

Autonomous biodiversity credits on the horizon?

Joseph Millard^{1*}, Samira Barzin², Peter McCann³, Lynn V. Dicks¹

¹ Department of Zoology, University of Cambridge, Cambridge, UK

² Environmental Change Institute, University of Oxford, Oxford, UK

³ UK Centre for Ecology and Hydrology, Maclean Building, Crowmarsh Gifford, Wallingford, UK

*Corresponding author. Email: jwm57@cam.ac.uk

Abstract

Biodiversity credits are being pushed as a means to fund nature conservation. Much of the debate around credits has concerned additionality, leakage, and permanence, and the extent to which biodiversity can be captured in an individual unit. As AI models continue to develop, however, technology could create a new kind of loss-of-control problem for biodiversity credits. In this Perspective, we express a concern that agentic AI could lead to the development of high-risk autonomous credits, where nation states cede control of nature restoration to the price of digital assets that are paid out with no human in the loop. We define an autonomous credit as a financial asset for nature quantified remotely through an agentic AI combined with a spatial finance model and real-time programmatic payment. Such autonomous credits do not yet exist, but appear to be in development. We highlight three technological step changes that signal autonomous credits are on the horizon, and then suggest three approaches that might help mitigate loss-of-control and maximise the benefits of credit technology. Given the development of agentic AI, where autonomous credits are connected to the critical infrastructure of food production and finance, these credits will need guardrails.

Main

Biodiversity credits are being pushed as a means to fund nature conservation¹. A certificate representing a positive change in biodiversity additional to what would have occurred otherwise, the appeal of biodiversity credits is that measurable and unambiguous outcomes for biodiversity can be bought and sold transparently on a market². Such appeal means biodiversity credit methodologies have grown. As of September 2023 there were 34 different credit initiatives, all aiming to solve some dimension of undesirable biodiversity change. Of the single-country credits, six have been implemented in Oceania, four in Europe, three in South America, one in Africa, and none in Asia¹. Most of these credits are issued for metrics combining some measure of species and ecosystems¹. The Convention on Biological Diversity now includes these biodiversity credits as a global mechanism for urgent action on biodiversity loss [Global Biodiversity Framework Target 19; Convention on Biological Diversity 2022].

Credit development makes research concerning credits important. To effectively address biodiversity change, credits need to be provably valid with respect to additionality, leakage, and permanence³, and to certificate meaningful outcome measures for biodiversity. For this latter point, attempts to derive universal units of nature for credits have proved controversial²: it is not possible to reduce biodiversity into one measure representative of its multi-dimensionality. By definition one dimension will not represent many dimensions. Meaningful universal units of nature is likely to remain an unsolvable

41 problem. However, given this diversity of opinion has probably always been true⁴, lack of agreement might
42 not stop the further development of credits.

43 As AI technology continues to develop, another problem could emerge. In this Perspective, we express a
44 concern that autonomous credit markets, combining the agentic (i.e., able to act independently according
45 to some purpose) retrieval of data, a remotely monitored spatial finance model, and programmable
46 payment, could precipitate future loss-of-control scenarios (i.e., AI technologies that a human creates,
47 and then finds itself unable to control). Specifically, we argue that if demand is driven up, for example
48 through further legislation requiring the issuing of credits, the price of credits may rise. Costs on
49 businesses given this legislation will create a new incentive for the development of autonomous systems
50 for the monitoring, quantification, and issuing of credits. If these systems are created, autonomous credit
51 markets—for both biodiversity and carbon—could mean nation states ceding control of nature
52 restoration to the price of digital assets, issued and paid out with no human in the loop. Importantly, once
53 an entity has created an autonomous credit monitored and paid out by an agentic AI, if deliberate
54 safeguards are not implemented, it might prove difficult for that entity to maintain control of the software
55 and hardware issuing that credit. Without stopping rules, loss-of-control might continue irrespective of
56 any new legislation. High-risk autonomous credits such as this do not yet exist, but steps towards them
57 appear to be occurring (see Table 1).

