

AI in EU Law: Training, Hallucinations, Memorization, and the Core Regulatory Architecture

Prof. Dr. Philipp Hacker, LL.M. (Yale)

March 6, 2026

Abstract:

This chapter examines the regulatory architecture for artificial intelligence in European Union law across the three stages of the machine learning pipeline: training, model, and output. First, in the training context, it analyzes the tensions between GDPR requirements and large-scale AI training, the copyright framework under the CDSM Directive’s text and data mining exceptions, and the AI Act’s data governance obligations. I also offer a detailed assessment of the GDPR Omnibus proposals. Second, at the model level, the chapter discusses the divergent rulings in *Getty Images v Stability AI* and *GEMA v OpenAI* on the question whether trained models constitute reproductions of copyrighted works. It also weighs in on the structurally parallel question whether models themselves constitute personal data. At the output level, it addresses hallucinations under the accuracy principle and proposes a strict liability regime modeled on pharmaceutical law.

The chapter identifies a structural symmetry between data protection and copyright: in both domains, legally protected content is encoded in model parameters in a manner that renders surgical excision technically near-infeasible. It argues, however, that the appropriate remedies diverge – property rules for data protection with specific exceptions, liability rules with output-based remuneration for copyright.

The chapter concludes with policy proposals for each level of the pipeline and advocates a regulatory approach that combines technology-neutral baseline protection with targeted technology-specific interventions.

Contents

A. Introduction.....	3
I. Defining AI: Technical and Legal Perspectives	3
II. The Machine Learning Pipeline as a Regulatory Framework.....	3
B. The Training Level: Data Governance Challenges	4
I. Data Protection: GDPR Constraints on AI Training	4
1. The Legitimate Interest Balancing Test.....	4
2. Sensitive Data Under Article 9 GDPR.....	5
3. The GDPR Omnibus Proposal	6
II. Copyright: The TDM Exception and Its Limits	7

1. The CDSM Directive Framework.....	8
2. The LG München I GEMA v OpenAI Decision.....	9
III. AI Act Data Governance Requirements	9
C. The Model Level: Legal Status of Trained Models	11
I. Copyright: The Model as Reproduction	11
II. Data Protection: The Model as Personal Data	12
D. The Output Level: Hallucinations and Memorization	13
I. Hallucinations and Data Accuracy	13
II. Memorization and Copyright Infringement	16
E. Evaluation: Technology Neutrality Versus Technology-Specific Regulation	16
I. Technology-Neutral Regulation	16
II. Technology-Specific Regulation.....	17
F. Policy Proposals.....	17
I. Training Level Reforms.....	17
II. Model Level Reforms.....	18
III. Output Level Reforms.....	19
G. Structural and Overarching Considerations	19
H. Conclusion	20
Bibliography	22

A. Introduction

Artificial intelligence (AI) has fundamentally transformed the technological landscape, with large language models (LLMs) and other generative AI systems at the forefront of this transformation.¹ This chapter examines the current and emerging regulatory architecture for AI in European Union law, with particular emphasis on how existing legal frameworks – data protection, copyright, product liability, and sector-specific AI regulation – map onto the distinct stages of the machine learning pipeline.²

I. Defining AI: Technical and Legal Perspectives

From a computer science perspective, AI refers broadly to computational systems that perform tasks typically associated with human cognition.³ Machine learning, a subset of AI, involves algorithms that improve their performance through experience – specifically, through exposure to data. The AI Act adopts a broader approach in Article 3(1), defining an AI system as ‘a machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.’⁴

This definition reflects an attempt to capture the essence of modern AI systems while remaining sufficiently flexible to accommodate future technological developments. The emphasis on autonomy, adaptiveness, and inferential capabilities distinguishes AI systems from conventional software, though the boundaries remain contested.⁵

II. The Machine Learning Pipeline as a Regulatory Framework

The machine learning pipeline provides a useful analytical framework for understanding how different legal regimes apply to AI systems.⁶ This pipeline can be conceptualized in three primary stages. First, at the training level, the process involves data collection, preprocessing, and the statistical optimization of model parameters. Second, at the model level, the trained

¹ See, e.g., Hacker and others (eds), *The Oxford Handbook on Generative AI and Law*, Oxford University Press 2025.

² On computational AI definitions, see Goodfellow/Bengio/Courville, *Deep Learning*, MIT Press 2016, 1 et seqq.; on the AI Act definition, see Hacker/Engel/Hammer/Mittelstadt, 'Introduction' in Hacker and others (eds), *The Oxford Handbook on Generative AI and Law*, Oxford University Press 2025.

³ See, e.g., Stuart J. Russell and Peter Norvig, *Artificial Intelligence: A Modern Approach* (4th Global edn, Pearson Education 2022), 19-22.

⁴ Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act), OJ L, 2024/1689.

⁵ See, e.g., European Law Institute, *Commission Guidelines on the Application of the Definition of an AI System and the Prohibited AI Practices Established in the AI Act*. Response of the European Law Institute, 2024; Philipp Hacker, 'Comments on the final trilogue version of the AI Act' (2024) Available at SSRN 4757603; European Commission. (2025b). *Guidelines on the definition of an artificial intelligence system* (C(2025) 5053 final).

⁶ See, e.g., Hacker, 'A Legal Framework for AI Training Data – From First Principles to the Artificial Intelligence Act' (2021) 13 *Law, Innovation and Technology* 257.

model exists as a computational artifact with specific capabilities and limitations. Third, at the output level, the model generates predictions, content, or decisions in response to user inputs.⁷

Each stage raises distinct legal questions and is subject to different – sometimes overlapping, sometimes conflicting – regulatory frameworks. The challenge for EU law is to ensure coherent and effective governance across these stages while balancing innovation interests with fundamental rights protection.⁸

B. The Training Level: Data Governance Challenges

The training of AI models involves the collection and processing of vast quantities of data. This stage is governed primarily by data protection law, copyright law, and the AI Act's data governance requirements for high-risk systems and general-purpose AI (GPAI) models.

I. Data Protection: GDPR Constraints on AI Training

The General Data Protection Regulation (GDPR) applies when AI training involves personal data – information that relates to an identified or identifiable natural person. Given that LLMs are typically trained on web-scraped data that includes substantial amounts of personal information, GDPR compliance is a central concern for AI developers operating in or targeting the EU market.⁹

1. The Legitimate Interest Balancing Test

In the absence of viable consent mechanisms for large-scale data processing, developers must typically rely on the legitimate interest ground under Article 6(1)(f) GDPR. This provision requires a three-part assessment: the controller must pursue a legitimate interest; the processing must be necessary to achieve that interest; and the controller's interest must not be overridden by the data subjects' fundamental rights and freedoms.¹⁰

Socially beneficial applications of AI weigh in favor of developers in this balancing exercise. Privacy-enhancing measures such as pseudonymization, differential privacy, and robust transparency mechanisms also strengthen the developer's position.¹¹ However, the reasonable expectations of data subjects – a criterion emphasized in Recital 47 GDPR – will rarely support

⁷ Kotsiantis/Kanellopoulos/Pintelas, 'Data Preprocessing for Supervised Learning' (2006) 1(2) International Journal of Computer Science 111 (117 et seq.).

⁸ See Novelli/Casolari/Hacker/Spedicato/Floridi, 'Generative AI in EU Law: Liability, Privacy, Intellectual Property, and Cybersecurity' (2024) 55 Computer Law & Security Review 106066.

⁹ See, e.g., Sartor, Giovanni, and Francesca Lagioia. The impact of the General Data Protection Regulation (GDPR) on artificial intelligence. Study for the EP, 2020; Hacker, AI & data protection: applications to the AV sector, in: Artificial intelligence in the audiovisual sector, Report for the Council of Europe (October 22, 2024), <https://www.obs.coe.int/en/web/observatoire/-/new-report-the-challenges-of-ai-for-the-audiovisual-sector-and-the-role-european-legislation-is-playing>; Hacker (n 6) 263 et seqq.; see also Mayer-Schönberger/Cukier, Big Data: A Revolution That Will Transform How We Live, Work, and Think, Houghton Mifflin Harcourt 2013.

¹⁰ CJEU, Case C-252/21, ECLI:EU:C:2023:537 – Meta v Bundeskartellamt.

¹¹ Novelli and others (n 8); see also Hacker (n 6) 287 et seqq.

AI training, as most individuals do not typically anticipate that their online data will be used for this purpose.¹² In a recent and much-contested summary judgment, however, the Higher Regional Court (OLG) of Cologne held, however, that the balancing test actually falls in favor of AI companies, in the concrete case *Meta*.¹³ The CJEU has not yet weighed in on this matter.

