

# Curse of rarity for autonomous vehicles

Henry X. Liu & Shuo Feng

 Check for updates

The curse of rarity—the rarity of safety-critical events in high-dimensional variable spaces—presents significant challenges in ensuring the safety of autonomous vehicles using deep learning. Looking at it from distinct perspectives, we identify three potential approaches for addressing the issue.

The concept of autonomous vehicles (AVs) has been around for about a century. Over the past two decades, AVs have attracted extensive attention from academic institutions, government agencies, professional organizations, and industries. By 2015, multiple companies had announced that they would mass-produce AVs before 2020 (Ref. 1). However, the reality has not lived up to expectations, and there are currently no commercially available SAE Level 4 (Ref. 2) AVs. One of the main reasons is the significant gap in safety performance of AVs<sup>1</sup>. This gap poses a major challenge as AVs struggle to effectively handle a multitude of rare safety-critical events, despite the accumulation of millions of testing miles on public roads. The occurrence of these events, characterized by a probability distribution resembling a long tail that is far from the head or central part of the distribution, is commonly referred to as the long-tail challenge for AV safety<sup>3,4</sup>. The catchphrase “long-tail challenge” for AV safety, however, is frequently used in a handwaving manner without a formal definition in the literature. This lack of understanding impedes progress in resolving the issue.

In this Comment, we uncover that the shape of the probability distribution for safety-critical occurrences, whether it exhibits a long tail or not, is not essential to the issue at hand. Instead, the primary challenge in defining the problem stems from the rareness of safety-critical situations in highly complex driving environments, which encompass various factors such as different weather conditions, diverse road infrastructures, and behavioral distinctions among road users. The safety-critical circumstances may arise due to a variety of reasons, such as misidentification of an unknown object or inaccurate prediction of nearby pedestrian’s movement, all of which have a low probability of occurrence. We term this challenge as the curse of rarity (CoR) and mathematically define CoR for a generic deep learning problem, which is commonly used for perception, behavior modelling, prediction and decision making in AVs. CoR emerges from the combination of rare occurrence of safety-critical situations and the vast number of variables involved, resulting in a compounding effect. Such an effect hinders the ability of deep learning models to perform safely in real-time<sup>5,6</sup>.

In the following, we elaborate the CoR in different AV tasks including perception, prediction, planning, as well as validation and verification. Based on these analyses, we discuss potential solutions towards addressing the CoR. We hope that this Comment can provide a better understanding of the safety challenges faced by the AV community, and a rigorous formulation of CoR can help accelerate the

development and deployment of AVs as well as other safety-critical autonomous systems.

## What is the curse of rarity?

The basic concept of CoR is that the occurrence probability for the events of interest in high-dimensional space is so rare that most available data contain very little information of the rare events. Therefore, it is hard for a deep-learning model to learn, since valuable information of rare events could be buried under a large amount of normal data. It becomes particularly challenging to improve safety performance because better safety performance also means a lower frequency of safety-critical events, which makes it more difficult for the deep-learning model to learn. An illustration of CoR can be found in Box 1.

## What challenges does it bring for autonomous vehicles?

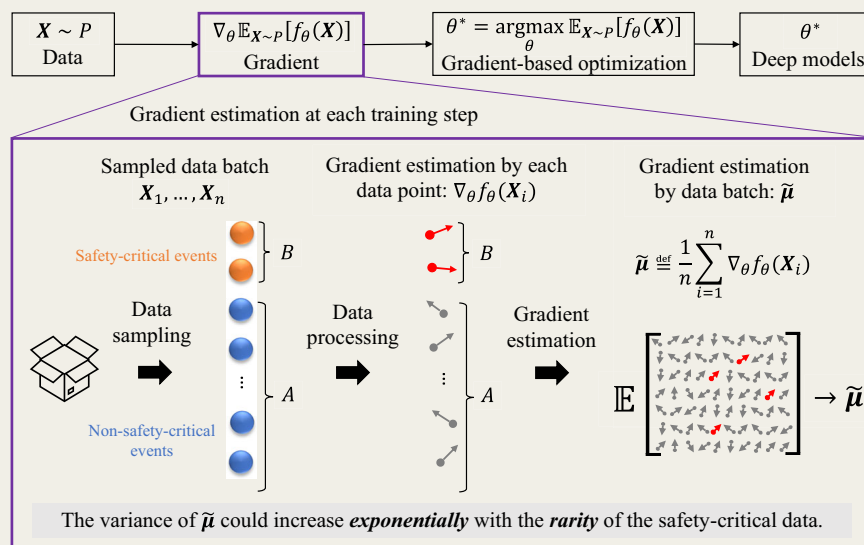
In this section, we elaborate on the CoR in various aspects of AVs including perception, prediction, planning, validation and verification.

