificial Intelligence: The Final Frontier

Wulf Kaal<sup>1</sup>

<sup>1</sup>Professor, University of St. Thomas School of Law, 1000 LaSalle Avenue, MSL 400, Minneapolis, MN 55403, USA. E-mail: [kaal8634@stthomas.edu](mailto:kaal8634@stthomas.edu)

### Article Info

Volume 5, Issue 1, January 2025

Received : 11 December 2024

Accepted : 05 January 2025

Published : 25 January 2025

doi: [10.51483/IJAIML.5.1.2025.37-57](https://doi.org/10.51483/IJAIML.5.1.2025.37-57)

### Abstract

Contemporary Artificial Intelligence (“AI”) systems, particularly Large Language Models (“LLMs”), face an imminent shortage of high-quality, human-generated textual data, a phenomenon often termed “data exhaustion”. This article examines the limitations of existing centralized data-annotation frameworks, highlighting critical issues such as bias, high computational overhead, and insufficiently adaptive infrastructures. Current market participants-including Scale AI, Appen, CloudFactory, and others-excel at rapidly scaling annotation services yet struggle with ethical sourcing, privacy compliance, and equitable compensation. In addition, legal and regulatory concerns, exemplified by stringent mandates such as the General Data Protection Regulation (“GDPR”), constrain the free flow of data essential for advanced AI research. As a corrective measure, decentralized data production paradigms are proposed, including the adoption of smart contracts, token-based incentives, and participatory governance through Decentralized Autonomous Organizations (“DAOs”). While existing decentralized initiatives-SingularityNET, Fetch.ai, Ocean Protocol, Numeraire, and DcentAI-offer incremental innovations in reputation management and stakeholder engagement, they fail to fully address the nuanced requirements of large-scale “Mechanical Turk”-style data creation. In contrast, the author proposes a Weighted Directed Acyclic Graph (“WDAG”) governance model which provides a multi-dimensional reputation framework, facilitating real-time validation of data contributions, adaptive ethical and legal compliance, and collaborative oversight by diverse community members. Findings suggest that such WDAG-centric systems can more effectively maintain data quality, ensure ethical alignment, and incentivize broad participation, thereby mitigating the looming data shortage and expanding AI’s societal benefits. Ultimately, successful implementation requires coordinated efforts among policymakers, industry practitioners, and civil society actors to sustain both the technological and ethical integrity of AI research. By integrating WDAG-based governance with emerging decentralized solutions, the AI community may realize a more equitable, scalable, and future-ready paradigm for data provisioning.

**Keywords:** *Web3, Artificial Intelligence, Decentralization, LLMs, Data Governance, Privacy, Smart Contracts, DAO, WDAG, Ethical AI, GDPR*

© 2025 Wulf Kaal. This is an open access article under the CC BY license (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

\* Corresponding author: Wulf Kaal, Professor, University of St. Thomas School of Law, 1000 LaSalle Avenue, MSL 400, Minneapolis, MN 55403, USA. E-mail: [kaal8634@stthomas.edu](mailto:kaal8634@stthomas.edu)

## 1. Introduction

Human-generated text for training Artificial Intelligence (AI) systems, especially Large Language Models (LLMs) is running out.<sup>1</sup> More specifically, the accessible reserves of publicly available human-created text could be exhausted by 2028, given current usage trajectories.<sup>2</sup> The phenomenon—often referred to as “data exhaustion”—arises largely from the exponential increase in the size of datasets required to develop increasingly sophisticated AI models.<sup>3</sup> Researchers estimate that the total effective stock of human-generated text may currently stand at approximately 300 trillion tokens, with the range varying from 100 trillion to 1 quadrillion tokens.<sup>4</sup>

Methodologically, this conclusion stems from analyses tracing the historical growth of training sets used for LLM development and comparing these data requirements against the finite nature of internet-sourced text.<sup>5</sup> Notably, while the internet contains a vast corpus of textual material, not all content meets quality thresholds suitable for model training, given issues such as redundancy, noise, or irrelevance.<sup>6</sup> As a result, these findings have raised questions about the sustainability of current AI research practices and prompted discussions on whether the trend in data utilization could hinder future technological progress.<sup>7</sup>

In light of these limitations, it is conceivable that AI research may shift toward alternative data generation methods, including synthetic text produced by other AI models.<sup>8</sup> Nonetheless, this approach risks “model collapse,” wherein reliance on AI-generated data, in lieu of content created by human authors, threatens to degrade overall performance.<sup>9</sup> Consequently, researchers and industry practitioners have begun to investigate strategies such as more efficient learning algorithms, targeted transfer learning from data-rich domains like academic literature and legal archives, or the cultivation of novel high-quality datasets.<sup>10</sup> These approaches highlight the importance of not merely acquiring additional data but using current reserves more judiciously and effectively. Thus, while the anticipated depletion of human-generated data poses a formidable obstacle, it also serves as an impetus for innovation, potentially reshaping the field of AI research to ensure continued advancements in spite of looming data constraints.<sup>11</sup>

Against the backdrop of dwindling text reserves, this article examines the market for data production specifically tailored to LLM training, its inherent shortcomings, and prospective frameworks for extending this market through decentralized data production mechanisms modeled on the “Mechanical Turk” concept.<sup>12</sup> In its current form, the data-production market largely depends on a narrow set of platforms and content aggregators that host most publicly available text; however, the limited incentives and lack of direct compensation structures for content creators present significant obstacles to scaling and sustaining high-quality data supply. Moreover, issues such as privacy, copyright, and the uneven global distribution of digital connectivity further exacerbate these shortcomings by curtailing the availability and diversity of human-generated textual content.

A decentralized data production model, akin to a “Mechanical Turk” design, proposes a system in which individual contributors are compensated for generating, refining, or annotating text data. Under such a model,

<sup>1</sup> AI Firms will Soon Exhaust Most of the Internet’s Data. *The Economist* (July 23, 2024). <https://www.economist.com/schools-brief/2024/07/23/ai-firms-will-soon-exhaust-most-of-the-internets-data>.

<sup>2</sup> *Id.* (suggesting with a 95% confidence interval, that human generated data for LLM model training will run out by 2028).

<sup>3</sup> Pablo Villalobos *et al.* (2024), Will We Run Out of Data? Limits of LLM Scaling Based on Human-Generated Data. *Epoch AI Blog* (June 6). <https://epoch.ai/blog/will-we-run-out-of-data-limits-of-llm-scaling-based-on-human-generated-data#:~:text=Our%2080%25%20confidence%20interval%20is,point%20between%202026%20and%202032>

<sup>4</sup> *Id.*

<sup>5</sup> Pablo Villalobos *et al.* (2024), Will We Run Out of Data? Limits of LLM Scaling Based on Human-Generated Data. *arXiv* (Preprint No. 2211.04325v2) (June 4). <https://arxiv.org/html/2211.04325v2>

<sup>6</sup> *Id.*

<sup>7</sup> AI Firms will Soon Exhaust Most of the Internet’s Data, *supra* note 2.

<sup>8</sup> Pablo Villalobos *et al.*, *supra* note 4.

<sup>9</sup> The Breakthrough AI Needs, *The Economist* (Sept 19, 2024). <https://www.economist.com/leaders/2024/09/19/the-breakthrough-ai-needs>

<sup>10</sup> How AI Models are Getting Smarter. *The Economist* (Aug 6, 2024). <https://www.economist.com/schools-brief/2024/08/06/how-ai-models-are-getting-smarter>

<sup>11</sup> *Id.*

<sup>12</sup> *Cf.* Amazon Mechanical Turk, [www.mturk.com](http://www.mturk.com) (last visited Jan 5, 2025).

participants from diverse linguistic and cultural backgrounds could collectively bolster the text repositories needed to train robust LLMs. This approach not only distributes the cost and labor of data creation and curation more equitably but also has the potential to enrich AI systems with heterogeneous linguistic structures and cultural nuances. Mechanisms like smart contracts or blockchain-based ledgers might serve as the infrastructure to manage contributions, ensure fair compensation, and maintain high-quality data standards.<sup>13</sup>

The promise of decentralized solutions lies in their ability to aggregate widely dispersed knowledge while minimizing the monopolistic control that might restrict the creation of novel or high-grade text resources. Nonetheless, their success critically depends on addressing challenges such as fraudulent submissions, content moderation, and ensuring alignment with ethical and legal frameworks. If implemented effectively, decentralized and incentive-driven data generation may offer a viable path toward alleviating data scarcity by creating a self-sustaining market that continually produces new, high-value textual content. Taken together, these considerations underscore the urgency of reimagining how data for LLM training is collected and valued, an endeavor that this article explores in detail.

## 2. Limitations of AI Development

The limitations of AI development—ranging from data quality and bias issues to technical and legal hurdles—underscore the intricate ecosystem in which modern AI systems are built and deployed. Addressing these challenges requires a concerted, interdisciplinary effort that involves computer scientists, legal experts, developers, and domain specialists. By recognizing the multifaceted nature of AI limitations, stakeholders can implement balanced solutions that preserve data integrity, safeguard personal and proprietary information, and ensure the longevity of AI innovations in both private and public sectors.

The limitations of AI—ranging from biased datasets to technical constraints—inhibit the technology's ability to deliver transformative benefits across sectors such as healthcare, finance, and disaster response. Biased or outdated training data can perpetuate inequities and undermine trust in AI-driven tools, while reliance on AI-generated data risks "model collapse," thereby eroding confidence in technological solutions. Further, stringent privacy regulations and inadequate infrastructure can hinder researchers' access to the rich, real-time data critical for accurate insights. These shortcomings diminish AI's utility for tackling global challenges, such as climate change and public health crises, and hamper efforts to shape ethically sound and equitable AI governance.

Moreover, the high costs associated with processing colossal datasets restrict such innovation to a handful of well-resourced institutions. This concentration of capability threatens to undermine the democratizing potential of AI by limiting widespread adoption and participation in research and development. Consequently, unless efforts are made to address biases, ensure ethical data stewardship, and foster robust infrastructures, AI may fail to serve as a unifying force for societal progress. Instead, it risks exacerbating existing inequalities, thereby impeding its capacity to fulfill the broad vision of improving human welfare on a global scale.