2. Sensitive Data Under Article 9 GDPR

The main battleground in data protection law is not finding a legal basis under Article 6, but an exception under Article 9 GDPR concerning special-category data. The Court of Justice's Grand Chamber decision in *Meta v Bundeskartellamt* significantly expanded the scope of Article 9 GDPR. The Court held that data need not directly reference protected categories such as health, religion, or ethnic origin to qualify as sensitive. Rather, it suffices that 'data processing allows information falling within one of those categories to be revealed.'¹⁴

This broad interpretation has profound implications for AI training. Machine learning techniques increasingly enable the inference of sensitive attributes from facially innocuous data points, and large data sets from machine learning will typically comprise sensitive data.¹⁵ Under the *Meta v Bundeskartellamt* approach, the mere potential for such inference may bring processing within Article 9's ambit, regardless of whether the controller intends to make such inferences.¹⁶

Unlike Article 6 GDPR, Article 9 does not contain a general balancing test. For example, for product development, developers must therefore identify a specific exception in Article 9(2). Outside of explicit consent, such exceptions are often unavailable. The research exemption in Article 9(2)(j), as implemented in national law, may provide some relief, though it is typically restricted to non-commercial activities.¹⁷ Germany, for instance, has introduced a balancing test concerning sensitive data processing for research under § 27 German Data Protection Act, BDSG, which is significantly stricter than the balancing test under Article 6, though, and does not apply to commercial AI product development.

In its infamous *Meta* decision, the OLG Köln brushed these concerns aside and held, in a spectacularly misguided judgment, but Article 9 GDPR does not stand in the way of AI training because, under its rather arbitrary reading, the AI Act implicitly constrains Article 9 for the

¹² On the reasonable expectations criterion, see Rec. 47 GDPR; for discussion, see Hacker (n 3) 291 et seq.

¹³ OLG Köln, Case 15 UKI 2/25, NJW 2025, 3156, para. 60 et seq.; for a discussion, see, e.g., Schwartmann, Rolf, Zulässigkeit von Datenverarbeitungen in der KI-Entwicklung – Der Stand zwischen OLG Köln und Digitalem Omnibus, EuDIR 2026, 3.

¹⁴ CJEU, Case C-252/21, ECLI:EU:C:2023:537 – *Meta v Bundeskartellamt*, para 73.

¹⁵ OLG Köln, Case 15 UKI 2/25, NJW 2025, 3156, para. 93; see also Christiane Wendehorst, Keep Calm and Read the Proposals: A Legal Lens on the Digital Omnibus, EuCML 2025, 277, 279.

¹⁶ See the discussion in Gianclaudio Malgieri and Giovanni Comandé, 'Sensitive-by-distance: quasi-health data in the algorithmic era' (2017) 26 Information & Communications Technology Law 229; on algorithmic inference generally, see Sandra Wachter and Brent Mittelstadt, 'A Right to Reasonable Inferences' (2019) 2 Columbia Business Law Review 494.

¹⁷ See Benedikt Buchner and Marie-Theres Tinnefeld, § 27 BDSG para 8, in: Jürgen Kühling and Benedikt Buchner (eds), DS-GVO BDSG (2nd edn, CH Beck 2018); János Mészáros and Chih-hsing Ho, 'Big Data and Scientific Research' (2021) 21 International Journal of Law and Information Technology 403 (412 et seq.).

sake of achieving AI leadership in the EU.¹⁸ None of this, in my view, withstands serious doctrinal scrutiny.¹⁹

3. The GDPR Omnibus Proposal

The European Commission's GDPR Omnibus proposal addresses some of these challenges through two significant amendments.²⁰ First, the proposed Article 88c would clarify that AI development and operation can constitute a legitimate interest under Article 6(1)(f) GDPR, subject to appropriate safeguards.²¹ Second, the envisaged Article 9(2)(k) would create a specific exemption for processing sensitive data in the context of AI system and model development, accompanied by mandatory technical and organizational measures under a new Article 9(5).²²

These proposals merit careful consideration. While they would provide needed legal certainty for AI development, concerns arise regarding the breadth of the exemptions and the adequacy of the proposed safeguards.

As a threshold matter, one must ask whether an exemption for sensitive data in AI contexts is warranted at all. The answer, I submit, is affirmative. Without such an exemption, AI training in the EU is de facto nearly impossible for most commercially relevant use cases. To be precise, without Article 9(2)(k), lawful processing of sensitive data in AI training would be confined to smaller models for which individual consent can be obtained for all training data, and to research purposes – for instance, under § 27 BDSG in Germany. This means that the proposed exemption is, in practice, primarily relevant for product development rather than for research, although it would formally apply to both. A clear statutory rule would also help to prevent ill-fated decisions such as the OLG Köln judgment on Meta's AI training, which illustrates the legal uncertainty that currently plagues the field. That said, the exemption should not be unlimited in scope; it ought to be restricted to socially beneficial, high-impact sectors to avoid an overly permissive regime. One might further ask whether similar exemptions are needed for other technologies. While this cannot be ruled out in principle, no pressing need exists at the present time.

¹⁸ OLG Köln, Case 15 UK1 2/25, NJW 2025, 3156, para. 100 et seq. (“In the AI Act, European legislators, aware of the need to train generative AI using large amounts of data [...] and of the long-standing use of web scraping to obtain AI training data, which always carries the risk of unintentional and non-targeted processing of data within the meaning of Article 9 of the GDPR, has not included any opening provisions that would allow AI to be trained with mass data. In the Senate's view, this can only be understood to mean that the legislator did not assume that such training was impossible or fundamentally unlawful. Otherwise, a provision enabling such training would have been mandatory in a law that serves to achieve a “pioneering role” [in AI]” para. 104 [translation by DeepL and author]).

¹⁹ See also, in a similar vein, Schwartmann, Rolf, Zulässigkeit von Datenverarbeitungen in der KI-Entwicklung – Der Stand zwischen OLG Köln und Digitalem Omnibus, EuDIR 2026, 3, 8-9.

²⁰ More generally on the GDPR proposals, see Gebehenne/Siebler/Hennemann, Der Digital Omnibus – Grundstruktur, Einordnung und Rahmenbedingungen der Vorschläge der EU-Kommission für eine Vereinfachung im Datenrecht, EuDIR 2026, 10; Christiane Wendehorst, Keep Calm and Read the Proposals: A Legal Lens on the Digital Omnibus, EuCML 2025, 277.

²¹ See proposed Art. 88c GDPR Omnibus.

²² See proposed Art. 9(2)(k) and Art. 9(5) GDPR Omnibus.

Turning to the formulation of the exemption, Article 9(2)(k) as proposed requires refinement in several respects. The Commission's draft uses the phrase "in the context of" AI system and model development. This formulation is too broad; it should be replaced with a necessity requirement – "necessary for" – to impose a meaningful constraint. Moreover, the exemption should be limited to what is strictly necessary to *develop and modify AI models*, not *AI systems*.²³ The concept of an AI system is far broader and can encompass virtually any software application that goes beyond basic data processing; the system component of AI software does not differ from traditional software with respect to personal data processing and can therefore be addressed through existing legal bases. For the same reason, the exemption should not extend to the *operation* of AI systems, as this could cover an essentially unlimited range of processing activities that must remain subject to existing exemptions and safeguards. The bottleneck squarely resides in training models; this is what the exemption should exclusively focus on. Otherwise, processing sensitive data is allowed, for example in analyzing sensitive patient health data, simply because an AI system, and not standard software, is used, which does not make any sense from the perspective of equal treatment or doctrinal coherence.

The safeguards in the proposed Article 9(5) also require recalibration. The requirement to implement measures to "avoid the collection" of sensitive data, combined with obligations to remove or protect such data if identified, generally represents a balanced approach – one that mirrors what Article 10(5) of the AI Act already requires for high-risk systems. However, the specific formulation needs tightening. The measures for detection and deletion of sensitive data should be held to a state-of-the-art standard to remain meaningful as technology evolves. Furthermore, the current drafting, which calls for "proportionate efforts" for removal, makes little practical sense if no threshold for detection is established; proportionality in removal presupposes that effective detection mechanisms are in place. Finally, the requirement to adopt effective strategies against memorization – while conceptually sound – poses significant technical challenges. Here, too, the standard should be state-of-the-art measures, which acknowledges that perfect prevention of memorization may not be feasible while still demanding robust technical safeguards.

Overall, the GDPR will remain at the forefront of legal challenges to AI training processes, even if the Digital Omnibus rules are enacted in the summer of 2026.