**Perception.** Deep learning methods have been extensively utilized in perception tasks to acquire information and extract pertinent knowledge from the surrounding environment. The problem of imbalanced data has been studied in perception tasks, where a small portion of object classes have a large number of samples, while the remaining classes have only a few samples<sup>3,7</sup>. However, this issue becomes particularly challenging for safety-critical perception tasks of AVs, as the imbalance ratios are much more severe, often exceeding  $10^6$  (Ref. 8). Existing approaches such as class rebalancing, information augmentation, and module improvement are inadequate in addressing this problem, as they can only handle a limited imbalance ratio, usually smaller than  $10^3$  (Ref. 7). This significant difference in magnitude fundamentally transforms the problem from an imbalanced data issue to the CoR problem. Moreover, the cumulative effects of a series of perception errors could be dangerous, even if each individual error appears insignificant. For example, an object misclassification in a single frame might be less of an issue, while multiple object misclassifications in a sequence of frames may lead to safety-critical outcomes. Since the occurrence probability of such a sequence is much lower than that of any individual error, the issue of CoR becomes even more severe.

**Behavior prediction and simulation.** AV’s high safety performance requirements necessitate precise behavior modeling and accurate prediction of surrounding road users. Even a minor error in predicting the behaviors of surrounding road users can be deemed unacceptable in safety-critical situations. For example, in a jaywalking scenario, precise prediction of pedestrian trajectories is crucial for AVs to avoid collisions. A small prediction error could result in either a false alarm or a missed alarm, leading to overly cautious driving decisions or overly confident decision that cause an accident. The same holds true for driving behavior simulation. Inaccuracies in simulations can lead to underestimation or overestimation of AV’s safety performance,

## BOX 1

### Curse of rarity for deep learning models



The key to deep learning is to obtain the optimal parameters  $\theta^*$  of neural networks by optimizing the expectations of the objective function  $f_{\theta}(\bullet)$  over the data  $X$  with an underlying distribution  $P$ . To solve this optimization problem, the most commonly used approach is based on the gradient descent (see Chapter 8 in ref. 19). Existing approaches estimate the gradient with Monte Carlo estimation<sup>20</sup> using a batch of data at each training step. However, the estimation variance could increase exponentially with the rarity of safety-critical events (see Theorem 1 in Methods), resulting in the curse of rarity. This analysis

is applicable to various deep-learning approaches. In the case of deep reinforcement learning, the data consists of the state-action pairs, and the objective function is based on the reward function (see our previous work<sup>5</sup> for details); whereas in deep supervised learning, the data comprises the labelled data, and the objective function is based on the loss function. More details can be found in the Methods section.

thereby misleading the development process<sup>9</sup>. To achieve the required level of safety, behavior prediction models must effectively handle rare events in high-dimensional driving environments, which are prone to the CoR.

**Decision making.** Deep learning techniques, such as deep imitation learning and deep reinforcement learning, have been applied in the decision-making process of AVs. However, when it comes to safety-critical scenarios, deep learning models suffers from the CoR due to the scarcity of real-world data. This scarcity may lead to severe variance in the estimation of policy gradients, thereby impeding the effectiveness of deep learning<sup>5</sup>. Another approach aiming to ensure the safety of decision-making involves using formal methods based on a set of assumptions. Typical assumptions include the availability of a system model, which may be characterized by bounded unknown dynamics and noise<sup>10</sup>. Due to the CoR, it is difficult to verify these assumptions to account for all rare safety-critical events in high-dimensional driving environments.