### 2.1. Technical Limitations

#### 2.1.1. Data Quality and Bias

AI models trained on publicly available data frequently mirror inherent biases present in the underlying datasets.<sup>14</sup> These biases can stem from demographic, cultural, and historical factors woven into internet-sourced text. For instance, if a dataset disproportionately represents a particular region or demographic group, the model may offer skewed performance, demonstrating suboptimal or inappropriate outputs when deployed in unfamiliar settings. Such limitations not only pose ethical dilemmas regarding equity and fairness in AI but also jeopardize the robustness of AI applications in real-world contexts. Mitigating bias requires proactive interventions, such as targeted sampling of underrepresented data sources, dynamic rebalancing of training sets, or post-training calibration procedures.<sup>15</sup> Failure to address data quality and bias risks perpetuating harmful stereotypes and undermines public trust in AI technologies.

---

<sup>13</sup> *Id.*

<sup>14</sup> Artificial Intelligence for Research and Scholarship. *Harv. Libr. Rsch. Guides*. <https://guides.library.harvard.edu/airesearch> (Nov. 8, 2024, 12:19 PM).

<sup>15</sup> *Id.*

### 2.1.2. Volume and Variety of Data

Although AI models thrive on large-scale datasets, they are vulnerable to “model collapse” if increasingly trained on synthetic data generated by other AI systems.<sup>16</sup> As AI-generated content proliferates online, it dilutes the overall diversity and originality of text available for subsequent training processes, potentially leading to a degradation in model performance over repeated generations. This phenomenon underscores the interdependence of data quality and quantity in sustaining model accuracy.<sup>17</sup> Consequently, balancing real and AI-generated content is critical to preserving a breadth of linguistic and conceptual variety that ensures robust, generalizable models.

### 2.1.3. Privacy and Security

The widespread reliance on publicly accessible data also triggers substantial privacy and legal concerns, as personal or sensitive information can inadvertently be included in training corpora. Such data usage may contravene privacy regulations like the European Union’s General Data Protection Regulation (GDPR) or risk running afoul of the emerging AI Act in Europe.<sup>18</sup> In addition, aggregated data harvested from multiple sources can encompass copyrighted material or trade secrets, presenting a legal risk that is not trivial to resolve. Compliance with GDPR, the AI Act, and other jurisdiction-specific legal frameworks necessitates rigorous data governance protocols, which may include anonymization, data minimization, and explicit consent where applicable. These measures can, however, complicate the collection and curation of robust datasets for AI training.

### 2.1.4. Real Time Data Challenges

The use of real-time data further compounds the complexity of AI model development. Because such data streams demand immediate processing at high velocity and volume, models must be consistently updated to maintain accuracy.<sup>19</sup> In sectors like finance and healthcare, the ability to ingest and analyze real-time data can dramatically influence decision-making and outcomes. Yet integrating dynamic data feeds into AI pipelines requires substantial infrastructure and robust system design. Latency, data throughput constraints, and the absence of fully automated mechanisms for continuous training and validation can collectively undermine the timeliness of insights derived from AI models.<sup>20</sup>

### 2.1.5. Data Timeliness and Relevance

Overreliance on static or outdated data remains a pressing limitation in AI training. In rapidly evolving fields—technology, law, or public policy—models trained on obsolete datasets fail to capture new trends, behaviors, or regulatory changes, ultimately reducing their predictive and explanatory power.<sup>21</sup> Consequently, AI systems may present recommendations or analyses that are no longer valid, an issue particularly acute where continuous updates are either not feasible or not prioritized. Scholars and practitioners increasingly advocate for iterative training regimes, incorporating fresh data in real time or near real time, to ensure models remain aligned with current realities.<sup>22</sup>

### 2.1.6. Technical Limitations

Finally, the vast quantity of publicly available data—despite its seeming abundance—can itself pose formidable

<sup>16</sup> Iliia Shumailov *et al.* (2023), The Curse of Recursion: Training on Generated Data Makes Models Forget. *arXiv* (Preprint arXiv:2305.17493v2) (May 31). <https://arxiv.org/abs/2305.17493v2>

<sup>17</sup> *Id.*

<sup>18</sup> Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation), 2016 O.J. (L 119) 1; Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonized rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act), 2024 O.J. (L series).

<sup>19</sup> Monika Steidl *et al.* (2023), The Pipeline for the Continuous Development of Artificial Intelligence Models—Current State of Research and Practice. *J. Sys. & Software*, May, at 1.

<sup>20</sup> *Id.*

<sup>21</sup> Michael R. Ortizo (2023), The Impending Death of Manual Bluebook Citations: How AI will Reshape Legal Education. *Syracuse J. Sci. & Tech. L. Blog* (Apr 17). <https://jost.syr.edu/the-impending-death-of-manual-bluebook-citations-how-ai-will-reshape-legal-education>

<sup>22</sup> *Id.*

computational challenges. Preprocessing massive datasets requires substantial memory, storage, and processing power.<sup>23</sup> Additionally, not all data will meet the stringent quality or format requirements essential for effective model training, prompting the need for sophisticated cleaning and transformation pipelines. Problems such as text duplication, multilingual inconsistencies, and unstructured data formatting necessitate more than cursory oversight, especially if models are to be deployed and improved on an ongoing basis. The interplay between voluminous data and the high computational overhead needed to process it thus underscores the importance of optimized data strategies and infrastructure investments.

## 2.2. Implications for Society and the Future of Humanity

If the limitations outlined above are not addressed, the transformative potential of AI—particularly in fields such as healthcare, education, environmental preservation, and economic development—may remain unrealized. AI systems that fail to overcome data bias, technical capacity constraints, and legal and ethical quandaries will likely yield diminished benefits or, worse, generate harmful outcomes.

In healthcare, for instance, biased or stale training data could lead to algorithms that misdiagnose underrepresented populations, reinforcing existing health disparities rather than alleviating them.<sup>24</sup> Furthermore, the risk of “model collapse” posed by reliance on AI-generated data not only jeopardizes research and innovation but also disrupts the broader public’s trust in technology—a critical resource in driving AI-assisted solutions for global challenges.<sup>25</sup>

In addressing climate change or other critical issues, AI’s capacity for large-scale modeling and simulation hinges on access to timely, diverse, and high-quality data. Should excessive legal restrictions or privacy regulations unduly limit data availability—without the counterbalancing adaptation of privacy-preserving techniques—the beneficial insights derived from such analysis may dwindle.<sup>26</sup> If compliance with data protection frameworks is conducted in a manner that overly constrains researchers’ access to relevant data, the broader societal gains of AI risk being offset by a stunted innovation ecosystem.

Likewise, real-time decision-making in finance, disaster response, and pandemic modeling depends on the ability to ingest, process, and analyze live data streams. If AI systems cannot incorporate dynamic data effectively—due to infrastructure limitations or prohibitive costs—the timeliness and accuracy of outcomes will degrade, undermining society’s capacity to adapt to rapidly-evolving events.<sup>27</sup> This shortfall resonates beyond immediate financial or logistical concerns; it also has humanitarian and ethical implications, as real-time AI systems increasingly shape resource distribution, crisis management, and public policy.

Moreover, the negative impact of outdated or static data extends beyond mere inefficiency. In domains like legal informatics or policy reform, AI-driven tools should ideally respond to new regulations and social norms in near real time. When reliant on defunct information, AI cannot identify evolving patterns or adapt to newly institute legal frameworks, leading to advice or decisions that are patently obsolete.<sup>28</sup> As a result, the broader project of harnessing AI for societal benefit becomes prone to stagnation, dampening the hope that AI can expedite equitable and forward-thinking governance.

Finally, the high computational burden inherent in processing enormous, heterogeneous datasets intensifies these challenges by increasing the resource divide between well-funded organizations and smaller public-interest entities. If only a handful of technology conglomerates can afford the infrastructure to glean insights from massive troves of data, innovation risks becoming concentrated in the hands of a select few.<sup>29</sup> Consequently, the democratizing promise of AI may falter, thwarting broader societal aspirations for inclusive economic development and collective problem-solving.

---

<sup>23</sup> Create a Document Processing Custom Model. *Microsoft Learn* (Dec. 12, 2024). <https://learn.microsoft.com/en-us/ai-builder/create-form-processing-model>

<sup>24</sup> Besse *et al.*, *supra* note 15.

<sup>25</sup> Shumailov *et al.*, *supra* note 17.

<sup>26</sup> Steidl *et al.*, *supra* note 20.

<sup>27</sup> Steidl *et al.*, *supra* note 20.

<sup>28</sup> Ortizo, *supra* note 22.

<sup>29</sup> Microsoft Learn, *supra* note 24.



Ultimately, each of these constraints—data bias, privacy concerns, real-time and up-to-date information requirements, and computational hurdles—illustrates why AI's potential to “heal society” and pioneer groundbreaking innovations is not guaranteed. Substantial, multidisciplinary efforts are necessary to foster systems that can continually learn from diverse, representative, and ethically sourced data in a legally compliant and computationally feasible manner. Failure to orchestrate these efforts not only risks squandering the promise of AI in tackling some of humanity's most pressing challenges but could also exacerbate existing inequalities and injustices. Addressing these limitations with resilience, foresight, and inclusivity may therefore be essential for ensuring AI's pivotal role in advancing human welfare over the coming decades.

### 3. Centralized AI Data Production

Startups are increasingly important in addressing the complex challenges inherent in AI development, including issues of data bias, data quality, real-time data processing, and compliance with regulatory frameworks. Certain enterprises specialize in the curation and refinement of datasets, offering AI developers access to high-quality, bias-reduced data sourced from a broad spectrum of demographic and geographic origins. These data-providing startups employ sophisticated sampling and validation techniques, aligning with the advocacy for more equitable and ethically sourced training data.<sup>30</sup> By incorporating perspectives from underrepresented groups, these companies strive to mitigate systemic biases in data sampling, thereby preventing the reinforcement of detrimental stereotypes and promoting the development of robust, fair AI systems.