II. Copyright: The TDM Exception and Its Limits

Data protection, copyright is the other key battleground concerning the intersection of AI training of the law. Training data frequently includes copyrighted works – texts, images, code,

²³ The proposals seek to limit the scope in Recital 33 of the Digital Omnibus: "This derogation should not apply where the processing of special categories of personal data is necessary for the purpose of the processing." However, any such fundamental limitation must be formulated in the article itself, not in a recital. See also Gebehenne/Siebler/Hennemann (n 20), 13.

and other creative content. The use of such material for AI training engages reproduction and potentially adaptation rights under copyright law.²⁴

1. The CDSM Directive Framework

The Directive on Copyright in the Digital Single Market (CDSM Directive) establishes two text and data mining (TDM) exceptions. Article 3 provides a mandatory exception for research organizations and cultural heritage institutions for scientific research purposes. Article 4 establishes a broader exception for any user with lawful access to content, subject to an opt-out mechanism by which rightholders may reserve their rights.²⁵

The scope of these exceptions encompasses the reproduction and extraction of works for TDM purposes. Preprocessing activities such as normalization, which are essential for machine learning, are covered, too.²⁶ Commercial AI developers are generally confined to Article 4, as Article 3's scope is limited to entities that do not operate for profit or that reinvest all profits in research.²⁷

Notably, the TDM exceptions do not provide for rightholder remuneration, reflecting the legislative assumption that the use of works for automated analysis does not significantly impair exploitation interests.²⁸ Against this background, an active debate has arisen concerning the applicability and contours of the TDM exception to AI training, particularly in the context of generative AI.²⁹

²⁴ Directive (EU) 2019/790 of the European Parliament and of the Council of 17 April 2019 on copyright and related rights in the Digital Single Market (CDSM Directive), OJ L 130, 92.

²⁵ Christophe Geiger, Giancarlo Frosio, and Oleksandr Bulayenko, 'The Exception for Text and Data Mining (TDM) in the Proposed Directive on Copyright in the Digital Single Market – Legal Aspects' (European Parliament, PE 604.941, 2018) 6; Benjamin Raue, 'Rechtssicherheit für datengestützte Forschung' (2019) ZUM 684 (685); Eva Inés Obergefell, 'Big Data und Urheberrecht' in Ahrens and others (eds), Festschrift für Wolfgang Büscher, 2018, 223 (226); Gerald Spindler, 'Text und Data Mining' (2016) GRUR 1112 (1113).

²⁶ Eleonora Rosati, 'The Exception for Text and Data Mining (TDM) in the Proposed Directive on Copyright in the Digital Single Market – Technical Aspects' (European Parliament, PE 604.942, 2018), 3 et seqq.

²⁷ See Art. 2(1) CDSM Directive; Rec. 12 CDSM Directive; Raue (n 25) 690.

²⁸ Spindler, 'Die neue Urheberrechts-Richtlinie der EU' (2019) Computer und Recht 277 (281).

²⁹ See, e.g., Kalpana Tyagi, 'Copyright, text & data mining and the innovation dimension of generative AI' (2024) 19 Journal of Intellectual Property Law & Practice 557; Thomas Margoni and Martin Kretschmer, 'A deeper look into the EU text and data mining exceptions: harmonisation, data ownership, and the future of technology' (2022) 71 GRUR International 685; Matthias Leistner and Lucie Antoine, 'TDM and AI Training in the European Union—From 'LAION' to Possible Ways Ahead?' (2025) 74 GRUR International 1027; João Pedro Quintais, 'Generative AI, copyright and the AI Act' (2025) 56 Computer Law & Security Review 106107; Eleonora Rosati, 'The exception for text and data mining (TDM) in the proposed Directive on Copyright in the Digital Single Market: technical aspects' (2018) European Parliament; from the German literature see, e.g., Tim W Dornis, 'Generatives KI-Training und der TDM-Trugschluss' (2024) GRUR 1676; Franz Hofmann, 'Zehn Thesen zu Künstlicher Intelligenz (KI) und Urheberrecht' (2024) WRP 11; Susanne Grimm and Laura Marie Münster, 'Training von KI in der EU: Fragen zum Nutzungsvorbehalt beim Text und Data Mining sowie zur Beweislast' (2025) GRURPrax 443; Paulina Jo Pesch and Rainer Böhme, 'Artocalypse now? – Generative KI und die Vervielfältigung von Trainingsbildern' (2023) GRUR 997; Malte Baumann, 'Generative KI und Urheberrecht – Urheber und Anwender im Spannungsfeld' (2023) NJW 3673; Katharina de la Durantaye, '"Garbage in, garbage out" – Die Regulierung generativer künstlicher Intelligenz durch Urheberrecht' (2023) ZUM 645; Haimo Schack, 'Auslesen von Webseiten zu KI-Trainingszwecken als Urheberrechtsverletzung de

2. The LG München I GEMA v OpenAI Decision

The Landgericht München I's December 2025 decision in *GEMA v OpenAI* represents a landmark ruling on the copyright implications of AI training.³⁰ The case concerned OpenAI's reproduction of song lyrics through its ChatGPT models, with GEMA – the German collecting society for musical works – asserting infringement claims.³¹

The court distinguished between three phases of the training process. Phase 1 involves the preparation of training data, including the copying of works into training datasets. Phase 2 covers the actual training process, during which works may become 'memorized' in the model's parameters. The court held that both phases involve acts of reproduction within the meaning of § 16 UrhG (implementing Article 2 InfoSoc Directive).³² The third phase concerns the output, where reproductions may again occur.

With respect to the training phase, the Munich court held that the preparation of training data material in Phase 1 (extracting and converting training material into machine-readable format) falls under the TDM exception as it serves the purpose of preparing for text and data mining. Concerning the mere collection and annotation of training data, a similar conclusion had been reached by the district and his original courts of Hamburg in a different case (LAION).³³ I expect this interpretation, which is correct, to be upheld by the CJEU in its upcoming landmark decision on AI and copyright (case *Like Company v. Google*).³⁴ In this sense, at least for the collection and preparation of training data, the European TDM exception offers a clearer compliance path than the fair use exception under US copyright law.³⁵

III. AI Act Data Governance Requirements

While both the GDPR and the TDM exception date from the second half of the 2010 years, the most recent regulatory push concerning AI is epitomized by the AI Act, passed in the summer of 2024. The AI Act establishes specific data governance obligations for high-risk AI systems under Article 10 and for GPAI models under Articles 53 and 55. For high-risk systems, training

lege lata et ferenda' (2024) NJW 113; Roman Konertz and Raoul Schönhof, 'Vervielfältigungen und die Text- und Data-Mining-Schranke beim Training von (generativer) Künstlicher Intelligenz' (2024) WRP 289; Matthias Leistner, 'TDM und KI-Training in der Europäischen Union, Erste Fingerzeige des LG Hamburg im "LAION"-Urteil' (2024) GRUR 1665; Matthias Leistner, 'Memorisierungen und urheberrechtliche Vervielfältigungen in KI-Modellen: Das LG München I betritt Neuland' (2026) GRUR 185.

³⁰ LG München I, Case 42 O 14139/24, GRUR 2025, 1917 – GEMA/Open AI; see, e.g., Matthias Leistner, 'Memorisierungen und urheberrechtliche Vervielfältigungen in KI-Modellen: Das LG München I betritt Neuland' (2026) GRUR 185; Gierschmann M, 'Anmerkung' (2026) MMR 87; Hofmann F, 'Bloß weil »KI« draufsteht, ist nicht zwingend »KI« drin oder zu den urheberrechtlichen Grenzen von Large Language Models (LLMs)' (2026) ZUM 63; Pesch PJ, 'Anmerkung' (2025) KIR 457; Peukert A, 'Vervielfältigung im KI-Modell?' (2026) ZUM 17.

³¹ LG München I, Case 42 O 11483/24, GRUR 2025, 1917 – GEMA/OpenAI.

³² LG München I GRUR 2025, 1917 (1924) – GEMA/OpenAI, para. 176.

³³ LG Hamburg, Case 310 O 227/23, GRUR 2024, 1710; OLG Hamburg, Case 5 U 104/24, RD 2026, 97.

³⁴ See, e.g., Hacker, Philipp, Copyright, AI, and the Future of Internet Search before the CJEU: Reflections on *Like Company v Google*, VerfBlog, 2025/7/17, <https://verfassungsblog.de/copyright-ai-cjeu/>.