**Verification and validation.** Verification and validation of safety performance play a crucial role in assessing the readiness of AVs for widespread deployment<sup>5</sup>. Prevailing approaches usually test AVs in the naturalistic driving environment through a combination of software simulation, closed test track, and on-road testing. Due to the CoR, however, hundreds of millions of miles would be required to evaluate the safety performance of AVs, which is impractical and inefficient<sup>8</sup>. To accelerate the process, various approaches have been developed, such as scenario-based approaches, which focus on testing AVs in purposely generated scenarios. Unfortunately, the complexity of generating spatiotemporally intricate safety-critical scenarios poses significant challenge due to the CoR. For example, it has been found that the importance-sampling-based approaches could suffer from a severe inefficiency owing to the dramatic variance for generating complex safety-critical scenarios<sup>6</sup>. As a result, many existing approaches are limited to handling short scenario segments with limited dynamic objects, failing to capture the full complexity and variability of real-world safety-critical events<sup>6</sup>.

## What are the potential solutions?

Based on the analyses and discussions above, we identify three potential approaches for solving the CoR problem, each addressing it from a distinct perspective. It is important to note that these approaches are not mutually exclusive, and combining these approaches holds immense potential in resolving the CoR issue and expediting the widespread deployment of AVs.

**Approach #1: Effective training with more rare event data.** The first approach focuses on data and aims to continually improve the handling of rare events by making better use of additional data. One potential method is to utilize exclusively the data associated with rare events, which could significantly reduce the estimation variance, as stated in Theorem 1 in Methods section. However, defining and identifying rare events are challenging, as they depend on problem-specific objective functions and suffer from the spatiotemporal complexity of safety-critical autonomous systems. More importantly, theoretical foundations that can guide the utilization of rare event data remain lacking. For AV safety validation tasks, tackling the CoR issue has been attempted by developing the dense deep reinforcement learning (D2RL) approach in our prior work<sup>5</sup>. Theoretical and experimental results show that D2RL can dramatically reduce the variance of the policy gradient estimation, a significant step towards addressing the CoR. Another crucial concern is how to gather or generate more rare event data. Tesla proposed the concept of shadow mode testing<sup>11</sup>, where rare events of interest are identified by comparing human driving behavior with autonomous driving behavior, but no details are given in the literature. Other than collecting data from naturalistic driving environment, various data augmentation methods have been developed to generate safety-critical scenarios<sup>12</sup>.

**Approach #2: Improving capabilities of generalization and reasoning.** The second approach centers around improving the generalization and reasoning capabilities of machine learning models to overcome the data insufficiency. Intuitively, as humans can learn to drive with limited experience (typically less than one hundred hours of training), future AI agents for AVs may be able to overcome the CoR without relying on extensive task-specific data. This requires an AI agent to possess both bottom-up reasoning (sensing data-driven) and top-down reasoning (cognition expectation-driven) capabilities<sup>13</sup>, bridging the information gap not found in the data. These requirements are in line with the development of artificial general intelligence (AGI). Recently, foundational models such as large language models (LLMs) and vision-language models (VLMs) have exhibited remarkable generalization and reasoning abilities in terms of natural language processing and visual comprehension and reasoning by employing techniques such as fully supervised fine-tuning, in-context learning, and chain of thought. By leveraging the extensive data available, LLMs and VLMs present a promising solution for enabling top-down reasoning to address the CoR issue<sup>14</sup>, although issues like hallucinations still need further investigations<sup>15</sup>.