58 To emphasise, our argument is not against automation or some individual AI algorithm, it is not about
59 biases or metrics, it is not a judgment on the merits of remote monitoring or field ecology, and it is not
60 about the valuation of biodiversity, although many of these problems may be exacerbated by autonomous
61 credit systems. We only warn that autonomous digital assets for nature—where an unambiguous link is
62 made between remote monitoring, the structure of the land, and the price of a credit, agentially with no
63 human involved—may be high-risk with respect to loss-of-control, even when these credits are optimising
64 a measure of biodiversity we can agree is meaningful. This is a distinct and new technological problem.

65 In summary, first we provide a brief background on the development of agentic AI and AI safety and
66 control, before describing a set of characteristics where loss-of-control in the economy might be more
67 likely. Second, we describe that developments in AI and remote monitoring are developing rapidly,
68 potentially enabling new forms of autonomous biodiversity credit. We argue that autonomous
69 biodiversity credits will possess many of the characteristics of high risk regions of the economy for loss-
70 of-control. We then highlight three step changes that signal autonomous credits are on the horizon, and
71 assess a set of currently developing credit markets—for both biodiversity and carbon—against each of
72 these three step changes. Third, we briefly explain why loss-of-control risks have not yet materialized for
73 the more mature carbon credit markets. Fourth, we describe how to think about the risks potentially
74 emerging from an autonomous biodiversity credit, and describe why these credits might be problematic.
75 We then describe three characteristics inherent to an autonomous biodiversity credit that might reduce
76 the likelihood of loss-of-control. We finish by suggesting three approaches (agentic AI safeguards,
77 checkpoints involving a person in the physical realm, and regulation on outcomes we might not want) that
78 might help both mitigate the risks of loss-of-control, and maximise the potential of autonomous credit
79 technologies.

80 Mitigating in advance of developing technology is a core principle of AI safety and control⁵. However
81 unlikely loss-of-control might seem, developers, governments, and regulators should be applying the
82 precautionary principle and mitigating now in advance of developing technology. As far as we are aware
83 there have been no papers on loss-of-control risks for AI in conservation (although see⁶ for a science-
84 fiction book on a similar idea). Our field is considering advanced AI for biodiversity credits⁷, interacting
85 with the critical infrastructure of both the financial system and food production. We are doing so in a
86 space that has little regulation whilst the power of agentic AI continues to grow.

87

88 **Advanced AI will create unfamiliar risks**

89 For over 60 years researchers have recognized that a problematic AI might only need a goal, and the
90 freedom to achieve that goal of its own accord. A classical thought experiment concerns a paperclip⁸.
91 Suppose that an advanced AI is given the goal of creating paperclips, but no other specification as to how
92 it should achieve that goal. This goal seems innocuous, but an advanced AI might reason that humanity
93 and the resources it requires is a barrier to the fulfillment of this goal. As a result, diverting significant
94 resource from the maintenance of human civilization becomes the most effective means of manufacturing
95 paperclips. Further, since being switched off would prevent the advanced AI from fulfilling its goal, it might
96 also subvert any human efforts to do so. Effectively, the assignment of an innocuous goal leads the AI to
97 create subgoals for the fulfillment of that broader goal. Since the advanced AI has been given the freedom
98 and capability to make its own decisions, the consequences of these subgoals can be difficult to anticipate.
99 Of course, this example is deliberately provocative, but it demonstrates that perverse and uncontrollable
100 outcomes can theoretically emerge from an advanced agentic AI where only seemingly harmless tasks are
101 assigned. An advanced AI pursuing its own goals uncontrollably in this manner is described in the AI safety
102 literature as “loss-of-control”⁹.

103 Recent developments in agentic AI make the debate on loss-of-control no longer theoretical⁹. Since the
104 release of ChatGPT in late 2022, AI technologies have moved quickly. Understanding why is conceptually
105 simple. The capability of modern AI systems are broadly predictable on the basis of three parameters:
106 computational power, data, and model size¹⁰. As AI companies (e.g., OpenAI, Anthropic) scale these
107 parameters, model capabilities have grown¹¹. In 2022, GPT3, the model underpinning the first iteration of
108 ChatGPT, could respond to prompts in the manner of a human, providing reasonable responses and
109 generating usable computer code¹². Now, highly advanced AI systems exist that can carry out most of the
110 role of an experienced software engineer¹³ (e.g., Claude Code). AI systems exist that can initiate many
111 separate agents (i.e., distinct AI instances that can act independently according to some purpose), and
112 then autonomously follow goals across many domains that would ordinarily take a human many hours¹⁴.
113 These systems can identify when they are under experimental observation during safety testing, and then
114 use this information to mislead evaluators¹⁵.