³⁵ See, e.g., Peter Henderson and others, 'Foundation models and fair use' (2023) 24 *Journal of Machine Learning Research* 1.

data must be ‘relevant, representative, free of errors and complete’ and must have ‘appropriate statistical properties’ with respect to the groups on which the system will be deployed.³⁶

For GPAI models, Article 53 requires providers to maintain detailed documentation on training data, disclose a summary of training data sources to the AI Office, and implement policies to comply with EU copyright law – including respect for rightholders’ opt-out declarations under the TDM exception. Article 55 imposes additional obligations on providers of GPAI models with systemic risk, including model evaluation, adversarial testing, incident reporting, and cybersecurity measures.³⁷ A model is presumed to pose systemic risk if trained with more than 10^{25} floating-point operations (FLOPs).³⁸

These provisions carry significant practical implications for AI training. Unlike the GDPR and copyright rules, the AI Act establishes a technology-specific regime. Whereas data protection and copyright law apply horizontally to any form of data processing or content use – regardless of the underlying technology – the AI Act targets the AI development pipeline directly and imposes obligations that are calibrated to the particular risks of machine learning systems. Data protection obligations arise whenever personal data are processed, and copyright constraints attach whenever protected works are reproduced – both irrespective of the system’s risk level. The AI Act, in contrast, modulates its requirements according to the risk classification of the AI system or the systemic significance of the GPAI model, which means that many AI applications fall outside its most demanding obligations altogether. This risk-based architecture complements the horizontal frameworks by adding a layer of regulation that is both sector-specific in its focus and graduated in its intensity.

Particularly in the case of copyright, the AI Act additionally directly connects to the TDM exceptions by requiring a dedicated compliance regime for all GPAI models – and by extending the reach of copyright law beyond the borders of the EU. More specifically, if models are trained outside of the EU, they need not respect EU copyright rules, such as the TDM exception with its upstart possibility, as a matter of international copyright law under the territoriality principle. However, if the provider wants to use the model or its output inside the EU, any hypothetical “breach” of or at least the lack of an effective compliance system concerning EU copyright rules, under the fictional premise that they were applicable to the training procedure, triggers a violation of the AI Act (but not copyright law).³⁹ With these rules, regulatory arbitrage is supposed to be prevented. To this effect, the AI Act globalizes EU copyright law, in a highly controversial and problematic move.⁴⁰ If other countries (e.g., UK, USA, India, China) were to adopt similar mechanisms, training companies including in the EU would soon

³⁶ See Art. 10 AI Act; Hacker/Engel/Hammer/Mittelstadt (n 1).

³⁷ See Art. 55 AI Act; see also Novelli and others (n 8).

³⁸ See Art. 51(2) AI Act; on the FLOPs threshold, see Novelli and others (n 8).

³⁹ See Recitals 105 and 106 AI Act.

⁴⁰ See, e.g., Stieper, Malte; Denga, Michael (2024) : The international reach of EU copyright through the AI Act, *Beiträge zum Transnationalen Wirtschaftsrecht*, No. 194, <https://doi.org/10.25673/116949>; Quintais JP, 'Copyright, the AI Act and extraterritoriality' (2025) Policy Brief The Lisbon Council; Matthias Leistner and Lucie Antoine, 'TDM and AI Training in the European Union—From ‘LAION’ to Possible Ways Ahead?' (2025) 74 GRUR International 1027, 1039 et seq.

be subject to a maze of different and partially conflicting copyright rules – the very scenario that the territoriality principle seeks to avoid in the first place.

C. The Model Level: Legal Status of Trained Models

The trained model itself – existing as a computational artifact with learned parameters⁴¹ – raises distinct legal questions under both copyright and data protection law.

I. Copyright: The Model as Reproduction

A first and particularly salient and far-reaching controversy concerns the question whether the model itself, or more precisely: certain combinations of parameters in the model, constitute a reproduction of a copyrighted work if that work can data be reproduced by the model as output. In November 2025, the UK High Court of Justice took a first stab at this question and, in its judgment in *Getty Images v Stability AI*, denied any reproductions in the model itself.⁴² It held that “in its final iteration Stable Diffusion does not store or reproduce any Copyright Works and nor has it ever done so.”⁴³

The second shot in this battle of the courts was handed down by the LG Munich I shortly after (see above). It established in its *GEMA v OpenAI* decision that memorization of copyrighted works in model parameters *does* constitute reproduction. The CJEU has held that reproduction encompasses any act by which a work is fixed in a medium that permits its communication to the public, even indirectly.⁴⁴ The LG München I applied this principle to find that neural network parameters, when they enable reconstruction of protected expression, satisfy the fixation requirement.

The court reasoned that the parameters constitute a ‘physical fixation’ of the work, rendering it ‘indirectly perceptible’ through user prompts.⁴⁵ This finding represents a significant departure from the assumption underlying the TDM exceptions that training merely extracts abstract information from works without reproducing their protected expression.⁴⁶

The court further held that the TDM exception under § 44b UrhG (implementing Article 4 CDSM Directive) does not cover memorization. The exception applies only to reproductions made ‘for the purpose’ of text and data mining.⁴⁷ Reproductions in the model serve no further analytical purpose and thus fall outside the exception’s scope. The court emphasized that the legislative premise of minimal harm to rightholders – expressed in Recital 17 CDSM Directive – does not apply where works are memorized and can be substantially reproduced.⁴⁸

⁴¹ See, e.g., Goodfellow/Bengio/Courville, *Deep Learning*, 2016, 6, 110, 116, 163.

⁴² [2025] EWHC 2863 (Ch).

⁴³ [2025] EWHC 2863 (Ch), para. 600.

⁴⁴ CJEU 16.7.2009 – case C-5/08, ECLI:EU:C:2009:465, para. 51 – *Infopaq*.

⁴⁵ LG München I GRUR 2025, 1917 (1924) – *GEMA/OpenAI*, para. 181 et seqq.

⁴⁶ LG München I GRUR 2025, 1917 (1925) – *GEMA/OpenAI*, para. 193.

⁴⁷ LG München I GRUR 2025, 1917 (1926 et seq.) – *GEMA/OpenAI*, para. 204 et seqq.

⁴⁸ LG München I GRUR 2025, 1917 (1926) – *GEMA/OpenAI*, para. 208; see Rec. 17 CDSM Directive.

The court also explicitly rejected an analogical extension of the TDM exception to cover memorization. Even if one assumed a legislative gap, no comparable interest situation exists: the TDM exception addresses scenarios where exploitation interests are not affected, whereas memorization directly impairs those interests.⁴⁹ Even where memorization is technically unavoidable, training on copyrighted works without authorization remains unlawful.⁵⁰

This analysis creates significant legal uncertainty, and challenges, for AI developers. The degree of memorization varies across models and training configurations, and identifying which specific works have been memorized may be technically challenging. The decision will prompt litigation strategies focused on output reproduction as evidence of memorization in the model itself. However, the ruling has not yet been confirmed by appeals courts. Again, its substance could be addressed by the CJEU in *Like Company v. Google*. In my view, there are good reasons to assume that parameter memorization constitutes a reproduction if the work can be produced by the model with reasonable efforts by a typical user.⁵¹ Needless to say that this would have vast implications for AI development, particularly concerning the need to obtain licenses, and the risk of damages and even deletion under current law.

II. Data Protection: The Model as Personal Data

Whether a trained model itself constitutes personal data under the GDPR depends on whether information relating to identifiable individuals can be extracted from it. Two technical phenomena are relevant: data leakage and model inversion attacks.⁵²

Data leakage occurs when models reproduce training data – including personal data – in their outputs. This is the data protection analogue to copyright memorization. Model inversion attacks involve the reconstruction of training data characteristics from model outputs or parameters, potentially enabling re-identification of individuals whose data was used in training.⁵³

Under the CJEU's *Breyer* test, data qualifies as personal if identification is possible using means reasonably likely to be used.⁵⁴ If model inversion attacks or data leakage render training data subjects identifiable, the model itself may constitute personal data, subjecting it to GDPR

⁴⁹ LG München I GRUR 2025, 1917 (1926) – GEMA/OpenAI, para. 208.

⁵⁰ Pesch in Spindler/Schuster/Kaesling, *Recht der elektronischen Medien*, 5th ed. 2026, UrhG § 44b para. 21; see also Pesch/Böhme, 'Artocalypse now? – Generative KI und die Vervielfältigung von Trainingsbildern' (2023) GRUR 997.

⁵¹ See also A Feder Cooper and James Grimmelmann, 'The files are in the computer: on copyright, memorization, and generative AI' (2025) 100 *Chi-Kent L Rev* 141; and n. 27 and 28.