**Approach #3: Reducing the occurrence of safety-critical events.** The third approach aims to mitigate the consequences of CoR on AV systems by reducing the occurrences of safety-critical events. Potentially, one can combine traditional model-based approaches with deep learning approaches, taking advantages of the strengths of both<sup>16</sup>. For example, formal methods have been developed to prevent unsafe behaviors of AVs based on abstract models, potentially leading to

defensive driving strategies. However, as discussed in ref. 10,17, multiple challenges need to be addressed to fully harness the potential of formal methods. Another approach is to enhance situational awareness by utilizing infrastructure-based sensors or cooperative awareness, aiding AVs in overcoming the limitations of their own onboard sensors. Nevertheless, effectively utilizing this additional information to achieve improved performance remains a challenging task, especially in safety-critical scenarios. Many existing approaches may even result in inferior perception and decision-making outcomes in such scenarios, due to the increased complexity and latency associated with gathering and integrating this extra information<sup>18</sup>.

## Methods

Let us consider a general deep learning problem that can be formulated as an optimization problem:

$$\max_{\theta} \mathbb{E}_P[f_{\theta}(\mathbf{X})], \quad (1)$$

where  $\theta \in \mathbb{R}^d$  denotes the parameters of a neural network,  $d$  is the dimension of the parameters,  $\mathbf{X} \in \Omega$  denotes the training data with an underlying distribution  $P$ , and  $f_{\theta}(\mathbf{X})$  denotes the objective function given the neural network  $\theta$  and training data  $\mathbf{X}$ . To optimize the objective function, the key is to estimate the gradient of the neural network parameters at each training iteration (see Chapter 8 in ref. 19) as

$$\boldsymbol{\mu} \stackrel{\text{def}}{=} \nabla_{\theta} \mathbb{E}_P[f_{\theta}(\mathbf{X})] \approx \tilde{\boldsymbol{\mu}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} f_{\theta}(\mathbf{X}_i), \mathbf{X}_i \sim P \quad (2)$$

where  $n$  denotes the number of training data samples at each iteration,  $\nabla_{\theta}$  denotes the gradient of parameters, and the approximation is obtained using the Monte Carlo method<sup>20</sup>. Let  $\tilde{\boldsymbol{\mu}}^{(k)}$  denote the  $k$ th component of  $\tilde{\boldsymbol{\mu}}$ , where  $k = 1, \dots, d$ . According to the Monte Carlo method,  $\tilde{\boldsymbol{\mu}}$  is an unbiased estimation of  $\boldsymbol{\mu}$ , that is,  $\mathbb{E}_P(\tilde{\boldsymbol{\mu}}) = \boldsymbol{\mu}$ . The variance of  $\tilde{\boldsymbol{\mu}}^{(k)}$  can be denoted as  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)})$ . To simplify the notations, we denote  $\mathbf{Y} \stackrel{\text{def}}{=} \nabla_{\theta} f_{\theta}(\mathbf{X})$  as a random vector where  $\mathbf{Y} = [Y_1, \dots, Y_d] \in \mathbb{R}^d$ , so  $\tilde{\boldsymbol{\mu}}$  in Eq. (2) can be represented as

$$\tilde{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{Y}_i. \quad (3)$$

Now let's focus on a special set of deep learning problems where only a very small portion of training data (safety-critical data) can contribute effectively to the gradient estimation, while a vast majority of training data (non-safety-critical data) contributes little. To be more specific, we can define normal events  $A \subset \Omega$  and critical but rare events  $B \subset \Omega$ , where  $A \cap B = \emptyset$  and  $A \cup B = \Omega$ . We can also define the corresponding indicator function  $\mathbb{I}_A(\mathbf{X})$ , where  $\mathbb{I}_A(\mathbf{X}) = 1$  if  $\mathbf{X}$  belongs to the set  $A$  and otherwise  $\mathbb{I}_A(\mathbf{X}) = 0$ .  $\mathbb{I}_B(\mathbf{X})$  can be defined similarly. Then, we can obtain a new estimator of the gradient that only utilizes the samples associated with the events  $B$  as

$$\hat{\boldsymbol{\mu}} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n (\mathbf{Y}_i \cdot \mathbb{I}_B(\mathbf{X}_i)), \quad (4)$$

where  $\hat{\boldsymbol{\mu}}^{(k)}$  denotes the  $k$ th component of  $\hat{\boldsymbol{\mu}}$  and the variance of  $\hat{\boldsymbol{\mu}}^{(k)}$  can be denoted as  $\sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$ .