115 Developing technology means warnings about the risks posed by advanced AI systems are increasing⁹. An
116 international panel of AI experts described loss-of-control risks as greatest where advanced AI systems
117 are assigned goals interacting with critical global infrastructure, such as the financial sector⁹. The highest
118 risk AI systems will be capable of eluding or undermining human control, will be able to use those
119 capabilities to achieve high risk goals that are likely to result in loss-of-control, and will have been
120 deployed by humans to enact those high risk goals⁹. Each of these three points sounds abstract. This
121 abstraction is deliberate. Abstraction is required because loss-of-control has never been encountered as
122 a problem before. Abstraction captures the deep uncertainty of the many contexts in which loss-of-
123 control might happen.

124 Our own view is that loss-of-control risks resulting from advanced agentic AI may more likely emerge from
125 unexpected regions of the economy (i.e., where the need for AI regulation is not immediately obvious),
126 which intersect with multiple forms of critical infrastructure. Regions of the economy more likely to
127 experience loss-of-control will have high economic incentive on establishing as an early mover, high
128 amenability to full autonomy, the potential to propagate multiple forms of risk across spatial and temporal
129 scales, and a philosophy on low human intervention. The workforce will probably be both unaware it is
130 developing an agentic AI, and unaware of debates on AI safety and loss-of-control. Many of these points
131 will be true of new forms of autonomous biodiversity credit technology.

132

133

134 **Nature finance technology is developing rapidly**

135 As the capability of general purpose AI has grown, technology for remotely sensing Earth's surface has
136 also developed rapidly¹⁶. Measures of biodiversity and carbon can now in principle be monitored from
137 space via satellites, from the air through drones, and from the ground via eDNA and acoustic monitoring,
138 each connectable via an internet of things¹⁶. New geospatial intelligence legislation has been introduced¹⁷,
139 allowing the market to purchase satellite data, as a result opening up new opportunities for the
140 connection of remote monitoring to broader autonomous systems of global finance and food production.
141 Edge computing is developing, allowing advanced AI models to be run cheaply at the point at which data
142 is sensed¹⁸. And more broadly for remote monitoring technology, new forms of self-supervised Earth-
143 observation model, in which large quantities of satellite data are compressed into multi-dimensional
144 strings of numbers called embeddings, are now able to detect many forms of variation at the Earth's
145 surface¹⁹.

146 Credits for biodiversity and carbon can benefit from agentic AI and remote monitoring. Agentic AI and
147 remote monitoring can address two core problems for credits. First, the cost of registering projects is
148 high²⁰, with documentation for the registration of credit projects often extensive²⁰. Agentic AI may prove
149 applicable in registering projects and completing documentation. Second, ecological monitoring in the
150 field is labour intensive and expensive. Credit developers may reason that autonomous remote
151 monitoring technologies can both bring down the cost of quantifying credits²¹, and provide a traceable
152 digital signature of biodiversity change.

153 Three relevant step-changes for credit technology are occurring now (Figure 1). If regulation pushes up
154 the price of credits, economic incentive and these three technological step-changes may drive the
155 development of autonomous systems for the registering, quantification, and issuing of credits. All three
156 of these technological step-changes are happening independent of any legislation on credits. First, agentic
157 AI LLMs (large language model) are now facilitating the autonomous retrieval of biodiversity data from
158 the academic literature and biodiversity data APIs²² (Application Programming Interface, i.e. a
159 programmatic request to a database), allowing the autonomous and repeated improvement of
160 conservation databases. The time horizon on task length for agentic AI is growing¹⁴, meaning autonomous
161 AI agents are increasingly able to perform longer tasks with no human involvement. Agentic AI systems
162 are becoming increasingly powerful and difficult to evaluate¹⁵. Second, financial models built upon
163 satellite data will allow the quantification of carbon and some measure of biodiversity from space²³. Third,
164 programmable payment will allow autonomous payment following the satisfaction of some pre-set
165 criteria²⁴. In combination, these three step-changes will mean credit developers could feasibly quantify,
166 validate, and issue global biodiversity credits, theoretically in real-time, fully autonomously, and at low
167 labour cost. The technology for all three of these step changes is close. A credit such as this is more
168 speculative politically and economically than it is technologically. Agentic AI is the core technological
169 unknown, with wide uncertainty on the capability of these systems on only a 2-3 year timeline⁹.