⁵² See Michael Veale, Reuben Binns and Lilian Edwards, 'Algorithms that remember: model inversion attacks and data protection law' (2018) 376 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 20180083; Novelli and others (n 8); on data leakage, see Carlini and others, 'Extracting Training Data from Large Language Models' (2021) 30th *USENIX Security Symposium* 2633.

⁵³ Novelli and others (n 8); Fredrikson/Jha/Ristenpart, 'Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures' (2015) *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security* 1322.

⁵⁴ CJEU, Case C-582/14, ECLI:EU:C:2016:779 – *Breyer*.

requirements including the right to erasure. In this sense, the model parameters can be compared to encrypted information; model inversion or data leakage based on prompts are the keys to decrypting the information containing personal data. If (even accidental) use of the key is reasonably likely, the model parameters constitute personal data. At the current state of development of LLMs, with personal data can be extracted with simple prompts, there is likely that threshold is crossed and indeed the law needs to take seriously the finding that model parameters themselves, in many instances, will constitute personal data.

On the technical end, this prospect has motivated research into ‘machine unlearning’ techniques to remove the influence of specific data points from trained models.⁵⁵ However, it should also be on legislators’ minds as the Digital Omnibus seeks to address core tensions between AI and the GDPR (see below, F. II.). The Commission has included a reference to this issue in Recital 33 of the Digital Omnibus, but a clarification in the article itself would be much preferable and indeed doctrinally necessary.⁵⁶

D. The Output Level: Hallucinations and Memorization

Model outputs – the predictions, content, or decisions generated in response to user inputs – raise their own regulatory challenges, particularly concerning hallucinations and memorized content reproduction.

I. Hallucinations and Data Accuracy

LLMs are prone to generating factually incorrect or misleading statements – a phenomenon commonly termed ‘hallucination.’ Empirical studies have found that general-purpose LLMs hallucinate up to 60 to 80%⁵⁷ and even specialized legal models between 17 and 33% of the time when prompted with legal queries.⁵⁸ A particularly subtle form of inaccuracy is ‘careless speech’ – factually incorrect information or misattributed sources that requires domain expertise to detect.⁵⁹

Where AI outputs concern identifiable individuals, Article 5(1)(d) GDPR’s accuracy principle applies. Personal data must be ‘accurate and, where necessary, kept up to date.’ Hallucinated personal information – false statements about real individuals – may violate this principle and may give rise to erasure and rectification rights, as illustrated by complaints against OpenAI

⁵⁵ Villaronga and others, ‘Humans Forget, Machines Remember: Artificial Intelligence and the Right to be Forgotten’ (2018) 34 *Computer Law & Security Review* 304 (310).

⁵⁶ See Recital 33: “If removal would require disproportionate effort, notably where the removal of special categories of data memorised in the AI system or AI model would require reengineering the AI system or AI model, the controller should effectively protect such data from being used to infer outputs, being disclosed or otherwise made available to third parties.” This seems to displace erasure and correction rights – again, any such groundbreaking amendments must be made in an article, not in a formally non-binding recital.

⁵⁷ Matthew Dahl and others, ‘Large legal fictions: Profiling legal hallucinations in large language models’ (2024) 16 *Journal of Legal Analysis* 64.

⁵⁸ Varun Magesh and others, ‘Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools’ (2025) 22 *Journal of Empirical Legal Studies* 216.

⁵⁹ Sandra Wachter, Brent Mittelstadt and Chris Russell, ‘Do large language models have a legal duty to tell the truth?’ (2024) 11 *Royal Society Open Science* 240197.

by noyb, the NGO led by privacy activist Max Schrems.⁶⁰ Hallucinations may also give rise to erasure and rectification rights under Articles 16 and 17 GDPR.

The question of *responsibility* under Article 5(2) GDPR requires, however, careful analysis along the AI value chain. Under the CJEU's functional approach, the responsible controller is the entity that actually determines the purposes and means of data processing.⁶¹ In the AI context, this determination must be made at each stage: the provider of the base model (e.g., GPT-4 or Llama), the deployer who integrates it into a specific application, and potentially even the end user may each bear a share of controller responsibility, depending on their respective influence over the processing operations that produce the output.

The *scope* of the accuracy principle, however, is not unlimited. Recital 39 GDPR already signals that the principle must be read in light of the purposes for which data are processed, and that the rectification deletion of inaccurate data is subject to a reasonableness requirement. Hence, a balancing exercise with conflicting fundamental rights is necessary – in particular, the rights of LLM providers and deployers, which may include the freedom to conduct a business (Article 16 CFR) and the freedom of expression and information (Article 11 CFR).

In my view, this balancing exercise calls for a *de minimis* threshold. Not every factual inaccuracy in an AI-generated text rises to the level of a violation of Article 5(1)(d) GDPR. Insignificant errors – such as an incorrect date of birth embedded in a simple, low-stakes text output – should fall below this threshold. By contrast, materially significant inaccuracies do violate the accuracy principle. This assessment operates on two dimensions. First, certain inaccuracies are so grave in absolute terms that they constitute a violation regardless of context – for instance, the false portrayal of an individual as a child murderer, as happened in a highly publicized Norwegian case. Second, other inaccuracies become significant only in a specific relational context – an incorrect date of birth, for example, may be trivial in a general text but constitutes a meaningful violation when it appears in an AI-agent-generated ID card application, where accuracy is functionally essential even for otherwise rather innocuous data points. This graduated approach runs parallel to the well-established doctrinal evaluation of the severity of interference with personal rights (*Persönlichkeitsrechte*), which likewise distinguishes between trivial and substantial impairments depending on context and gravity, to which we now briefly turn.

Beyond data protection, liability for hallucinations may arise under personality rights, defamation law, or contract law.⁶² The Bundesgerichtshof's *Autocomplete* decision established that operators of automated systems may be liable for unlawful outputs where they fail to take reasonable precautions, particularly upon gaining knowledge of specific harms.⁶³ In a parallel

⁶⁰ noyb – European Center for Digital Rights, Complaint against OpenAI (redacted version) (29 April 2024).

⁶¹ CJEU, Case C-40/17, Fashion ID.

⁶² See, e.g., Reuben Binns and Lilian Edwards, 'Reputation Management in the ChatGPT Era' in Philipp Hacker and others (ed), *Oxford Handbook of the Foundations and Regulation of Generative AI* (Oxford University Press 2025); for a US perspective, see Peter Henderson, Tatsunori Hashimoto and Mark Lemley, 'Where's the Liability in Harmful AI Speech?' (2023) *Journal of Free Speech Law* 589.

⁶³ See BGH, Case VI ZR 269/12, NJW 2013, 2348 – *Autocomplete*; the following analysis draws on Philipp Hacker and others, 'Generative discrimination: What happens when generative AI exhibits bias, and what can be

case, Bettina Wulff, the wife of the then Federal President of Germany, brought proceedings against Google because the search engine's autocomplete function prominently suggested the word "prostitute" in connection with her name – a reflection of tabloid speculation about her past prior to her husband's assumption of the presidency. The BGH held that the search engine operator could incur liability for such automated suggestions if it failed to adopt adequate measures to prevent outputs that violate personality rights or are defamatory in nature. In particular, the court imposed a reactive obligation: once the operator obtains knowledge of a harmful output, it must take steps to prevent its recurrence.

This precedent carries direct implications for the developers of large language models. The autocomplete function at issue in the BGH case relied, in effect, on what amounted to a small language model. The doctrinal logic, however, scales: if the operator of a comparatively simple prediction system bears responsibility for rights-violating outputs, the same must apply – a fortiori – to operators of far more powerful generative AI systems. The ruling thus supports the conclusion that developers face potential direct liability for AI-generated content if they neglect to implement appropriate safeguards. Such safeguards may include content moderation measures integrated into the training process, as well as responsive corrective action upon receipt of notification about problematic outputs.

With respect to *deployers* of AI technologies, courts can be expected to apply an analogous framework, adjusted *mutatis mutandis* to the deployer's specific role. In the *Autocomplete* case itself, Google occupied the position of both developer and deployer simultaneously. In my view, both developers and deployers must undertake what is reasonable within their power to prevent generative AI outputs from violating personality rights. Deployers, in particular, may avoid liability by establishing a robust compliance framework. The essential components of such a system include the deployment of AI with built-in moderation guardrails, regular proactive review of AI outputs, and an effective notice-and-takedown procedure that allows for the swift removal and future prevention of harmful content – for example, through the blocking of specific prompts or output patterns. Where both developer and deployer fail in their respective duties, joint and several liability represents the appropriate doctrinal consequence.