Then we have the following theorem, and the proof can be found at the end of Methods.

## Theorem 1:

If the set  $A$  satisfies the following condition:

$$\mathbb{E}_P[\mathbf{Y} \cdot \mathbb{I}_A(\mathbf{X})] = \mathbf{0}, \quad (5)$$

we have the following properties:

- (1)  $\mathbb{E}_P(\tilde{\boldsymbol{\mu}}) = \mathbb{E}_P(\hat{\boldsymbol{\mu}}) = \boldsymbol{\mu}$ ;
- (2)  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}) \geq \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$ ; and
- (3)  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}) \geq 10^r \cdot \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$ , with the assumption

$$\mathbb{E}_P(Y_k^2 \cdot \mathbb{I}_B(\mathbf{X})) = \mathbb{E}_P(Y_k^2) \cdot \mathbb{E}_P(\mathbb{I}_B(\mathbf{X})), k=1, \dots, d, \quad (6)$$

where  $r \stackrel{\text{def}}{=} -\log_{10}[\mathbb{E}_P(\mathbb{I}_B(\mathbf{X}))]$  is defined as the rarity of the events  $B$  in all samples with the sampling distribution  $P$ .

**Remark 1.** The condition in Eq. (5) indicates that the non-safety-critical data ( $\mathbb{I}_A(\mathbf{X})=1$ ) contributes little to the gradient. Taking the AV safety testing task as an example (see ref. 5 for details), the key is to learn a deep model to control background vehicles to conduct adversarial maneuvers. In this case, the non-safety-critical data that could be identified by safety metrics usually contains no information for learning such adversarial maneuvers, so the condition could be satisfied. We note that the condition is primarily for the theoretical analysis to be clean and is not strictly required in practice. For example, if  $\mathbb{E}_P[\mathbf{Y} \cdot \mathbb{I}_A(\mathbf{X})]$  is a near-zero value and dramatically smaller than  $\mathbb{E}_P[\mathbf{Y} \cdot \mathbb{I}_B(\mathbf{X})]$ , we can still find that the variance of  $\tilde{\boldsymbol{\mu}}$  increases dramatically with the rarity of safety-critical events.

**Remark 2.** Defining and identifying the events  $A$  and  $B$  are non-trivial and dependent on specific deep learning tasks. An important aspect of these definitions is the approximate fulfillment of the condition stated in Eq. (5), as explained in Remark 1. To illustrate, in the context of AV safety testing, we have chosen safety-critical states as events  $B$  and non-safety-critical states as events  $A$  (see ref. 5 for details). The definitions will vary across different AV tasks, warranting further exploration.

**Remark 3.** The assumption in Eq. (6) can be satisfied if all  $Y_k^2, k=1, \dots, d$  are independent of the events  $B$ . For deep learning approaches, the gradient  $\mathbf{Y}_{\text{def}} = \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}}(\mathbf{X})$  is mainly determined by the parameters  $\boldsymbol{\theta}$  of neural networks. As the parameters are usually randomly initiated,  $\mathbf{Y}$  could have an uncertainty that is approximately independent of the events  $A$  and  $B$ , particularly at the beginning of the learning process. Therefore, the assumption could be approximately satisfied particularly at the beginning of the learning process, so the CoR hinders the effectiveness of learning from the very beginning. Again, we note that the assumption is primarily for the theoretical analysis to be clean and is not strictly required in practice.

**Remark 4.** The third property suggests that the variance  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)})$  will grow exponentially with the rarity of the events  $B$ , provided that  $\sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$  does not decrease exponentially with the rarity. As the estimator  $\hat{\boldsymbol{\mu}}$  is primarily focused on estimating the gradient using safety-critical events, its variance  $\sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$  will not be affected significantly by the rarity.