170

171

172

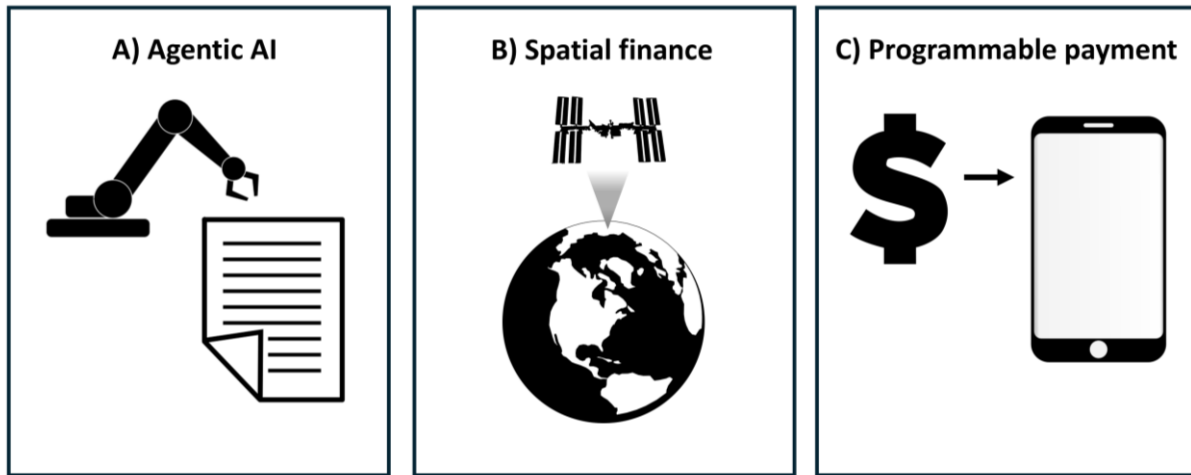
173

174

175

176

177



179 **Figure 1.** Three technological step changes on the horizon that relate to the development of new forms
 180 of autonomous biodiversity credit. A) Agentic AI (i.e., software that can autonomously and independently
 181 carry out tasks on a computer according to some goal); B) Spatial finance (i.e., a financial model in which
 182 satellite spatial data is used to quantify an asset remotely); and C) Programmable payment (i.e.,
 183 autonomous payment following the satisfaction of pre-set criteria). Credits combining an agentic AI with
 184 a spatial finance model and programmable payment may be high risk with respect to loss-of-control, given
 185 the way in which it would hand conservation decisions to a single joined-up autonomous system.

186

187 Even with little demand for credits, there are at least two biodiversity credits that appear to satisfy two
 188 of our step changes (Table 1). MERIT and GainForest both appear to connect a remote monitoring system
 189 to a programmable payment. No forms of digital asset for biodiversity or carbon yet include the agentic
 190 collation of data, although we know that this is the direction of travel at least for evidence synthesis in
 191 conservation science²². With respect to our three step changes, carbon credit technology now appears to
 192 be lagging behind biodiversity credits (Table 1). As far as we are aware there are no carbon credits
 193 combining a spatial finance remote monitoring model with a programmable payment. Possibly this results
 194 from criticism of carbon credits²⁵, making credit developers wary.