Simultaneously, other startups are innovating in the realm of real-time data processing and computational challenges through advanced infrastructural solutions. This includes the development of edge computing systems designed to reduce latency in handling vast data streams, addressing the accuracy issues.<sup>31</sup> Additionally, companies focusing on privacy-preserving technologies are adopting methodologies like federated learning and homomorphic encryption, facilitating adherence to stringent data protection laws such as the General Data Protection Regulation (GDPR) and the AI Act.<sup>32</sup> Through decentralizing data processing across secure nodes, these startups help circumvent the privacy and security issues associated with centralized data management, thus reducing the risk of data breaches.

Collectively, these entrepreneurial efforts encapsulate a multidisciplinary approach combining data science, technological innovation, and regulatory compliance, positioning themselves as the vanguard of developing AI solutions that are both technologically advanced and socially responsible.

#### 3.1. Market Overview

The production of data for AI constitutes a critical pillar of AI-driven technological advancement, centering on the creation, annotation, and management of datasets required to train AI models across various domains. Recent developments in autonomous vehicles, healthcare, finance, and other data-intensive sectors have amplified the demand for large volumes of high-quality, diverse data. This burgeoning market can be broadly categorized into four main areas of focus:

1. Data labeling and annotation, essential for supervised learning models;
2. Data curation and management, which includes cleaning and normalization processes to ensure data quality;
3. Synthetic data generation, employed to compensate for sensitive or sparsely available real-world data; and
4. Data privacy and compliance solutions, which cater to regulations such as the GDPR.<sup>33</sup>

Within this ecosystem, data labeling and annotation stand out as indispensable to the development of supervised learning models. AI model performance hinges on meticulously labeled data that accurately captures relevant features and classifications. Additionally, data curation and management strategies minimize noise and inconsistencies, thereby enhancing the utility of raw datasets. Synthetic data generation methods further

<sup>30</sup> Michael Kearns and Aaron Roth (2020), Ethical Algorithm Design. *ACM SIGecom Exch.*, Dec., at 31.

<sup>31</sup> Pedro Garcia Lopez *et al.* (2015), Edge-Centric Computing: Vision and Challenges. *ACM SIGCOMM Comput. Comm'n Rev.*, Oct., at 37.

<sup>32</sup> Michèle Finck (2019), Smart Contracts as a Form of Solely Automated Processing Under the GDPR, *Int'l Data Priv. L.*, 9, 78.

<sup>33</sup> Regulation (EU) 2016/679 of the European Parliament and of the Council, 2016 O.J. (L 119) 1.

diversify and augment real-world data, helping organizations avoid privacy pitfalls or logistical constraints tied to sensitive information. Meanwhile, data privacy and compliance frameworks ensure that these processes align with legally mandated standards. Thus, safeguarding individuals' rights and mitigating corporate liability. Against this backdrop, the market landscape features several key players, each offering distinctive approaches to the challenges of data production, governance, and innovation.

### 3.2. Market Participants

An understanding of the contemporary market for AI training data—particularly as it evolves into 2025 and beyond—necessitates a closer examination of the key companies currently shaping this domain. While the AI ecosystem is vast, these specific market players exemplify the diverse approaches and innovations critical to meeting rapidly escalating demands for high-quality, large-scale datasets. Their distinct business models range from sophisticated annotation platforms to ethically driven labor-sourcing strategies, providing a microcosm of how this sector is striving to address both technical and social imperatives. Furthermore, these companies occupy pivotal positions in the broader AI supply chain, supporting the development of cutting-edge applications in sectors as varied as autonomous driving, healthcare, finance, and robotics.

Taken together, Scale AI, Appen, Labelbox, CloudFactory, Sama, Hive, and V7 Labs illustrate the multifaceted priorities that will shape the future of AI training data. For instance, the insistence on swift, accurate annotation is vital to enabling real-time applications,<sup>34</sup> while ethical and impact sourcing strategies highlight ongoing commitments to equitable labor practices.<sup>35</sup> Efforts to automate workflow processes and integrate advanced analytics further underscores the forward-thinking nature of this market.<sup>36</sup> As AI continues to permeate more sectors and the need for reliable, representative data expands, these companies' services and innovations will remain at the forefront, signaling important directions for the industry's ongoing development.

#### 3.2.1. Scale AI

**Origin:** Founded in 2016 by Alexandr Wang and Lucy Guo in San Francisco, California, Scale AI initially concentrated on data annotation services before expanding its portfolio.

**Market Position:** The company's capacity for handling large-scale, intricate datasets—particularly in autonomous driving—has positioned it as a leader in the AI data annotation sphere.<sup>37</sup> Further strengthening its market hold, Scale AI offers comprehensive data management tools, including "Nucleus" and "Launch," enabling end-to-end solutions for training, validating, and deploying AI models.

#### 3.2.2. Appen

**Origin:** Established in 1996 in Sydney, Australia, Appen now maintains a significant presence in the United States.

**Market Position:** Appen has amassed an extensive client base that includes major technology firms such as Microsoft and Amazon. Its longstanding experience in the field, coupled with its capacity to provide both automated and human-powered data annotation services, secures its relevance across multiple linguistic and cultural contexts.<sup>38</sup> This versatility allows Appen to excel in projects demanding scale, accuracy, and complexity.

#### 3.2.3. Labelbox

**Origin:** Founded in 2018 in San Francisco, California, Labelbox entered the market with a focus on simplifying and accelerating data labeling.

<sup>34</sup> Jan-Christoph Klie *et al.* (2024), Analyzing Dataset Annotation Quality Management in the Wild. *Computational Linguistics*, 50, 817.

<sup>35</sup> The Sama Team, Responsible AI in the Age of Generative Models. *Sama*. <https://www.sama.com/ebook/responsible-ai-in-the-age-of-generative-models>, (last visited Jan. 5, 2025).

<sup>36</sup> Labelbox, The Labelbox Guide to Labeling Automation. *Labelbox*. <https://labelbox.com/learn/library/Labeling-Automation-Guide/>, (last visited Jan. 5, 2025).

<sup>37</sup> Marc Vartabedian and Belle Lin (2024), Can 1\$ Billion Turn Startup Scale AI into an AI Data Juggernaut?. *Wall Street J.* (June 28). [https://www.wsj.com/articles/can-1-billion-turn-startup-scale-ai-into-an-ai-data-juggernaut-2417ccc2?st=QXVcb9&reflink=desktopwebshare\\_permalink](https://www.wsj.com/articles/can-1-billion-turn-startup-scale-ai-into-an-ai-data-juggernaut-2417ccc2?st=QXVcb9&reflink=desktopwebshare_permalink); Berber Jin (2024), The 27-Year-Old Billionaire Whose Army Does AI's Dirty Work. *Wall Street J.* (Sep 20). [https://www.wsj.com/tech/ai/alexandr-wang-scale-ai-d7c6efd7?st=bhTC1n&reflink=desktopwebshare\\_permalink](https://www.wsj.com/tech/ai/alexandr-wang-scale-ai-d7c6efd7?st=bhTC1n&reflink=desktopwebshare_permalink)

<sup>38</sup> Appen, <https://www.appen.com/> (last visited Jan. 5, 2025).

**Market Position:** Renowned for its user-friendly platform tailored to a variety of industry-specific needs—particularly computer vision—the company caters to both smaller teams and large enterprises.<sup>39</sup> Its enterprise-oriented approach includes features for collaborative workflows, advanced analytics, and integration with existing data pipelines.

#### 3.2.4. CloudFactory

**Origin:** Launched in 2010 with headquarters in Kathmandu, Nepal, CloudFactory has since established operations in the United States.

**Market Position:** CloudFactory's market edge lies in its commitment to "impact sourcing" and ethical labor practices, recruiting and training workers in developing regions to perform data annotation tasks.<sup>40</sup> This model not only addresses the need for reliable data annotation services but also contributes to broader social objectives, particularly within the healthcare and agriculture industries.

#### 3.2.5. Sama

**Origin:** Established in 2008 in San Francisco, California, with prominent operations in Africa, Sama has forged a unique path in providing ethical data annotation services.

**Market Position:** By centering its workforce in underserved communities, Sama ensures high-quality data annotation while supporting socio-economic improvements. Its computer vision services, specifically deployed in tech and automotive sectors, underscore its dedication to efficient labeling workflows combined with social impact.<sup>41</sup>

#### 3.2.6. Hive

**Origin:** Founded in 2014 in San Francisco, California, Hive emphasizes speed and accuracy in AI-driven data annotation solutions.

**Market Position:** While it serves several verticals, Hive's particular focus on autonomous driving and robotics allows it to showcase specialized expertise. By leveraging automation tools and an expansive annotation workforce, Hive delivers rapid turnarounds without compromising on data precision.<sup>42</sup>

#### 3.2.7. V7 Labs

**Origin:** Launched in 2018 in London, UK, V7 Labs focuses on end-to-end annotation automation for images and videos.

**Market Position:** Particularly well-known in the medical imaging domain, V7 Labs' technology-driven approach emphasizes workflow optimization and machine learning-assisted labeling to cut costs and reduce labor intensity.<sup>43</sup> The platform thus appeals to organizations seeking to harness state-of-the-art automation without sacrificing overall data quality.

### 3.3. Market Dynamics

The AI data production market is marked by several evolving trends. First, competition on quality and speed

<sup>39</sup> Josh Constine (2018), A Pickaxe for the AI Gold Rush, Labelbox Sells Training Data Software. *Tech Crunch* (July 30). <https://techcrunch.com/2018/07/30/labelbox/>; Smart One. Ai, *Getting Started with LabelBox (A Comprehensive Guide)*, Smart One. Ai Blog (Oct 26, 2023). <https://smartone.syspark.net/blog/getting-started-with-labelbox-comprehensive-guide/>

<sup>40</sup> James Muldoon *et al.* (2023), The Poverty of Ethical AI: Impact Sourcing and AI Supply Chains, *AI & Soc'y*. <https://doi.org/10.1007/s00146-023-01824-9>; Humans in the Loop, Data Labeling Companies with a Social Impact. *Humans in the Loop*, 10. <https://humansintheloop.org/data-labeling-companies-with-a-social-impact/> (last visited Jan. 3, 2025).