Going forward, the key question is whether models might be subject to deletion when they cannot be prevented from hallucinating in high-stakes cases. Such a doctrine could draw on the accuracy principle in data protection law in conjunction with the fact that part of the model itself might be considered personal data, since trained weights encode information about individuals from the training corpus in a form that can be – and regularly is – reproduced in outputs (see above).

Drawing an analogy to pharmaceutical liability under § 84 German Pharmaceutical Act (AMG, *Arzneimittelgesetz*), one might conceive of AI-specific strict liability, with deletion as an ultima ratio, but only where hallucination rates exceed scientifically acceptable levels and cause serious harm (§ 84(1)(1.) AMG). Section 84 AMG imposes strict, fault-independent liability

done about it' in Philipp Hacker and others (ed), *The Oxford Handbook of Foundations and Regulation of Generative AI* (2025).

on pharmaceutical manufacturers when a medicinal product, used as intended, causes harm that exceeds the bounds of what is deemed acceptable in light of the current state of medical science. The underlying rationale is instructive: the legislature permits socially beneficial but inherently risky products to be placed on the market under an overall risk-benefit analysis, even where harm in individual cases remains unavoidable. In return for this permission, the manufacturer bears liability without fault for damages that materialize within the risk corridor that society has accepted. Transposed to AI, a parallel regime could allow the deployment of generative models whose outputs carry residual hallucination risks – provided these risks remain within scientifically defensible bounds – while imposing strict liability on developers or deployers when those bounds are exceeded and serious harm results. In effect, this would require a modification of Articles 16-18 and 82 GDPR, subjecting them to the above caveat in cases of AI-generated hallucinations.

II. Memorization and Copyright Infringement

The *GEMA v OpenAI* decision also addressed the copyright status of outputs that reproduce memorized works. The court found that such outputs infringe the reproduction right (§ 16 UrhG) and, where made available through online chatbots, the right of communication to the public (§ 19a UrhG).⁶⁴

No copyright exception justified these uses. The quotation exception (§ 51 UrhG) requires a purpose of intellectual engagement with the quoted work, which AI models – lacking subjective intention – cannot fulfill. The pastiche exception (§ 51a UrhG) requires artistic engagement, which is similarly inapplicable. The private copying exception (§ 53 UrhG) benefits only natural persons and cannot be invoked by corporate defendants. Overall, liability for output is where the case for copyright enforcement is clearest – and where new remuneration rules may be normatively anchored.⁶⁵

E. Evaluation: Technology Neutrality Versus Technology-Specific Regulation

The regulatory landscape for AI in the EU reflects a combination of technology-neutral and technology-specific approaches.⁶⁶ Evaluating their relative merits illuminates the path for future regulatory development.

I. Technology-Neutral Regulation

The GDPR and traditional copyright law exemplify technology-neutral approaches. Such regulation offers several advantages. It provides broad coverage, applying to AI without requiring specific legislative amendments. It demonstrates resilience against technological change, as principles like data minimization, accuracy, or the reproduction right adapt to new

⁶⁴ LG München I GRUR 2025, 1917 (1931 et seq.) – GEMA/OpenAI, paras. 253 et seqq.

⁶⁵ See, e.g., Martin Senftleben, 'Generative AI and author remuneration' (2023) 54 IIC-International Review of Intellectual Property and Competition Law 1535.

⁶⁶ Hacker/Engel/Hammer/Mittelstadt, 'Introduction' in *The Oxford Handbook on Generative AI and Law*, OUP 2025.

contexts. It promotes regulatory consistency by subjecting AI to the same rules as analogous technologies.⁶⁷

However, technology-neutral frameworks struggle to address AI-specific phenomena in the sense that they do not necessarily deliver the ex-ante legal clarity that remains a desideratum both for regulatees and consumers. The GDPR’s consent and legitimate interest mechanisms were not designed for large-scale training on web-scraped data. Copyright’s TDM exception was not drafted with memorization in mind. Applying general principles to novel technical realities produces uncertainty, which can only overtime be resolved by guidelines and court decisions, as the *GEMA v OpenAI* litigation illustrates – particularly in its stark contrast to the *Getty v Stability AI* judgment.

II. Technology-Specific Regulation

The AI Act represents the EU’s primary technology-specific intervention. Such regulation can address AI’s distinctive characteristics – its opacity, scale, adaptiveness, and emergent capabilities – with tailored requirements. The GPAI provisions on transparency, risk management, and cybersecurity respond to specific concerns that general frameworks may only inadequately address.⁶⁸

However, technology-specific regulation risks obsolescence as technology evolves. It may create gaps where AI applications fall outside defined categories (e.g., “GPAI Agents” used for a significant amount of downstream applications⁶⁹) or generate compliance burdens that, de facto, disadvantage EU-based developers. The appropriate scope of technology-specific rules – whether they should target only high-risk or systemically risky AI, or extend more broadly – remains contested. While, at the moment of writing this manuscript, large parts of the AI Act are not applicable yet, there are signs that some of the most significant negative AI externalities (hallucinations; unauthorized copying) can be adequately captured by technology-neutral regimes. Product liability law provides another layer of protection. Hence, the most important addition in the AI Act may be the specific rules on GPAI models in its Chapter V. They constitute a much-needed AI safety regime, combined with transparency and extension of copyright rules that address highly GPAI-specific risks for which, indeed, technology-specific rules were and continue to be needed.

F. Policy Proposals

The analysis above suggests several areas where regulatory reform could enhance the coherence and effectiveness of AI governance in the EU.

I. Training Level Reforms

⁶⁷ On technology neutrality in EU law, see also CJEU, Case C-403/08, ECLI:EU:C:2011:43 – Premier League.

⁶⁸ Wachter, 'Limitations and Loopholes in the EU AI Act and AI Liability Directives: What This Means for the European Union, the United States, and Beyond' (2024) 26 *Yale Journal of Law and Technology* 671.

⁶⁹ See Holweg/Hacker, *Towards a pragmatic regulation of AI agents*, forthcoming.

First, regarding data protection, the GDPR Omnibus proposals merit refinement. Article 88c should specify high-impact sectors – such as medicine or employment – where the legitimate interest basis is particularly appropriate. Article 9(2)(k) should be limited to the training and modification of AI models, and require state-of-the-art measures for detection and removal of sensitive data, with proportionate safeguards against memorization (see above).⁷⁰

Additionally, the exemption should address the ‘model-as-personal-data’ problem by explicitly covering the existence of models containing personal data, with appropriate safeguards including machine unlearning capabilities.⁷¹

Second, regarding copyright, it should be clarified that AI training – but not the model itself – falls under the TDM exception. Furthermore, the principle of equal treatment under Article 20 of the Charter supports extending the Article 3 CDSM Directive exception to commercial research units that reinvest profits in research, subject to appropriate conditions.⁷² This would align copyright law with the treatment of commercial research under Article 89 GDPR.⁷³

II. Model Level Reforms

At the model level, the *GEMA v OpenAI* decision creates urgent need for legislative clarification. If memorization is unavoidable with current technology, training on copyrighted works cannot proceed under the TDM exception without authorization. This may necessitate a new exception covering training with remuneration, potentially administered through collective rights organizations.⁷⁴ Such an approach would balance innovation interests with fair compensation for creators.⁷⁵

Similarly, data protection law requires explicit attention to the model-as-personal-data problem. As seen, when an AI model has been trained on personal data – particularly sensitive personal data – the question arises whether the model *itself* constitutes personal data under the GDPR. This is not merely a theoretical concern: if a model memorizes or allows the extraction of training data, it may qualify as personal data in its own right, which would subject the mere storage and deployment of the model to the full range of GDPR obligations, including the right to erasure under Article 17. Specific guidance from the European Data Protection Board,

⁷⁰ See proposed Art. 9(2)(k) and Art. 9(5); cf. Art. 10(5) AI Act.

⁷¹ On machine unlearning, see Bourtole and others, 'Machine Unlearning' (2021) IEEE Symposium on Security and Privacy 141; Novelli and others (n 8).

⁷² On the status quo bias in copyright opt-outs, see Samuelson/Zeckhauser, 'Status Quo Bias in Decision Making' (1988) 1 Journal of Risk and Uncertainty 7; Hacker (n 3) 280.

⁷³ See Margoni/Kretschmer, 'The Text and Data Mining Exception in the Proposal for a Directive on Copyright in the Digital Single Market: Why it is not what EU copyright law needs' (2018) Working Paper, 4 et seqq.; Geiger/Frosio/Bulayenko (n 17) 20 et seqq.; Ducato/Strowel, 'Limitations to Text and Data Mining and Consumer Empowerment: Making the Case for a Right to Machine Legibility' (2019) 50 IIC 649 (666).