*Proof of Theorem 1.*

- (1) Proof of  $\mathbb{E}_P(\tilde{\boldsymbol{\mu}}) = \mathbb{E}_P(\hat{\boldsymbol{\mu}}) = \boldsymbol{\mu}$ :

$$\begin{aligned} \mathbb{E}_P(\tilde{\boldsymbol{\mu}}) &= \mathbb{E}_P\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{Y}(\mathbf{X}_i) \cdot \mathbb{I}_B(\mathbf{X}_i))\right) = \mathbb{E}_P(\mathbf{Y}(\mathbf{X}_i) \cdot \mathbb{I}_B(\mathbf{X}_i)) = \boldsymbol{\mu} \\ &= \mathbb{E}_P(\hat{\boldsymbol{\mu}}). \end{aligned}$$

*End of proof.*

- (2) Proof of  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}) \geq \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$ :

$$\begin{aligned} \sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}) &= \text{Var}_P[\mathbf{Y}_k \cdot \mathbb{I}_B(\mathbf{X})] = \mathbb{E}_P[Y_k^2 \cdot \mathbb{I}_B(\mathbf{X})] - \mathbb{E}_P^2[\mathbf{Y}_k \cdot \mathbb{I}_B(\mathbf{X})] \\ &= \mathbb{E}_P[Y_k^2 \cdot \mathbb{I}_B(\mathbf{X})] - \mathbb{E}_P^2(Y_k) \leq \mathbb{E}_P[Y_k^2 \cdot \mathbb{I}_B(\mathbf{X})] + \mathbb{E}_P[Y_k^2 \cdot \mathbb{I}_A(\mathbf{X})] \\ &\quad - \mathbb{E}_P^2(Y_k) = \mathbb{E}_P[Y_k^2] - \mathbb{E}_P^2(Y_k) = \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)}). \end{aligned}$$

*End of proof.*

- (3) Proof of  $\sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}) \geq 10^r \cdot \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)})$ :

$$\begin{aligned} \sigma_P^2(\hat{\boldsymbol{\mu}}^{(k)}) &= \mathbb{E}_P[Y_k^2 \cdot \mathbb{I}_B(\mathbf{X})] - \mathbb{E}_P^2(Y_k) = \mathbb{E}_P(Y_k^2) \cdot \mathbb{E}_P(\mathbb{I}_B(\mathbf{X})) - \mathbb{E}_P^2(Y_k) \\ &\leq \mathbb{E}_P(Y_k^2) \cdot \mathbb{E}_P(\mathbb{I}_B(\mathbf{X})) - \mathbb{E}_P^2(Y_k) \cdot \mathbb{E}_P(\mathbb{I}_B(\mathbf{X})) = \mathbb{E}_P(\mathbb{I}_B(\mathbf{X})) \cdot \\ &\quad [\mathbb{E}_P(Y_k^2) - \mathbb{E}_P^2(Y_k)] = 10^{-r} \cdot \sigma_P^2(\tilde{\boldsymbol{\mu}}^{(k)}). \end{aligned}$$

*End of proof.*

**Henry X. Liu** <sup>1,2</sup> & **Shuo Feng** <sup>3</sup>

<sup>1</sup>Department of Civil and Environmental Engineering, University of Michigan, Ann Arbor, MI, USA. <sup>2</sup>Mcity, University of Michigan, Ann Arbor, MI, USA. <sup>3</sup>Department of Automation, Beijing National Research Center for Information Science and Technology, Tsinghua University, Beijing, China. ✉e-mail: [henryliu@umich.edu](mailto:henryliu@umich.edu); [fshuo@tsinghua.edu.cn](mailto:fshuo@tsinghua.edu.cn)

Received: 21 December 2023; Accepted: 27 May 2024;