<sup>41</sup> Muldoon *et al.* (2021), *supra* note 41; Kyle Wiggers, Sama Aims to Bring Greater Equality to Crowd-Labeling of Datasets with New \$70M. *Venture Beat* (Nov 4). <https://venturebeat.com/ai/sama-aims-to-bring-greater-equality-to-crowd-labeling-of-datasets-with-new-70m/>; Forbes, Wendy Gonzalez, CEO Sama Profile, Forbes: Bus. Council. <https://councils.forbes.com/profile/Wendy-Gonzalez-CEO-Sama/94a9e932-acfb-4289-986e-4dbe69f41f75> (last visited Jan 3, 2025).

<sup>42</sup> Anna Yoo Jeong Ha *et al.* (2024), Organic or Diffused: Can We Distinguish Human Art from AI-generated Images?. *In Proc. ACM SIGSAC Conf. Comput. and Comm'n Sec.*, 4822, 4822-4836 (Dec). <https://doi.org/10.1145/3658644.3670306>; Hive (2024), Announcing Hive's Partnership with the Defense Innovation Unit, *Hive Blog* (Dec 4). <https://thehive.ai/blog/announcing-hives-partnership-with-the-defense-innovation-unit>; Hive (2024), Expanding Our CSAM Detection API. *Hive Blog* (Nov 21). <https://thehive.ai/blog/expanding-our-csam-detection-api>

<sup>43</sup> V7 Labs. <https://www.v7labs.com/>; Cong-Tai Nguyen *et al.* (2023), A Comprehensive Web Annotation Application for Organ Image Segmentation and Predictive Inference. *In Proc. 12<sup>th</sup> Int'l Symp. on Info. and Comm'n Tech.*, 850 (Dec 7). <https://doi.org/10.1145/3628797.3628955>.



incentivizes firms to refine their annotation processes to deliver faster, more accurate results. Second, innovation in automation is driving the partial or full transition from manual labeling to algorithmic solutions that minimize human intervention, thereby reducing both cost and error rates. Third, ethical and sustainable practices are becoming central, as stakeholders and regulators increasingly scrutinize how data is sourced, processed, and managed.<sup>44</sup> Finally, expansion into new verticals is augmenting the market's scope; companies once specialized in, for example, autonomous vehicles are diversifying into healthcare diagnostics, industrial robotics, and more.

As AI's influence expands into numerous fields, the demand for sophisticated, heterogeneous datasets will only intensify. Continuous learning environments, where models are routinely updated with new data, further underscore the value of reliable data production pipelines. Moreover, privacy and compliance considerations—exemplified by global regulations such as the GDPR—will likely play an increasingly influential role, necessitating careful data governance practices that align with both legal and ethical principles. Consequently, the AI data production sector is poised for robust, sustained growth as firms race to innovate and address emerging complexities in this critical domain of AI development.

### 3.4. AI Data Optimization

As AI permeates numerous sectors—ranging from healthcare and finance to autonomous vehicles—the need for high-quality, secure, and ethically sourced data has never been more critical. It is within this context that key industry players—exemplified by Scale AI—have emerged to address the dual challenges of ensuring data accuracy and maintaining compliance with evolving legal and ethical frameworks. Their methodologies incorporate human oversight to mitigate algorithmic bias, proprietary approaches for synthetic data generation, and secure systems aligned with guidelines under regulations like the European GDPR.

The demand for real-time data integration and scalable infrastructures has prompted organizations to invest in advanced data-handling solutions. These solutions not only enable continuous learning—thereby ensuring AI models remain up-to-date and contextually relevant—but also reinforce ethical practices through consent-driven data acquisition and transparent sourcing.

In the following sections, the author surveys seven key focal areas that define contemporary AI data strategies—from high-quality data curation and proprietary solutions to privacy practices, scalability, and collaborations—underscoring how companies like Scale AI, along with their counterparts, respond to escalating demands for reliable, performance-oriented AI datasets. By investing in user-centric data frameworks and forging cross-industry partnerships, these providers are poised to facilitate breakthroughs that extend far beyond traditional applications of AI technology.

- 1. High-Quality Data Curation:** Scale AI enhances data quality through its Data Engine, which involves human-in-the-loop processes to ensure accurate labeling and annotation. This approach helps in curating datasets that are not only vast but also of high quality, which is crucial for training effective AI models. By integrating human oversight, they mitigate biases and errors inherent in purely automated data collection.
- 2. Proprietary and Custom Data Solutions:** Companies like Scale AI focus on the creation of proprietary datasets tailored to specific industry needs. This includes generating synthetic data where real-world data is scarce or sensitive, ensuring that AI models can be trained on scenarios that might not be publicly available or ethical to collect. This addresses both the volume and variety limitations by providing unique, industry-specific data.
- 3. Data Privacy and Security:** Addressing privacy concerns, Scale AI and similar providers offer secure environments for data handling. They use techniques like federated learning, where models are trained across multiple decentralized devices or servers, holding the data locally, thereby reducing privacy risks associated with centralized data storage. This helps in complying with regulations like GDPR while still providing valuable data for AI training.
- 4. Scalability and Real-Time Data Integration:** To tackle the challenge of real-time data, these providers have developed systems that can scale data processing and integration. This ensures that AI models can learn

---

<sup>44</sup> Regulation (EU) 2016/679, *supra* note 34.

from the latest data streams, enhancing model relevance and timeliness. Scale AI's platforms are designed to handle large volumes of data efficiently, ensuring models can evolve with the changing data landscape.

5. **Ethical Data Practices:** Scale AI and others are developing frameworks to ensure data is ethically obtained, with considerations for consent and transparency. This includes initiatives like safety evaluations and alignment labs to ensure AI models respect ethical boundaries while training.
6. **Partnerships and Collaborations:** By partnering with various industries and academic institutions, these providers can access proprietary data that would otherwise be siloed. This collaboration not only enriches the dataset but also ensures that the data reflects practical, real-world applications, thereby enhancing model performance across diverse scenarios.
7. **Technical Infrastructure:** Investment in robust technical infrastructure ensures that the data can be processed at scale. This includes advancements in data storage, retrieval, and processing technologies to manage the vast amounts of data required for AI training, addressing technical limitations like computational resources.

### 3.5. Customer Relationships

Specialized data annotation and curation providers, such as Scale AI and its competitors provide mission critical services to AI companies. Therefore, Scale AI and its competitors are frequently retained by both well-established enterprises and nascent start-ups seeking to bolster their AI model development through enhanced data accuracy and volume. For instance, major technology firms like Microsoft and Amazon regularly engage companies such as Appen to refine and expand their AI capabilities.<sup>45</sup> Similarly, emerging ventures in areas like autonomous vehicles or health technology depend on advanced annotation services for model training, thereby broadening the market scope of these service providers.

Notably, the autonomous driving sector exemplifies this trend: entities like Scale AI and Hive are hired specifically to handle the intricate data annotation tasks necessary for training sophisticated computer vision algorithms.<sup>46</sup> Meanwhile, CloudFactory targets areas such as healthcare and agriculture, furnishing image and video annotation services that cater to the precise requirements of these specialized fields.<sup>47</sup> Enterprises with broader AI initiatives—spanning various verticals—turn to Labelbox and V7 Labs for scalable, industry-specific datasets and annotation solutions.<sup>48</sup>

Client relationships may originate through direct engagement, partnerships, or market-driven networking efforts. Frequently, organizations circulate Requests for Proposals (“RFPs”) to solicit competitive bids, enabling providers like Scale AI to demonstrate the efficacy and adaptability of their services.<sup>49</sup> Industry events, conferences, and technology expos also serve as crucial networking forums for these annotation firms to cultivate business contacts. Once potential collaborations are identified, demonstrations and trial periods allow clients to assess data quality, annotation speed, and overall compatibility with project objectives.

In their decision-making process, clients weigh several factors: they prioritize exacting standards of annotation quality and accuracy, seek providers capable of efficiently scaling to accommodate large volumes of data, and remain sensitive to cost-effectiveness.<sup>50</sup> Moreover, compliance with regulatory frameworks such as the GDPR necessitates robust data privacy and ethical sourcing protocols—elements highlighted by providers like CloudFactory and Sama. Customization, which may be indispensable for heavily regulated domains such as medical imaging, further sways client preferences toward flexible annotation platforms like V7 Labs.<sup>51</sup>

<sup>45</sup> Appen, *supra* note 39.

<sup>46</sup> Scale (2025), Automotive Data Engine. *Scale*. <https://scale.com/automotive> (last visited Jan 1); Hive.Ai, *Enterprise Data Annotation*, Hive.Ai. <https://thehive.ai/data-labeling> (last visited Jan 3, 2025).

<sup>47</sup> Cloud Factory, Explore Use Cases. *Cloud Factory*. <https://www.cloudfactory.com/use-cases> (last visited Jan 3, 2025).

<sup>48</sup> Labelbox, The Most Complete Data Labeling Solution. *Labelbox*. <https://labelbox.com/product/annotate/> (last visited Jan 3, 2025); V7, Data Labeling. Lightning-Fast, Pixel-Perfect, V7. <https://www.v7labs.com/darwin> (last visited Jan 3, 2025).

<sup>49</sup> Scale, RFP Evaluation Assistant. *Scale*. <https://scale.com/enterprise/prebuilt-applications/rfp-evaluation-assistant> (last visited Jan. 4, 2025).

<sup>50</sup> Klie *et al.*, *supra* note 35; Verified Market Research, Data Annotation Outsourcing Market. *Verified Market Research*. <https://www.verifiedmarketresearch.com/product/data-annotation-outsourcing-market> (last visited Jan 4, 2025).

<sup>51</sup> Casimir Rajnerowicz (2023), How to Do Voxel Annotations in V7. *V7 Blog* (Sep 20). <https://www.v7labs.com/blog/multi-planar-annotations>; V7, Annotate DICOM & NIFTI Files. *V7 Guides*. <https://docs.v7labs.com/docs/annotate-dicom-nifti-files> (last visited Jan 3, 2025).

Ultimately, enterprises and start-ups alike gravitate toward providers who excel in accuracy, scalability, cost-efficiency, compliance, and industry-focused customization. As a result, Scale AI and comparable firms occupy a central role in shaping the future trajectory of AI by delivering indispensable datasets that fuel ongoing innovation, expand AI's operational reach, and ensure model robustness across an ever-growing array of applications.