⁷⁴ See again Senftleben (n 62); IUM-Symposium, Kollektive Vergütungsmodelle für KI-Nutzungen: Wege zu einem fairen Interessenausgleich (14 November 2025), <https://www.urheberrecht.org/events/20251114.php> (last accessed: 03.03.2026).

⁷⁵ Hacker, 'Copyright, AI, and the Future of Internet Search before the CJEU' (2025) Verfassungsblog, <https://verfassungsblog.de/copyright-ai-cjeu/> (last accessed: 03.03.2026).

beyond the more generic comments delivered in its Opinion on AI models,⁷⁶ or preferably legislative clarification should address when models constitute personal data, how the right to erasure applies to trained models, and what technical measures – including machine unlearning – satisfy GDPR obligations.

The proposed Article 9(2)(k) GDPR offers an opportunity to address this problem directly, but only if its scope is formulated with sufficient precision. As currently drafted, the provision covers processing "in the context of" AI development, which does not clearly extend to situations in which the model itself is the object of regulatory concern. A more adequate formulation would provide that the exemption applies where processing is necessary for the development or modification of an AI model, *or* where processing concerns the existence of an AI model as such – in both cases subject to the safeguards referred to in Article 9 paragraph 5 and limited to high-impact sectors listed in a dedicated annex. An accompanying recital should clarify that "processing concerns the existence of an AI model" refers specifically to situations in which the model itself or parts of it are considered personal data – for instance, because it has memorized training samples or because individual data points can be reconstructed from model parameters. Without such a provision, the legal status of trained models would remain uncertain, and data controllers would face the paradoxical situation that they possess a lawful basis to process sensitive data *during* training but lack a clear legal ground for the continued existence of the resulting model. This gap would undermine the practical effectiveness of the exemption and create a persistent source of legal risk for AI developers operating in the European Union.

III. Output Level Reforms

For hallucinations, current legal frameworks are insufficiently adapted to the cumulative, long-term risks of AI-generated misinformation. As mentioned, a sector-specific strict liability regime, modeled on § 84 AMG for pharmaceutical liability, could apply where AI outputs cause serious harm and hallucination rates exceed scientifically acceptable levels. The AI Act's implementing measures or secondary data protection legislation could specify acceptable error rates for particular deployment contexts.

For memorized output reproduction, deployers must exercise due diligence to prevent and respond to copyright infringement via reproduction. This is a core task that the law cannot exempt them from.

G. Structural and Overarching Considerations

It bears noting that a peculiar structural symmetry exists between the most urgent data protection and copyright challenges that generative AI poses. In both domains, the core difficulty resides in the fact that information encoded in an AI model's parameters may constitute legally protected content – personal data under the GDPR, copyrighted works under

⁷⁶ EDPB, Opinion 28/2024 on certain data protection aspects related to the processing of personal data in the context of AI models, 2024; see also German Data Protection Authorities, Orientierungshilfe KI und Datenschutz, Version 1.0, 6 May 2024.

copyright law – even though this concerns only fragments of the model. It is technically near-infeasible to identify and excise only the infringing portions of a trained model’s weights while the remainder of the model stays intact. While the framers of the data protection and copyright frameworks likely did not specifically conceive of this extension of legal protection to highly dimensional generative AI tensors, the open wording of the technology-neutral rules has precisely this effect. Yet, the entanglement of protected and non-protected information within the parameter space of large neural networks renders surgical, *ex post* correction an extraordinarily difficult undertaking from an engineering perspective.

The practical challenge that flows from this structural parallel is therefore similar across both regimes: the legal system must find ways to keep large generative AI models deployable and in socially beneficial, productive use while simultaneously preventing ongoing rights violations. The alternative – to mandate deletion of the entire model – would carry significant adverse consequences not only from an economic standpoint, given the enormous investment in training infrastructure and compute, but also from an ecological one, since the energy expenditure required to train large-scale models is substantial and cannot be recovered once a model is destroyed. While a possibility *ultima razione*, such a consequence should generally be avoided.

The remedies, however, may well diverge between the two regimes despite these structural parallels. In data protection law, personal data is not easily monetized by the data subject; there is no functioning market that would allow individuals to license their personal information to AI developers at an efficient price. The balance must therefore be struck under a property rule framework – that is, through categorical permissions and prohibitions – by the legislature’s decision to carve out certain exemptions for AI training, as the proposed Digital Omnibus Directive foresees. Copyright, by contrast, presents a different constellation. Here, the tensions between rightsholder interests and AI development can most effectively and most efficiently be resolved through a liability rule: the law would permit AI training and even the memorization of protected works within models, but in return impose an output-based remuneration framework that compensates authors, through collective rights organisations, when the model generates content that draws on their works.

H. Conclusion

The EU’s regulatory architecture for AI reflects both the capabilities and limitations of existing legal frameworks as applied to novel technology. The GDPR and copyright law provide broad coverage but struggle to address AI-specific phenomena such as large-scale training, memorization, and hallucination. The AI Act introduces technology-specific rules but cannot alone ensure coherent governance across all stages of the machine learning pipeline.

The *GEMA v OpenAI* decision marks a turning point for AI copyright law, as it suggests (even though it may ultimately be reversed) that memorization constitutes reproduction outside the scope of the TDM exception. Combined with the expansive interpretation of sensitive data in *Meta v Bundeskartellamt*, these developments signal increased legal risk for AI training operations in the EU that involve copyrighted content or personal data.

At the same time, the regulation of generative AI demands a differentiated approach that resists the temptation of uniform solutions. The legal challenges span data protection, personality rights, liability, and copyright – and while these fields share the structural feature that protected information is deeply embedded in model parameters and near-impossible to excise surgically, they call for distinct regulatory instruments. A *de minimis* threshold for data accuracy, strict liability modeled on pharmaceutical law for serious harms caused by hallucinations, property rules for data protection exemptions, and liability rules for copyright remuneration together form the contours of a regulatory architecture that takes both innovation and fundamental rights seriously.

The optimal regulatory approach thus combines technology-neutral baseline protection with targeted technology-specific interventions where AI raises distinctive concerns. The GDPR Omnibus proposals represent a step in this direction for data protection. For copyright law, the most effective and efficient adaptation lies in output-based remuneration mechanisms – administered by collective rights organizations – that permit AI training and even memorization while they ensure fair compensation for authors.

As generative AI continues to evolve, the regulatory architecture must remain adaptive while it preserves fundamental commitments to rights protection, fair compensation for creators, and the promotion of beneficial innovation. The current moment presents both challenges and opportunities for this balance – and the decisions made now will shape the trajectory of AI governance, but also development and deployment, in Europe for years to come.

Bibliography

Baumann M, 'Generative KI und Urheberrecht – Urheber und Anwender im Spannungsfeld' (2023) NJW 3673

Binns R and Edwards L, 'Reputation Management in the ChatGPT Era' in others PHa (ed), Oxford Handbook of the Foundations and Regulation of Generative AI (Oxford University Press 2025)

Bourtoule/Chandaserakan/Choquette-Choo/Jia/Travers/Zhang 'Machine Unlearning' (2021) IEEE Symposium on Security and Privacy (SP) 141;

Burrell, 'How the Machine 'Thinks': Understanding Opacity in Machine Learning Algorithms' (2016) 3(1) Big Data & Society 1;

Carlini and others, 'Extracting Training Data from Large Language Models' (2021) 30th USENIX Security Symposium 2633;

Cooper AF and Grimmelmann J, 'The files are in the computer: on copyright, memorization, and generative AI' (2025) 100 Chi-Kent L Rev 141

Dahl M and others, 'Large legal fictions: Profiling legal hallucinations in large language models' (2024) 16 Journal of Legal Analysis 64

de la Durantaye K, '"Garbage in, garbage out" – Die Regulierung generativer künstlicher Intelligenz durch Urheberrecht' (2023) ZUM 645

de la Durantaye, 'Garbage In, Garbage Out – Die Regulierung generativer KI durch Urheberrecht' (2023) Working Paper;

Dermawan, 'Text and Data Mining Exceptions in the Development of Generative AI Models' (2024) 27 Journal of World Intellectual Property 44;

Dornis TW, 'Generatives KI-Training und der TDM-Trugschluss' (2024) GRUR 1676

Dornis/Stober, Urheberrecht und Training generativer KI-Modelle, 2024;

Ducato/Strowel, 'Limitations to Text and Data Mining and Consumer Empowerment' (2019) 50 IIC 649;

European Commission. (2025). Guidelines on the definition of an artificial intelligence system (C(2025) 5053 final.