Published online: 05 June 2024

## References

1. Safe driving cars. *Nat. Mach. Intell.* **4**, 95–96 (2022).
2. Society of Automotive Engineers. Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles. J3026-202104 Available at [https://www.sae.org/standards/content/j3016\\_202104/](https://www.sae.org/standards/content/j3016_202104/) (2021).
3. Zhang, Y., Kang, B., Hooi, B., Yan, S. & Feng, J. Deep long-tailed learning: a survey. *IEEE Trans. Pattern Anal. Mach. Intell.* **45**, 10795–10816 (2023).
4. Wang, J. et al. Parallel vision for long-tail regularization: initial results from IVFC autonomous driving testing. *IEEE Trans. Intell. Veh.* **7**, 286–299 (2022).
5. Feng, S. et al. Dense reinforcement learning for safety validation of autonomous vehicles. *Nature* **615**, 620–627 (2023).
6. Feng, S., Yan, X., Sun, H., Feng, Y. & Liu, H. X. Intelligent driving intelligence test for autonomous vehicles with naturalistic and adversarial environment. *Nat. Commun.* **12**, 748 (2021).
7. Johnson, J. M. & Khoshgoftaar, T. M. Survey on deep learning with class imbalance. *J. Big Data* **6**, 1–54 (2019).
8. Kalra, N. & Paddock, S. M. Driving to safety: how many miles of driving would it take to demonstrate autonomous vehicle reliability? *Transp. Res. Part A Policy Pract.* **94**, 182–193 (2016).
9. Yan, X. et al. Learning naturalistic driving environment with statistical realism. *Nat. Commun.* **14**, 2037 (2023).
10. Brunke, L. et al. Safe learning in robotics: from learning-based control to safe reinforcement learning. *Annu. Rev. Control Robot. Auton. Syst.* **5**, 411–444 (2021).
11. Karpathy, A. Tesla Inc. System and method for obtaining training data. U.S. Patent Application 17/250,825 Available at <https://patents.google.com/patent/US20210271259A1/en> (2021).
12. Wang, J. et al. AdvSim: Generating safety-critical scenarios for self-driving vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 9909–9918 [https://openaccess.thecvf.com/content/CVPR2021/html/Wang\\_AdvSim\\_Generating\\_Safety-Critical\\_Scenarios\\_for\\_Self-Driving\\_Vehicles\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Wang_AdvSim_Generating_Safety-Critical_Scenarios_for_Self-Driving_Vehicles_CVPR_2021_paper.html) (2021).
13. Cummings, M. L. Rethinking the maturity of artificial intelligence in safety-critical settings. *AI Mag.* **42**, 6–15 (2021).
14. Tian, X. et al. DriveVLM: The convergence of autonomous driving and large vision-language models. Preprint at: <https://arxiv.org/abs/2402.12289> (2024).
15. Kandpal, N., Deng, H., Roberts, A., Wallace, E. & Raffel, C. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, 15696–15707 <https://proceedings.mlr.press/v202/kandpal23a.html> (2023).
16. Krasowski, H. et al. Provably safe reinforcement learning: Conceptual analysis, survey, and benchmarking. *Trans. Mach. Learn. Res.*, 1–38 available at <https://openreview.net/pdf?id=mcN0ezbnzO> (2023).

17. Seshia, S. A., Sadigh, D. & Sastry, S. S. Toward verified artificial intelligence. *Commun. ACM* **65**, 46–55 (2022).
18. Bai, Z. et al. Infrastructure-based object detection and tracking for cooperative driving automation: a survey. In *2022 IEEE Intelligent Vehicles Symposium*, 1366–1373 <https://doi.org/10.1109/IV51971.2022.9827461> (IEEE, Aachen, Germany 2022).
19. Goodfellow, I., Bengio, Y. & Courville, A. Deep learning. available at <https://mitpress.mit.edu/9780262035613/deep-learning/> (MIT Press, 2016).
20. Owen, A. B. Monte Carlo Theory, Methods and Examples. Preprint at <https://artowen.su.domains/mc/> (2013).

## Acknowledgements

This research was partially funded by the US Department of Transportation (USDOT) Region 5 University Transportation Center: Center for Connected and Automated Transportation (CCAT) of the University of Michigan (#69A3551747105) and the National Science Foundation (CMMI #2223517). Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the official policy or position of the US government.

## Author contributions

H.X.L. and S.F. equally contributed to the preparation of the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Correspondence** and requests for materials should be addressed to Henry X. Liu or Shuo Feng.

**Peer review information** *Nature Communications* thanks Matthias Althoff, Fredrik Warg and Colin Paterson for their contribution to the peer review of this work.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024