**Overview of AI Data Startup Client Rosters:** Below is an overview of the clients who hire AI data production companies. To preview, clients choose based on a combination of these factors, often through a process involving evaluations, trials, and negotiations to find the best fit for their AI data production needs.

- **Tech Giants and Startups:** Companies like Google, Microsoft, and Amazon<sup>52</sup> hire for data annotation to expand or refine their AI capabilities. Startups in AI, particularly in fields like autonomous vehicles or health tech, also seek these services for model training.
- **Automotive and Robotics Industries:** Entities like Scale AI and Hive are notably hired by the autonomous driving sector for complex data annotation needs.<sup>53</sup>
- **Healthcare and Agriculture:** CloudFactory targets these sectors for data annotation services, focusing on image and video data.<sup>54</sup>
- **Enterprises with AI Projects:** Labelbox and V7 Labs cater to enterprises that need scalable, industry-specific data annotation solutions.<sup>55</sup>

### 3.5.1. How Client Relationships are Created

- **Direct Engagement:** Clients often directly approach these companies, especially if they're looking for specialized services. This can involve RFPs (Request for Proposals) where companies like Scale AI might pitch their comprehensive solutions.
- **Partnerships:** Companies like Scale AI have partnerships with major tech players, which can lead to client relationships through referrals or bundled services.<sup>56</sup>
- **Market Presence and Networking:** Through industry events, conferences, and tech expos, these companies establish visibility, leading to networking opportunities and client acquisition.
- **Demonstrations and Trials:** Offering trial periods or demonstrations of their platforms allows potential clients to evaluate the service quality, speed, and fit for their data needs before committing.

### 3.5.2. Client Decision-Making Process

- **Quality and Accuracy:** Clients prioritize the accuracy and quality of annotations, which is crucial for effective model training, as evident from the competitive dynamics.<sup>57</sup>
- **Scalability and Speed:** The ability to handle large volumes of data quickly is a key decision factor, particularly for time-sensitive projects like those in autonomous driving.

<sup>52</sup> Will Henshall (2024), Side Hustle or Scam? What to Know about Data Annotation Work. *Time* (Apr 2). <https://time.com/6962608/data-annotation-legit-tech-jobs-ai/>; Otto Kässi *et al.*, How Many Online Workers are there in the World? A Data-Driven Assessment, *Open Rsch. Eur.* <https://doi.org/10.12688/openreseurope.13639.4>.

<sup>53</sup> Scale, *supra* note 47; Hive.AI, *supra* note 47.

<sup>54</sup> Cloud Factory, *supra* note 48, at *Medical Devices*. <https://www.cloudfactory.com/medical-devices> (last visited Jan 4, 2025); *Medical Research and Pathology*. <https://www.cloudfactory.com/medical-research> (last visited Jan 4, 2025); *Medical Imaging*. <https://www.cloudfactory.com/medical-ai-imaging> (last visited Jan 4, 2025); *Precision Agriculture*. <https://www.cloudfactory.com/precision-agriculture> (last visited Jan 4, 2025).

<sup>55</sup> Labelbox, *supra* note 49; V7, *supra* note 49.

<sup>56</sup> Scale, *Customers*, Scale. <https://scale.com/customers> (last visited Jan 4, 2025); Riley de León (2024), Amazon, Meta Back Scale AI in \$1 billion Funding Deal that Values Firm at \$14 billion. *CNBC* (May 21). <https://www.cnbc.com/2024/05/21/amazon-meta-back-scale-ai-in-1-billion-funding-deal.html>

<sup>57</sup> Nick Tubis (2023), The Secret Sauce of AI Companies: The Importance of Human-Annotated Data in High-Performing Machine Learning Models, *Forbes* (Aug. 3). <https://www.forbes.com/councils/forbescommunicationscouncil/2023/08/03/the-secret-sauce-of-ai-companies-the-importance-of-human-annotated-data-in-high-performing-machine-learning-models/>; Gonçalo Ribeiro (2024), The Critical Role of Data Quality in AI. *Forbes* (Aug 1). <https://www.forbes.com/councils/forbestechcouncil/2024/08/01/the-critical-role-of-data-quality-in-ai/>; Alec Burmania *et al.* (2016), Increasing the Reliability of Crowdsourcing Evaluations Using Online Quality Assessment. *IEEE Transactions on Affective Computing*, 7(374), 374-375. <https://doi.org/10.1109/TAFFC.2015.2493525>.

- **Cost and Efficiency:** Clients consider the cost-effectiveness of the service, balancing between human labor and automation as offered by companies like Hive or V7 Labs for reducing expenses.
- **Compliance and Ethics:** With regulations like GDPR, clients look for providers who ensure data privacy and ethical sourcing practices, as emphasized by CloudFactory and Sama.
- **Customization:** The ability to tailor services to specific industry needs or use cases influences decisions, especially in sectors with unique data requirements like medical imaging.<sup>58</sup>

#### 4. Shortcomings of Centralized Data Optimization for AI models

Recent analyses reveal significant theoretical and practical shortcomings inherent in the centralized frameworks currently employed by leading data-annotation companies such as Scale AI, Appen, Hive, V7 Labs, CloudFactory, and Sama. These challenges involve issues of bias, ethical data sourcing, and data diversity, each of which can undermine the equitability and generalizability of resulting AI models.<sup>59</sup> The difficulty of balancing nuanced human oversight with rapid, automated scaling, coupled with labor costs and the complexities of GDPR compliance, likewise poses formidable obstacles to real-world AI application.<sup>60</sup>

Inherent theoretical concerns, including the propagation of annotator biases and the limited extent of data diversity, threaten the legitimacy and effectiveness of AI outputs across varied demographic and cultural contexts.<sup>61</sup> In addition, practical challenges—ranging from high human-in-the-loop costs and the need to ensure data privacy, to reconciling domain-specific innovations with industry-wide standardization—demonstrate the ongoing tension between speed, scale, and quality of data annotation. These interwoven theoretical and practical shortcomings contextualize the formidable task of refining AI data production processes for the next generation of AI system upgrades. Ultimately, they point to a need for collaborative approaches, improved governance mechanisms, and innovative technological solutions that can manage the delicate balance between efficiency, consistency, and ethical responsibility in AI data provisioning.

Below are summaries in bullet point format of the theoretical and practical shortcomings of AI data production for AI model upgrading via traditional centralized companies:

##### 4.1. Theoretical Shortcomings

- **Bias in Data Annotation:** Despite the emphasis on high-quality data by companies like Scale AI and Appen, there's a theoretical risk that biases inherent in data annotators or automated systems might be propagated into AI models. This can lead to AI that does not perform equitably across different demographic groups or scenarios.<sup>62</sup>
- **Ethical Data Sourcing:** While entities like CloudFactory and Sama focus on ethical and impact sourcing, there's a theoretical concern about whether the data collected or annotated respects the cultural, ethical, and privacy norms of the regions from where it originates. This could impact the global applicability of AI models trained on such data.<sup>63</sup>

<sup>58</sup> Michael Weber *et al.* (2024), *Orchestration Logics for Artificial Intelligence Platforms: From Raw Data to Industry-Specific Applications*. *Special Issue Info. Sys. J.*, 1. <https://doi.org/10.1111/isj.12567>

<sup>59</sup> Kearns and Roth, *supra* note 31; Muldoon *et al.*, *supra* note 41; Caitlin Kuhlman *et al.* (2020), *No Computation Without Representation: Avoiding Data and Algorithm Biases Through Diversity*. *arXiv* (Feb). <https://doi.org/10.48550/arXiv.2002.11836>

<sup>60</sup> Marc Vartabedian (2024), *AI Can Take the Slog Out of Compliance Work, but Executives Not Ready to Fully Trust it*. *Wall Street J.* (Dec 17). <https://www.wsj.com/articles/ai-can-take-the-slog-out-of-compliance-work-but-executives-not-ready-to-fully-trust-it-7cd60a16>; Antonio Zappulla (2024), *Comment: Business Leaders Risk Sleepwalking Towards AI Misuse*. *Reuters* (Nov 19). <https://www.reuters.com/sustainability/society-equity/comment-business-leaders-risk-sleepwalking-towards-ai-misuse-2024-11-19/>; Tijana Žuniæ Mariæ *et al.* (2024), *Comparing the Obligations Under the GDPR and AI Act—Where is the Overlap?*. *Zunic Law* (Oct 22). <https://zuniclaw.com/en/comparing-the-obligations-under-the-gdpr-and-ai-act-where-is-the-overlap/>

<sup>61</sup> See Kearns and Roth, *supra* note 31; Muldoon *et al.*, *supra* note 41; Caitlin Kuhlman *et al.*, *supra* note 60.

<sup>62</sup> Ninareh Mehrabi *et al.* (2021), *A Survey on Bias and Fairness in Machine Learning*. *ACM Computing Surv. Art.*, 54(115), 1. <https://doi.org/10.1145/3457607>

<sup>63</sup> Emily Bender and Timnit Gebru *et al.* (2021), *On the Dangers of Stochastic Parrots: Can Language Models be Too Big?. in Proc. 2021 ACM Conf. on Fairness, Accountability, and Transparency (FAccT'21)*, Mar. 3-10, Virtual Event, Canada, at 610. <https://doi.org/10.1145/3442188.3445922>; Jin Yang *et al.* (2024), *Problematic Tokens: Tokenizer Bias in Large Language Models*. *arXiv* (Preprint arXiv:2406.11214) (Nov 14). <https://doi.org/10.48550/arXiv.2406.11214>

- **Data Diversity:** Although the market emphasizes the need for diverse and high-quality data, there's a theoretical limitation in capturing true diversity, especially in less represented areas or languages, potentially leading to AI models with limited generalizability.<sup>64</sup>

#### 4.2. Practical Shortcomings

- **Scalability and Speed:** The push towards automation for speed and quality, as seen with companies like Hive and V7 Labs, introduces practical challenges. Automation can sometimes result in less nuanced data labeling, potentially missing complex human judgments or context that are critical for certain AI applications.<sup>65</sup>
- **Human-in-the-Loop Costs:** The reliance on human annotators, even with companies striving for ethical labor practices like CloudFactory, involves significant costs, both financial and in terms of time, which can slow down the pace of AI model upgrades.<sup>66</sup>
- **Data Privacy and Compliance:** Ensuring compliance with regulations like GDPR is a practical challenge. As data is centralized in these companies, there's always a risk of data breaches or misuse, which could lead to legal issues or loss of trust in AI technologies.<sup>67</sup>
- **Quality Control:** With the scale at which data annotation occurs, maintaining consistent quality across different annotators or automated systems remains a practical challenge. Small errors in annotation can have significant impacts on model performance, especially in critical applications like autonomous driving, where Scale AI has a strong presence.<sup>68</sup>
- **Innovation vs. Standardization:** While there's innovation in automation and data handling, there's a practical tension between creating standardized data procedures for consistency across industries and innovating for specific use cases, which might limit the adaptability of AI models.<sup>69</sup>

These shortcomings highlight the complexities of scaling AI data production while ensuring that the data quality, diversity, and ethical considerations meet the evolving demands of AI model training and application.