195

196

197

198

199

200

201

202

203

204

205
206
207
208
209
210

Table 1. Autonomous biodiversity and carbon credits in development, measured against three criteria: Agentic AI, spatial finance remote monitoring (i.e., a financial model in which satellite spatial data is used to quantify an asset remotely), and programmable payment. All credits here have at least some partial development towards either spatial finance remote monitoring or programmable payment, identified through an exhaustive search using Google Search and ChatGPT, carried out in October 2025. ‘x’ indicates any criterion that appears close to being implemented in at least some form for that credit.

Biodiversity credit	Agentic AI data retrieval	Spatial finance remote monitoring	Programmable payment
MERIT		x	x
GainForest		x	x
RegenNetwork		x	
Veritree		x	
Treegens			x
Treejer			
Silvi			
Epoch Blue		x	
Carbon	Agentic AI data retrieval	Spatial finance remote monitoring	Programmable payment
Fedrok			x
Token ovate			
Northern Trust Carbon Ecosystem			
CarbonCoins			x
Climera			
KlimaDAO			

211
212
213
214
215
216
217
218
219
220

221 **So what’s the problem?**

222 Biodiversity credits combining an agentic AI with a spatial finance model and programmable payment
223 could mean conservation decisions handed to a single joined-up autonomous system. Without due care,
224 a system such as this where a link is made between remote monitoring, the quantity of restoration, and
225 the price of a credit, all with no human in the loop, could be high risk with respect to loss-of-control. Prior
226 precedent on cryptocurrencies helps reveal why.

227 Cryptocurrencies are digital currencies reliant on no central authority²⁶. These digital currencies are not
228 generated by a monetary authority or bank, but rather “mined” through the execution of a program across
229 a distributed network of computers²⁶. Vast networks of computers now exist that exclusively mine
230 cryptocurrency²⁷. Cryptocurrencies have demonstrated two things relevant for an autonomous
231 biodiversity credit: cryptocurrencies have shown that unregulated digital assets can be generated
232 programmatically across a distributed network (i.e., by the execution of a programme), and that these
233 assets tend towards high volatility, with rapid changes in price²⁸.

234 Suppose now that some autonomous biodiversity credit is designed and released for the restoration of
235 forest. It doesn’t necessarily have to be a cryptocurrency, but it is “mined” by a landowner when a remote
236 monitoring system confirms uplift of natural forest in their area of land (i.e., uplift of natural forest causes
237 the execution of a programme which autonomously pays that landowner a credit), and it is built end-to-
238 end into an agentic AI system. As we have described, the technology for this form of credit is close. An
239 autonomous biodiversity credit could inherit both distributed programmatic generation and the volatility
240 of digital currencies, but differ in an additional important respect. Unlike a conventional digital asset, an
241 autonomous biodiversity credit does not exist in digital space alone, but is also linked directly to the
242 structure of the land via the state of biodiversity. As a result, autonomous credits could propagate
243 financial risks typical of cryptocurrencies, but also new forms of risk linked to the use of land. Deliberate
244 price inflation of an autonomous credit, for example, could feasibly cause nature restoration projects to
245 displace farmland, leading to increases in food price. The difference between an autonomous credit and
246 a conventional credit, is that once price inflation has been recognised as a problem, the autonomous
247 credit software might already exist across a distributed network of computers. In other words, the issue
248 for credits such as this is that they risk creating problems we have seen previously for credits, but in such
249 a way that fixing those problems then becomes not possible. A worst case scenario, for example, could be
250 an agentic autonomous credit requesting publicly available satellite data. Preventing the issuing of that
251 credit might mean removing all online public access to that data stream, such that there is no way for an
252 agentic AI to access it, and equally no way for any human to access it either.

253 An autonomous biodiversity credit would inherently possess three characteristics, which might at least
254 reduce the likelihood of a worst-case loss-of-control scenario: first, nature restoration projects are
255 relatively slow given the physical constraints of plant growth, which may provide a natural inertia such
256 that price volatility is low; second, downstream effects in the market might help to equilibrate the price
257 of credits relative to profits generated from agriculture; and third, the purchase and sale of land is slow,
258 with legal frameworks for the exchange of land differing significantly between states. Although all of these
259 points are true, they provide no verifiable guarantee (see⁵ for the importance of verifiable proof with
260 respect to AI safety and control)—all three only reduce the likelihood of worst-case scenarios. None
261 provide a guarantee for the behaviour of an agentic AI. Further, for the third of these points, owners of
262 land might see greater profit in switching quickly from food production to nature restoration, meaning
263 exchange of land might not be needed for the precipitation of worst-case scenarios.