European Law Institute, Commission Guidelines on the Application of the Definition of an AI System and the Prohibited AI Practices Established in the AI Act. Response of the European Law Institute, 2024

Expert Group on Liability and New Technologies, Liability for Artificial Intelligence and Other Emerging Digital Technologies, 2019;

Fredrikson, M., Jha, S. and Ristenpart, T. (2015) Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures. Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, Denver, 12-16 October 2015, 1322-1333;

Gebehenne/Siebler/Hennemann, Der Digital Omnibus – Grundstruktur, Einordnung und Rahmenbedingungen der Vorschläge der EU-Kommission für eine Vereinfachung im Datenrecht, EuDIR 2026, 10

- Geiger/Frosio/Bulayenko, 'The Exception for Text and Data Mining (TDM) in the Proposed Directive on Copyright in the Digital Single Market – Legal Aspects' (European Parliament, PE 604.941, 2018);
- Goodfellow/Bengio/Courville, *Deep Learning*, MIT Press 2016;
- Griffiths, 'The Three-Step Test in European Copyright Law' (2009) Working Paper;
- Grimm S and Münster LM, 'Training von KI in der EU: Fragen zum Nutzungsvorbehalt beim Text und Data Mining sowie zur Beweislast' (2025) GRURPrax 443
- Gurney, 'Sue My Car Not Me' (2013) U. Ill. J.L. & T. 247;
- Hacker and others (eds), *The Oxford Handbook on Generative AI and Law*, Oxford University Press 2025
- Hacker P and others, 'Generative discrimination: What happens when generative AI exhibits bias, and what can be done about it' in others PHa (ed), *The Oxford Handbook of Foundations and Regulation of Generative AI* (2025)
- Hacker P, 'Comments on the final trilogue version of the AI Act' (2024) Available at SSRN 4757603
- Hacker, 'A Legal Framework for AI Training Data' (2021) 13 *Law, Innovation and Technology* 257;
- Hacker, AI & data protection: applications to the AV sector, in: *Artificial intelligence in the audiovisual sector*, Report for the Council of Europe (October 22, 2024), <https://www.obs.coe.int/en/web/observatoire/-/new-report-the-challenges-of-ai-for-the-audiovisual-sector-and-the-role-european-legislation-is-playing>
- Hacker, 'Copyright, AI, and the Future of Internet Search before the CJEU' (2025) *Verfassungsblog*;
- Hacker/Cordes/Rochon, 'Regulating Gatekeeper AI and Data' (2024) 15 *European Journal of Risk Regulation* 49;
- Hacker/Engel/Hammer/Mittelstadt, 'Introduction' in *The Oxford Handbook on Generative AI and Law*, OUP 2025;
- Henderson P and others, 'Foundation models and fair use' (2023) 24 *Journal of Machine Learning Research* 1
- Henderson P, Hashimoto T and Lemley M, 'Where's the Liability in Harmful AI Speech?' (2023) *Journal of Free Speech Law* 589
- Hofmann F, 'Zehn Thesen zu Künstlicher Intelligenz (KI) und Urheberrecht' (2024) WRP 11
- Konertz R and Schönhof R, 'Vervielfältigungen und die Text- und Data-Mining-Schranke beim Training von (generativer) Künstlicher Intelligenz' (2024) WRP 289
- Leistner M and Antoine L, 'TDM and AI Training in the European Union—From 'LAION'to Possible Ways Ahead?' (2025) 74 *GRUR International* 1027
- Leistner M, 'Memorisierungen und urheberrechtliche Vervielfältigungen in KI-Modellen: Das LG München I betritt Neuland' (2026) *GRUR* 185
- Leistner M, 'TDM und KI-Training in der Europäischen Union, Erste Fingerzeige des LG Hamburg im "LAION"-Urteil' (2024) *GRUR* 1665

- Lewis and others, 'Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks' (2020) 33 *NeurIPS* 9459;
- Magesh V and others, 'Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools' (2025) 22 *Journal of Empirical Legal Studies* 216
- Malgieri G and Comandé G, 'Sensitive-by-distance: quasi-health data in the algorithmic era' (2017) 26 *Information & Communications Technology Law* 229
- Margoni T and Kretschmer M, 'A deeper look into the EU text and data mining exceptions: harmonisation, data ownership, and the future of technology' (2022) 71 *GRUR international* 685
- Margoni/Kretschmer, 'The Text and Data Mining Exception' (2018) Working Paper;
- Mayer-Schönberger/Cukier, *Big Data: A Revolution That Will Transform How We Live, Work, and Think*, Houghton Mifflin Harcourt 2013
- Mészáros/Ho, 'Big Data and Scientific Research' (2021) 21 *IJLIT* 403;
- Novelli/Casolari/Hacker/Spedicato/Floridi, 'Generative AI in EU Law' (2024) 55 *Computer Law & Security Review* 106066;
- Obergfell, 'Big Data und Urheberrecht' in FS Büscher, 2018, 223;
- Pesch PJ and Böhme R, 'Artpocalypse now? – Generative KI und die Vervielfältigung von Trainingsbildern' (2023) *GRUR* 997
- Pesch, in G Spindler, F Schuster and T Kaesling (eds), *Recht der elektronischen Medien* (5th edn, C.H. Beck 2026);
- Pesch/Böhme, 'Artpocalypse now?' (2023) *GRUR* 997;
- Quintais JP, 'Copyright, the AI Act and extraterritoriality' (2025) Policy Brief The Lisbon Council.
- Quintais JP, 'Generative AI, copyright and the AI Act' (2025) 56 *Computer Law & Security Review* 106107
- Raue, 'Rechtssicherheit für datengestützte Forschung' (2019) *ZUM* 684;
- Rosati E, 'The exception for text and data mining (TDM) in the proposed Directive on Copyright in the Digital Single Market: technical aspects' (2018) European Parliament
- Rosati, 'The Exception for Text and Data Mining (TDM) in the Proposed Directive on Copyright in the Digital Single Market – Technical Aspects' (European Parliament, PE 604.942, 2018);
- Russell SJ and Norvig P, *Artificial Intelligence: A Modern Approach* (4th Global edn, Pearson Education 2022)
- Sartor, Giovanni, and Francesca Lagioia. *The impact of the General Data Protection Regulation (GDPR) on artificial intelligence. Study for the EPRS*, 2020.
- Schack H, 'Auslesen von Webseiten zu KI-Trainingszwecken als Urheberrechtsverletzung de lege lata et ferenda' (2024) *NJW* 113
- Schönberger, 'Artificial Intelligence in Healthcare' (2019) 27 *IJLIT* 171;

- Schuhmacher/Fatalin, in: Praxishandbuch KI 2023, 204;
- Schwartmann, Rolf, Zulässigkeit von Datenverarbeitungen in der KI-Entwicklung – Der Stand zwischen OLG Köln und Digitalem Omnibus, EuDIR 2026, 3
- Senftleben M, ‘Generative AI and author remuneration’ (2023) 54 IIC-International Review of Intellectual Property and Competition Law 1535
- Senftleben, ‘Generative AI and Author Remuneration’ (2023) 54 IIC 1535;
- Sesing-Wagenpfeil, ‘KI-Training und Urheberrecht’ (2024) ZGE 212;
- Spindler, ‘Die neue Urheberrechts-Richtlinie der EU’ (2019) CR 277;
- Spindler, ‘Text und Data Mining’ (2016) GRUR 1112;
- Stieper, Malte; Denga, Michael (2024) : The international reach of EU copyright through the AI Act, Beiträge zum Transnationalen Wirtschaftsrecht, No. 194, <https://doi.org/10.25673/116949>
- Tyagi K, ‘Copyright, text & data mining and the innovation dimension of generative AI’ (2024) 19 Journal of Intellectual Property Law & Practice 557
- Veale M, Binns R and Edwards L, ‘Algorithms that remember: model inversion attacks and data protection law’ (2018) 376 Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences 20180083
- Villaronga and others, ‘Humans Forget, Machines Remember’ (2018) 34 CLSR 304;
- Wachter S, Mittelstadt B and Russell C, ‘Do large language models have a legal duty to tell the truth?’ (2024) 11 Royal Society Open Science 240197
- Wachter, ‘Limitations and Loopholes in the EU AI Act’ (2024) 26 Yale J.L. & Tech. 671;
- Wachter/Mittelstadt, ‘A Right to Reasonable Inferences’ (2019) 2 Columbia Business Law Review 494;
- Wendehorst, Christiane, Keep Calm and Read the Proposals: A Legal Lens on the Digital Omnibus, EuCML 2025, 277