### 5. Decentralized Market for AI Training Data

As the AI data production market continues to evolve, startups grounded in decentralized data production logic possess a distinctive capacity to address the industry's principal imperatives—quality, speed, automation, ethical sourcing, and expansion to new verticals—while simultaneously tackling longstanding issues of compensation and governance. Their reliance on blockchain-based solutions and community-driven protocols naturally aligns with the heightened focus on quality and speed, since transparent smart contracts incentivize contributors to provide timely, high-fidelity data annotations in exchange for tokenized rewards. This structure empowers diverse global workforces to compete on both efficiency and accuracy, ensuring that AI models benefit from broader cultural and linguistic inputs. In a related vein, automation is facilitated via algorithmic moderation of data contributions: smart contracts automatically release payments to contributors once preset quality thresholds are met, thereby decreasing the risk of human error, reducing administrative overhead, and accelerating data-labeling cycles. These attributes address the industry's continual pursuit of innovative annotation methodologies, positioning decentralized platforms at the forefront of AI data production.

In parallel, decentralized startups are uniquely positioned to champion ethical and sustainable practices within the marketplace. By implementing transparent governance mechanisms—often in the form of Decentralized Autonomous Organizations (“DAOs”)—these ventures can equitably compensate data contributors based on demonstrated expertise and reputation, thus mitigating exploitative labor practices

<sup>64</sup> Cf. Mehrabi, *supra* note 63; Kuhkman *et al.*, *supra* note 60; Klie *et al.*, *supra* note 35; Villalobos *et al.*, *supra* note 6; Epoch AI, *supra* note 4.

<sup>65</sup> Bender and Gebru *et al.*, *supra* note 64.

<sup>66</sup> Muldoon *et al.*, *supra* note 41; Klie *et al.*, *supra* note 35.

<sup>67</sup> Giovanni Sartor and Francesca Lagioia (2020), The Impact of the General Data Protection Regulation (GDPR) on Artificial Intelligence. *Panel for Future Sci. & Tech.* (Eur. Parl. Rsch. Serv., Study PE 641.530). <https://data.europa.eu/doi/10.2861/293>.

<sup>68</sup> Klie *et al.*, *supra* note 35.

<sup>69</sup> Weber, *supra* note 59.



sometimes associated with centralized data-annotation services.<sup>70</sup> Token-based reward systems, dynamic pricing models, and voting rights for reputable contributors collectively foster trust and a sense of shared ownership, reinforcing more sustainable “impact sourcing” approaches. This focus on equitable, community-led governance further aligns with the expansion into new verticals where stakeholder trust and compliance with global regulations—such as the European Union’s GDPR<sup>71</sup>—are paramount. Consequently, decentralized frameworks enable these startups to serve specialized data needs in emerging sectors, from medical diagnostics to financial analytics, by leveraging transparent, privacy-preserving methods of contribution and compensation. In so doing, decentralized firms are poised not merely to accommodate the intensifying demands of continuous learning and real-time data streams, but to thrive within them, propelling the AI data production market toward a more equitable and efficient future.

By integrating these multifaceted strategies—transparent token-based compensation, reputation-driven incentives, dynamic pricing, efficient blockchain-based transactions, privacy-by-design principles, and participatory governance—startups can overcome the structural limitations inherent in centralized AI data compensation. The inherent transparency and equity of a decentralized system not only draw more diverse and committed contributors but also uphold ethical data practices in an increasingly regulated environment. As AI applications permeate wider societal and economic spheres, a well-structured decentralized reputation governance model stands to support sustainable growth, stronger data integrity, and broader opportunities for stakeholder engagement.

## 5.1. Reputation Governance

In contemporary AI learning models, many centralized systems suffer from inadequate compensation frameworks and inefficient payment procedures that fail to reward data contributors in proportion to the value of their work. Such shortfalls not only discourage meaningful participation but also create inequities in the broader AI development ecosystem. To mitigate these challenges, the implementation of a decentralized reputation governance model—ideally rooted in blockchain and smart-contract technologies—can offer fairer, more transparent compensation structures. The following strategic measures outline how startups and emerging enterprises can adopt this approach to enhance contributor engagement, ensure data quality, and comply with evolving legal standards.

### 5.1.1. Transparent Compensation Mechanism

**a. Token-Based Rewards:** A central premise of decentralized compensation is the establishment of a native cryptocurrency or token system. Contributors receive tokens commensurate with their participation in data creation or annotation. Over time, these tokens may appreciate in value depending on market demand, thus offering potential financial gains to participants and incentivizing sustained, high-caliber engagement. Moreover, because no single authority holds exclusive control over the tokens, the market regulates their value, ensuring a more democratic approach to compensation.<sup>72</sup>

**b. Smart Contracts:** Smart contracts, hosted on a blockchain network, automate and expedite payments. Once predefined conditions—such as data quality thresholds, volume of contributions, or specific task completions—are met, payments are instantly disbursed. This autonomy obviates the need for intermediaries like banks or centralized platforms, reducing transaction delays, operational costs, and the risk of inconsistent payouts.<sup>73</sup>

### 5.1.2. Reputation-Based Incentives

**a. Reputation Score:** A core component of decentralized governance is the assignment of reputation points. Contributors gain reputation based on the verifiable quality, timeliness, and accuracy of their work. High-

<sup>70</sup> Cf. Youssef El Faqr *et al.* (2020), An Overview of Decentralized Autonomous Organizations on the Blockchain. *Proc. 16<sup>th</sup> Int’l Symp. on Open Collaboration*, 1, Article No. 11. <https://doi.org/10.1145/3412569.3412579>

<sup>71</sup> See *Regulation 2016/679*, *supra* note 34 (emphasizing robust data protection standards that decentralized models may address through transparent data handling and contributor oversight).

<sup>72</sup> See generally Satoshi Nakamoto (2008), Bitcoin: A Peer-to-Peer Electronic Cash System. <https://bitcoin.org/bitcoin.pdf> (discussing foundational concepts of decentralized digital currencies).

<sup>73</sup> Cf. Smart Contracts Alliance (2016), Smart Contracts: 12 Use Cases for Business & Beyond (Chamber of Digit. Com., in Collaboration with Deloitte, Dec.). [https://d3h0qzni6h08fz.cloudfront.net/Smart-Contracts-12-Use-Cases-for-Business-and-Beyond\\_Chamber-of-Digital-Commerce.pdf](https://d3h0qzni6h08fz.cloudfront.net/Smart-Contracts-12-Use-Cases-for-Business-and-Beyond_Chamber-of-Digital-Commerce.pdf)

reputation users can be rewarded with premium compensation rates or access to specialized, higher-value tasks, thereby encouraging ongoing improvement and adherence to community-driven standards.

**b. Voting Rights:** Beyond financial rewards, contributors with superior reputation may obtain governance privileges, such as voting on funding allocations, policy amendments, or strategic platform partnerships.<sup>74</sup> This structure democratizes platform oversight and aligns the most invested, quality-focused participants with key decision-making processes.<sup>75</sup>

### 5.1.3. Decentralized and Fair Pricing

**a. Market-Driven Pricing:** A decentralized marketplace empowers contributors to negotiate compensation directly through bidding or fixed-price listings. This mechanism fosters a more equitable environment in which supply-and-demand forces calibrate the value of data contributions, protecting participants from exploitative rate-setting by singular, centralized entities.

**b. Dynamic Compensation Models:** Algorithms adjusting compensation based on criteria such as task complexity, data type, or submission urgency further refine the market-driven paradigm. Contributors performing especially intricate or specialized work receive higher remuneration, promoting merit-based pay structures and encouraging expertise development across the contributor base.

### 5.1.4. Reducing Transaction Costs and Delays

**a. Blockchain for Efficiency:** Decentralized payment systems eliminate redundant intermediaries. By recording every transaction on a secure, distributed ledger, blockchain mechanisms enable near-instant settlements that minimize overhead costs. This innovation is particularly salient for international workforces, bypassing conventional currency conversion fees and banking delays.

**b. Microtransactions:** High transaction costs in centralized models often preclude micro-payments for smaller tasks. In contrast, blockchain-based ledgers can process tiny financial transfers efficiently, facilitating fair compensation for incremental work and expanding the pool of potential contributors willing to perform micro-tasks that accumulate into large-scale data sets.

### 5.1.5. Privacy and Compliance

**a. Decentralized Data Handling:** To uphold individual contributors' rights, decentralized systems provide personal control over submitted data. Contributors can choose to obscure identifiable markers or limit access through robust encryption protocols. Such measures inherently support alignment with global privacy regulations—such as the European GDPR—because privacy by design is woven into the platform's architecture.<sup>76</sup>

**b. Transparent Compliance:** Blockchain's public, tamper-evident record enables platforms to demonstrate conformance with legal and ethical standards. Stakeholders can track how data is utilized, who has access, and how rewards are distributed, thus fostering confidence among regulators, contributors, and end users.

### 5.1.6. Community Governance

**a. Decentralized Governance (DAO):** Platforms may establish a DAO, wherein pivotal decisions—ranging from adjusting minimum compensation rates to resolving disputes—are collectively ratified by token holders. DAOs have garnered scholarly and industry attention as a means of ensuring inclusivity, transparency, and accountability in digital ecosystems.<sup>77</sup>

<sup>74</sup> See, e.g., Anisha Pandey (2024), Hamster Kombat Enters New Era, Announces the DAO. *Coinspeaker* (Dec 10, 4:13 PM). <https://www.coinspeaker.com/hamster-kombat-enters-new-era-announces-the-dao/>

<sup>75</sup> See Vitalik Buterin (2017), The Meaning of Decentralization. *Medium* (Feb 6). <https://medium.com/@VitalikButerin/the-meaning-of-decentralization-a0c92b76a274>. (articulating principles of decentralized decision-making).