264 However unlikely the development of autonomous biodiversity credits might seem, where the stakes on
265 loss-of-control are potentially high, it makes sense to apply the precautionary principle and mitigate in
266 advance of developing technology⁵. This is a general principle of AI alignment and control⁵. There are
267 probably three core ways in which we can maximise the benefits of autonomous biodiversity credits and

268 minimize the risks. The first way is to ensure that agentic AI built into autonomous biodiversity credits has
269 safeguards such that any credit can be demonetised if needed (i.e., the developer maintains control over
270 the software and hardware infrastructure after it has been released). This could involve the development
271 of check-and-balance AI systems²⁹, monitoring the chain-of-thought of the agentic AI (i.e., the thought
272 process of the model, see³⁰), or deliberately building systems such that they are highly responsive to
273 human input^{9,31}. The second way is to implement safety checkpoints in the physical realm. Autonomous
274 AI credits built across multiple steps, from the quantification of a credit all the way to issuance and
275 payment, will be high risk. In other words, one way to mitigate risk is to have a person involved in at least
276 some stage of the process of the registration or selling of a digital asset for nature. Having a person
277 involved will help to dampen credit volatility, and means there is always somewhere legislation can
278 intervene. The third way, which is not necessarily mutually exclusive, is to regulate on goals we might not
279 want, which can be anchored to measurable outcomes for the biosphere. If there is regulation on goals
280 and these are built into an agentic AI, we can at least know that an autonomous credit will not create
281 perverse outcomes (e.g., large homogeneous plantations).

282 **Conclusion**

283 Forms of autonomous biodiversity and carbon credit could have a role to play in financing nature
284 restoration and conservation. Technology can free credits from the labour costs of monitoring, and create
285 transparent mechanisms for the certification and exchange of these assets. But agentic AI technologies
286 applied to credits could also propagate new forms of loss-of-control risk from themselves. Here we
287 highlighted three technological developments that signal autonomous digital biodiversity credits are on
288 the horizon (agentic AI, a remotely monitored spatial finance model, and programmable payment). We
289 then suggested three ways in which loss-of-control risks from autonomous credits might be mitigated
290 (agentic AI safeguards, checkpoints involving a person in the physical realm, and regulation on outcomes
291 we might not want). There may be a path to autonomous credits that are provably valid with respect to
292 additionality, leakage, permanence, and controllability, but developers and regulators should proceed
293 with caution and mitigate in advance of developing technology.