<sup>76</sup> See GDPR, Regulation 2016/679, *supra* note 34 (emphasizing privacy safeguards in data management).

<sup>77</sup> DAOs have garnered scholarly and industry attention as a means of ensuring inclusivity, transparency, and accountability in digital ecosystems. See, e.g., Cristiano Bellavitis *et al.* (2022), The Rise of Decentralized Autonomous Organizations (DAOs): A First Empirical Glimpse. *Venture Capital*, 25, 187; Samer Hassan and Primavera De Filippi (2021), Decentralized Autonomous Organization. *Internet Pol'y Rev.*, 10. <https://doi.org/10.14763/2021.2.1556>; Wulf A. Kaal (2021), A Decentralized Autonomous Organization (DAO) of DAOs (unpublished manuscript) (Mar 9). <http://dx.doi.org/10.2139/ssrn.3799320>; Wulf A. Kaal (2024), DAO Fallacies: Common Myths and Uses for Decentralized Autonomous

**b. Feedback Loop.** Open communication channels—surveys, message boards, or decentralized chat applications—allow contributors to propose adjustments to compensation or governance frameworks. Over time, this cyclical exchange of feedback underpins the platform’s evolution, ensuring that operational and strategic policies remain responsive to the community’s experiences and needs.

## 5.2. Key Competitors

A new generation of startups is emerging to decentralize AI data creation and governance by leveraging blockchain technologies, native tokens, and community-driven protocols. Their core objective is to establish transparent, merit-based ecosystems in which contributors receive equitable compensation for producing or annotating data, as opposed to relying on centralized intermediaries.

### 5.2.1. SingularityNET

**Overview:** SingularityNET is an AI marketplace built on blockchain, designed to permit decentralized creation, sharing, and monetization of AI services. By leveraging a distributed network architecture, the platform aspires to democratize AI-related activities, offering developers and users an open environment in which they can seamlessly transact and collaborate.

**Approach:** At the core of SingularityNET is the native AGI cryptocurrency, which incentivizes active participation and rewards both service providers and end users. Payments and revenue-sharing agreements are encoded in smart contracts, substantially reducing the influence of centralized intermediaries that typically impose or withhold payments. This model not only fosters a level of financial autonomy but also enables transparent auditing of transactions by all stakeholders.<sup>78</sup>

**Reputation System:** A distinguishing feature of SingularityNET is its reputation framework, whereby each AI service provider is evaluated based on user feedback. Service providers exhibiting consistent quality receive higher visibility on the platform, which in turn amplifies their earnings potential. Consequently, the system aligns contributor incentives with service quality, strengthening trust and fostering a meritocratic environment.

### 5.2.2. Fetch.ai

**Overview:** Fetch.ai targets the creation of a decentralized digital economy wherein autonomous software agents can perform tasks such as data annotation. These agents thus become a crucial node in the broader AI supply chain, generating or refining datasets for various machine learning models.

**Approach:** Central to the Fetch.ai ecosystem is the FET token, facilitating low-latency micro transactions among agents and users. By minimizing transaction costs, Fetch.ai promotes real-time compensation for data contribution, thereby attracting a broader user base that may otherwise be deterred by complex or delayed payment processes.<sup>79</sup>

**Reputation System:** Contributors’ software agents accrue reputations based on performance metrics such as annotation accuracy, reliability, and adherence to deadlines. Agents demonstrating superior service quality are subsequently prioritized for future tasks, thereby enjoying enhanced compensation prospects. This feedback loop facilitates a self-regulating system that tends to reward excellence in data provision.

### 5.2.3. Ocean Protocol

**Overview:** Ocean Protocol is designed as a decentralized marketplace for the secure exchange of data assets. It employs blockchain technologies to ensure transparent recordkeeping, thereby mitigating many privacy and trust concerns typically associated with large-scale data-sharing arrangements.

**Approach:** The platform’s OCEAN tokens underpin data transactions, enabling a direct compensation mechanism for data providers. By tokenizing data assets, Ocean Protocol encourages market participants to

---

Organizations. in *Decentralized Autonomous Organizations*, 91 (Sven Van Kerckhoven and Usman W. Chohan (Eds.), Routledge. <http://dx.doi.org/10.2139/ssrn.4067783>; cf. Vitalik Buterin (2024), DAOs, DACs, DAs and More: An Incomplete Terminology Guide (May 6). <https://blog.ethereum.org/2014/05/06/daos-dacs-das-and-more-an-incomplete-terminology-guide>

<sup>78</sup> See SingularityNET Project Overview. <https://singularitynet.io/> (last visited Jan. 5, 2025).

<sup>79</sup> See Fetch.ai Project Overview. <https://fetch.ai/> (last visited Jan. 5, 2025).

share data sets with a quantifiable, transparent metric of value, which in turn catalyzes broader collaboration among data consumers and AI developers.<sup>80</sup>

**Reputation System:** Although Ocean Protocol does not prioritize a reputation mechanism akin to other platforms, its decentralized governance relies substantially on market dynamics. Data providers offering demonstrably higher-quality datasets are naturally rewarded through increased demand, establishing an indirect quality-driven feedback system.

#### 5.2.4. *Numeraire*

**Overview:** Numeraire (NMR) is a platform where data scientists are incentivized to predict stock market outcomes. While primarily focused on financial data, the project's decentralized design principles and compensation structures can be extended to broader AI-related data tasks.

**Approach:** Participants receive NMR tokens for accurate predictions, thereby fostering a merit-driven economy within the platform. This token-based framework underscores the project's commitment to decentralizing decision-making and compensation in data-centric tasks.<sup>81</sup>

**Reputation System:** NMR employs a staking mechanism in which contributors place a portion of their tokens as collateral for their predictions. Effective predictions enhance participants' reputations, improving both their credibility and earning potential. This model aligns token-based economic incentives with the reliability and performance of user contributions.

#### 5.2.5. *DcentAI*

**Overview:** DcentAI focuses on decentralized AI governance by actively involving the broader community in core decision-making processes. A key objective is to ensure that data contributions and model training activities align with stakeholder values rather than solely with corporate directives.

**Approach:** Utilizing a native token for compensating data contributors, DcentAI automates reward distribution through smart contracts, thus diminishing reliance on centralized authorities. The platform aims to engender a more equitable framework for participants, with transparent tracking of how data contributions feed into AI model enhancements.<sup>82</sup>

**Reputation System:** Contributors enhance their standing within DcentAI's governance model via their level of involvement, quality of contributions, and adherence to platform rules. Over time, high-reputation members wield greater influence in shaping platform policies, including determining compensation schemes and adjudicating disputes.

### 5.3. *Shortcomings*

Although decentralized approaches to AI data creation and governance offer significant promise, their transformative potential is tempered by four interrelated challenges that must be addressed to achieve robust market adoption. First, adoption is hindered by the user expertise required to navigate blockchain infrastructure. Potential contributors or data consumers are often deterred by the need to manage cryptocurrency tokens, maintain secure digital wallets, and understand the underlying technology. This skill barrier restricts broader involvement, particularly among users unaccustomed to advanced digital assets. Consequently, developers of decentralized AI platforms must balance intuitive platform design against the need to uphold security, integrity, and trust within the ecosystem.

Second, scalability underscores the operational feasibility of decentralized platforms. To process and annotate vast quantities of data—especially in real-time or near real-time—systems must sustain high throughput while maintaining the decentralized ethos of consensus-based governance. This requirement is especially daunting as existing blockchain infrastructures typically grapple with limits on transaction speed and network capacity. Continuous innovation in consensus algorithms, such as Proof of Stake or layer-two scaling solutions, is thus critical. Without such advances, data pipelines will be prone to bottlenecks that

---

<sup>80</sup> See Ocean Protocol Overview. <https://oceanprotocol.com/> (last visited Jan. 5, 2025).

<sup>81</sup> See Numeraire Project Overview. <https://numer.ai/> (last visited Jan. 5, 2025).

<sup>82</sup> See DcentAI Blog. *Medium*. <https://medium.com/@dcentai> (last visited Jan. 3, 2025).

diminish platform reliability and customer confidence, limiting the appeal of decentralized approaches in large-scale commercial or research deployments.

Third, regulation and compliance requirements, particularly those articulated by the European Union's GDPR, pose substantial legal and ethical complexities. Decentralized ecosystems must adopt privacy-preserving techniques, including encryption protocols and minimized data retention, while also establishing governance frameworks that transparently manage consent and data ownership rights. Moreover, decentralized governance models, such as DAOs, introduce uncertainties regarding liability and legal accountability when personal data is exchanged across international boundaries. A mismatch between platform governance and regulatory mandates can result in legal liabilities, undermining user trust and impeding the broader adoption of decentralized AI platforms.

Finally, competition emerges from two angles: on the one hand, peer decentralized projects vie for market share; on the other, entrenched centralized technology corporations—often with extensive resources—can rapidly pivot to offer competing data annotation and governance services. Well-capitalized incumbents may leverage established customer bases and brand recognition, enabling them to outpace nascent decentralized platforms in user acquisition and marketing. As a result, decentralized startups must demonstrate clear advantages—such as fair compensation, heightened transparency, or superior data protection—to persuade clients that they warrant adoption over mainstream centralized solutions.

Despite these barriers, decentralized AI platforms project an appealing vision of equitable data governance, heightened transparency, and participatory decision-making. Through fair compensation models, contributors from diverse regions and skill levels can be rewarded in proportion to their contributions, mitigating concerns about exploitative or opaque data-sourcing practices. The public, tamper-evident nature of blockchain transactions further facilitates openness in data handling, enabling contributors and users to trace and verify data provenance. Likewise, reputation-based evaluation mechanisms can ensure that high-quality contributions receive the greatest visibility and rewards, thus incentivizing excellence and safeguarding data integrity.