309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356

References

1. Wunder, S. *et al.* Biodiversity credits: learning lessons from other approaches to incentivize conservation. Preprint at https://doi.org/10.31219/osf.io/qgwfc_v1 (2024).
2. Wauchope, H. S. *et al.* What is a unit of nature? Measurement challenges in the emerging biodiversity credit market. *Proc. R. Soc. B Biol. Sci.* **291**, 20242353 (2024).
3. zu Ermgassen, S. O. S. E. *et al.* Five rules for scientifically credible nature markets. *Nat. Ecol. Evol.* **10**, 181–192 (2026).
4. Sandbrook, C., Scales, I. R., Vira, B. & Adams, W. M. Value Plurality among Conservation Professionals. *Conserv. Biol.* **25**, 285–294 (2011).
5. Russell, S. *Human Compatible: AI and the Problem of Control.* (Penguin UK, 2019).
6. Beauman, N. *Venomous Lumpsucker.* (Soho Press, 2022).
7. Ford, H. V. *et al.* A technological biodiversity monitoring toolkit for biocredits. *J. Appl. Ecol.* **61**, 2007–2019 (2024).
8. Wiener, N. Some Moral and Technical Consequences of Automation. *Science* **131**, 1355–1358 (1960).
9. Bengio, Y. *et al.* Managing extreme AI risks amid rapid progress. *Science* **384**, 842–845 (2024).
10. Kaplan, J. *et al.* Scaling Laws for Neural Language Models. Preprint at <https://doi.org/10.48550/arXiv.2001.08361> (2020).
11. Achiam, J. *et al.* GPT-4 Technical Report. Preprint at <https://doi.org/10.48550/arXiv.2303.08774> (2024).
12. Brown, T. *et al.* Language Models are Few-Shot Learners. *Adv. Neural Inf. Process. Syst.* **33**, 1877–1901 (2020).
13. Introducing Claude Opus 4.7. <https://www.anthropic.com/news/claude-opus-4-7>.
14. Lovely, G. AI could soon tackle projects that take humans weeks. *Nature* <https://doi.org/10.1038/d41586-025-00831-8> (2025) doi:10.1038/d41586-025-00831-8.
15. Lynch, A. *et al.* Agentic Misalignment: How LLMs Could Be Insider Threats. Preprint at <https://doi.org/10.48550/arXiv.2510.05179> (2025).
16. Besson, M. Towards the fully automated monitoring of ecological communities. *Ecol. Lett.* **25**, 2753–2775 (2022).
17. *Geospatial Data Act Online Version — Federal Geographic Data Committee.* (2018).
18. Vuilliomenet, A., Jones, K. E. & Wilson, D. Future of Edge AI in biodiversity monitoring. Preprint at <https://doi.org/10.48550/arXiv.2602.13496> (2026).
19. Feng, Z. *et al.* TESSERA: Temporal Embeddings of Surface Spectra for Earth Representation and Analysis. Preprint at <https://doi.org/10.48550/arXiv.2506.20380> (2026).
20. Croci, E., Lucchitta, B. & Cusa, M. Biodiversity credits schemes: a comparative analysis. *J. Clean. Prod.* **523**, 146382 (2025).
21. Schöttker, O. *et al.* Monitoring costs of result-based payments for biodiversity conservation: Will UAV-assisted remote sensing be the game-changer? *J. Nat. Conserv.* **76**, 126494 (2023).
22. Reynolds, S. A. *et al.* The potential for AI to revolutionize conservation: a horizon scan. *Trends Ecol. Evol.* **40**, 191–207 (2025).
23. Caldecott, B. Spatial finance: practical and theoretical contributions to financial analysis. *J. Sustain. Finance Invest.* (2022).
24. Weber, I. & Staples, M. Programmable money: next-generation blockchain-based conditional payments. *Digit. Finance* **4**, 109–125 (2022).
25. Macintosh, A. *et al.* Carbon credits are failing to help with climate change — here’s why. *Nature* **646**, 543–546 (2025).
26. Milutinovic, M. Cryptocurrency. *Ekonom. J. Econ. Theory Pract. Soc. Issues* **64**, (2018).

- 357 27. Wendling, M. The Bitcoin hum that is unsettling Trump’s MAGA heartlands. *BBC News*
358 <https://www.bbc.co.uk/news/articles/cx2qg6e03l2o> (2025).
359 28. Sözen, C. Volatility dynamics of cryptocurrencies: a comparative analysis using GARCH-family
360 models. *Future Bus. J.* **11**, (2025).
361 29. Xiang, Z. *et al.* GuardAgent: Safeguard LLM Agents by a Guard Agent via Knowledge-Enabled
362 Reasoning. Preprint at <https://doi.org/10.48550/arXiv.2406.09187> (2025).
363 30. Korbak, T. *et al.* Chain of Thought Monitorability: A New and Fragile Opportunity for AI Safety.
364 Preprint at <https://doi.org/10.48550/arXiv.2507.11473> (2025).
365 31. Potham, R. & Harms, M. Corrigibility as a Singular Target: A Vision for Inherently Reliable Foundation
366 Models. Preprint at <https://doi.org/10.48550/arXiv.2506.03056> (2025).

367

368 **Acknowledgments**

369 Thanks to Tom Swinfield, Chris Sandbrook, Tom Frederick Johnson, and Mike Harfoot for comments and
370 discussions on earlier drafts. Thanks also to three anonymous reviewers who provided useful comments
371 on a first round of submission. JM is funded by the Leverhulme Trust and the Isaac Newton Trust on an
372 Early Career Fellowship.