If decentralized AI ventures can reconcile these innovations with large-scale adoption, rigorous security standards, and regulatory compliance, they may successfully challenge existing centralized paradigms. Achieving this balance could redefine the landscape of AI data management, ushering in an ecosystem characterized by heightened user agency, stronger safeguards for personal data, and more robust, equitable models of compensation and collaboration. Ultimately, the success of these platforms will depend on their capacity to communicate tangible benefits, align with legal frameworks, and provide user experiences that rival or surpass those of established centralized competitors.

#### **5.4. Reputation System Shortcomings**

Several of the aforementioned decentralized platforms—namely SingularityNET, Fetch.ai, Ocean Protocol, Numeraire, and DcentAI—have each introduced reputation systems aimed at improving data integrity and building trust in AI services. While these platforms offer incremental innovations in how contributors are evaluated and rewarded, they remain structurally insufficient for a fully decentralized “Mechanical Turk” model of large-scale AI dataset creation. Such a model involves highly contextual, iterative, and ethically sensitive tasks, including domain-specific data annotation and continuous compliance with evolving legal standards.

In contrast, the proposed Weighted Directed Acyclic Graph (“WDAG”) governance approach<sup>83</sup> brings a more dynamic, context-sensitive solution. By enabling every contributor action and piece of data to be linked to on-chain governance elements in a structured graph, the WDAG model supports real-time community oversight, iterative updates, and granular reputation tracking. This layered system is especially suited to AI dataset governance, which demands both thorough traceability—so that datasets are accurately and ethically sourced—and flexibility—so that new ethical guidelines or regulatory changes can be swiftly integrated. By establishing a forum of weighted nodes and directed edges, WDAG provides a multi-dimensional view of contributors' work, ensuring that reputation is not limited to a single numeric score. Instead, it reflects a

---

<sup>83</sup> Wulf A. Kaal (2025), AI Governance Via Web3 Reputation System, *Stan. J. Blockchain L. & Pol'y*, 8(1).



broader, evolving consensus on data quality, ethical compliance, and alignment with sector-specific requirements.

Thus, while existing platforms have laid important groundwork, they lack the depth of feedback loops, adaptive validation, and community-driven layering found in the WDAG system. Consequently, the WDAG approach offers a more robust, transparent, and ethically focused mechanism for governing the complex process of AI dataset creation.

Below I examine each of the proposed systems and their potential shortcomings vis-a-vis the proposed system:

#### 5.4.1. SingularityNET

**Limitation:** SingularityNET's reputation system effectively measures the quality of AI outputs by focusing on service-level performance metrics. However, it places insufficient emphasis on the granular requirements of dataset creation—namely, accuracy, consistency, and contextual relevance. AI training datasets often demand extensive domain-specific input and iterative checks, features that traditional service-oriented reputation scores may not adequately capture.

**Comparison with Proposed System:** By incorporating the WDAG approach, each annotation, dataset, or contribution can be directly tied to a specific governance element, thereby enabling a traceable and context-rich validation process.<sup>84</sup> This depth of monitoring exceeds the simplified rating or review model in SingularityNET's ecosystem. As a result, data contributors would be incentivized to meet precise standards, ensuring that subsequent AI models are built upon high-quality, well-documented data.

#### 5.4.2. Fetch.ai

**Limitation:** Fetch.ai's system centers on metrics such as annotation accuracy and task completion speed. While these metrics can be effective at a basic operational level, they do not incorporate ongoing adjustments for ethical, legal, or community-driven standards—factors that frequently evolve in AI governance. As a consequence, the platform risks reinforcing biases or outdated practices if metrics are insufficiently aligned with changes in societal expectations or regulatory norms.

**Comparison with Proposed System:** The WDAG-based framework integrates Decentralized Community Governance through validation pools and smart contracts, allowing community members to collectively vet data quality and conformance with emerging legal and ethical requirements.<sup>85</sup> This mechanism ensures real-time recalibration of quality criteria, a capability lacking in Fetch.ai's current model.

#### 5.4.3. Ocean Protocol

**Limitation:** Ocean Protocol employs a largely market-driven reputation system in which price signals and demand reflect perceived dataset value. Though effective in encouraging providers to supply high-quality data, this indirect feedback loop does not thoroughly account for the individual expertise or consistency of data annotation work. Market forces can lag behind real-time changes in best practices or ethical standards, leaving critical gaps in data quality and compliance.

**Comparison with Proposed System:** The proposed WDAG system designates Reputation (REP) tokens that directly reward and evaluate individual contributions, thereby highlighting merit at the contributor level.<sup>86</sup> Rather than relying on broad market sentiment, each participant's record is transparently documented and validated against on-chain governance elements, ensuring merit-based progression and more immediate responsiveness to changes in quality standards.

#### 5.4.5. Numeraire

**Limitation:** Numeraire features a *staking* and prediction-based reputation mechanism well-suited to financial modeling. However, this configuration is not readily transferable to the comprehensive vetting required for diverse AI datasets, which may encompass tasks ranging from medical data annotations to text-based content

---

<sup>84</sup> *Id.*

<sup>85</sup> *Id.*

<sup>86</sup> *Id.*

moderation. The narrow focus on predictive accuracy overlooks the broader ethical and contextual concerns that characterize AI governance.

**Comparison with Proposed System:** The proposed WDAG-centric model takes advantage of staking concepts but integrates them into a multi-dimensional framework, ensuring that tasks involving ethical sourcing, regulatory compliance, or subjective expert judgment receive equal scrutiny.<sup>87</sup> In effect, contributors are evaluated on numerous parameters, including but not limited to predictive accuracy, thereby ensuring a more holistic assessment of data quality.

#### 5.4.6. *DcentAI*

**Limitation:** *DcentAI*'s governance approach, emphasizing community-driven decision-making, represents a step forward relative to strictly market-based or performance-only schemes. Nonetheless, its reputation model lacks the *depth* required to track and enforce *dataset quality* across multiple domains and evolving standards. The current framework may struggle to capture interdependencies between various data facets—such as privacy concerns, domain-specific regulations, and real-time ethical updates—within a single reputation score.

**Comparison with Proposed System:** In the WDAG-based framework, each contribution or post undergoes layered scrutiny within a dynamic citation system, thereby establishing precedent and providing context for subsequent validations.<sup>88</sup> The validation pools, combined with smart contracts and WDAG structures, generate a flexible environment where the community can incorporate new ethical guidelines and legal requirements as they arise, updating contributor reputations accordingly.

#### 5.4.7. *Key Advantages of the Proposed WDAG Governance Model*

1. **Granularity and Specificity:** The WDAG architecture links each data contribution to a defined set of governance, ethical, or regulatory metrics. This granular mapping outperforms simpler rating or staking systems, enabling contributors to demonstrate specialized expertise while maintaining transparency and traceability in their work.
2. **Adaptive Real-Time Adjustments:** The proposed model's reliance on validation pools and on-chain governance ensures that reputational adjustments can occur in real time, accommodating rapid shifts in societal, regulatory, or domain-specific requirements. Unlike static feedback systems, the WDAG approach allows for continuous recalibration and alignment with the latest AI governance standards.
3. **Robust Community Consensus:** Decentralized forums and validation pools provide a collective decision-making mechanism that extends beyond individual performance scores. Because AI dataset creation affects stakeholders ranging from software engineers to domain experts and end users, community consensus is critical for ensuring fairness, accuracy, and ethical compliance.
4. **Scalability and Flexibility:** WDAGs support complex relationships and large volumes of contributions without collapsing into inefficiency or opacity. The system can scale alongside the growth of AI data demands, incorporating more contributors, projects, and domains without losing clarity or consistency in its reputation metrics.

## 6. Conclusion

The impending scarcity of high-quality, human-generated textual data for large-scale AI training underscores the critical need for innovative mechanisms that extend beyond conventional data sourcing and centralized annotation models. As illustrated in the foregoing discussion, the exponential demands of LLMs and other AI systems have outpaced existing reserves of publicly accessible text, necessitating more refined approaches that emphasize quality, sustainability, and compliance with ethical and legal standards. While centralized providers and startups have contributed significantly to improving data annotation, curation, and real-time processing, their models often struggle with issues of bias, high operational costs, and limited capacity for dynamic updates.

---

<sup>87</sup> *Id.*

<sup>88</sup> *Id.*

In response, decentralized frameworks—particularly those leveraging blockchain technologies, smart contracts, and distributed governance—present a promising alternative. They address the core challenges of equitable compensation, real-time quality control, robust privacy protections, and scalable data production pipelines. By rewarding contributors through tokenized incentives and reputation-based mechanisms, these decentralized systems foster a participatory and transparent environment that can rapidly adapt to evolving legal regulations (such as the GDPR) and emerging ethical imperatives. Moreover, their inherent flexibility in handling domain-specific tasks aligns well with the diverse, context-sensitive needs of LLMs and other advanced AI architectures.

Nevertheless, as this article has shown, even existing decentralized solutions—exemplified by SingularityNET, Fetch.ai, Ocean Protocol, Numeraire, and DcentAI—do not fully meet the rigorous requirements of a “Mechanical Turk” for AI dataset creation. In contrast, a WDAG governance model offers greater depth of validation, iterative feedback loops, and a multi-dimensional approach to contributor reputation. By aligning each contribution with specific governance elements and enabling collective oversight through validation pools, WDAG systems better ensure data quality, ethical compliance, and stakeholder engagement. This refined model suggests a path toward not only mitigating the limitations of current AI data markets but also laying the foundation for a more equitable, scalable, and future-ready ecosystem.

Ultimately, addressing the challenges outlined herein—data bias, high computational overhead, stagnant or outdated datasets, privacy and security constraints, and inequitable compensation structures—requires concerted, interdisciplinary efforts. Governments, corporations, researchers, and civil society organizations must collaborate to implement robust, decentralized data governance models that can adapt swiftly to technological and societal changes. Such an approach holds the potential to democratize AI innovation, reduce monopolistic data practices, and sustain the ever-growing demands of LLMs and other AI systems. In so doing, the AI research community may find sustainable, ethically grounded, and legally compliant pathways to tap into humanity’s collective knowledge—ensuring that AI can continue to evolve and serve the common good for decades to come.