

AIREG-BENCH: BENCHMARKING LANGUAGE MODELS THAT ASSESS AI REGULATION COMPLIANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

As governments move to regulate AI, there is growing interest in using Large Language Models (LLMs) to assess whether or not an AI system complies with a given AI Regulation (AIR). However, there is presently no way to benchmark the performance of LLMs at this task. To fill this void, we introduce **AIReg-Bench**: the first benchmark dataset designed to test how well LLMs can assess compliance with the EU AI Act (AIA). We created this dataset through a two-step process: (1) by prompting an LLM with carefully structured instructions, we generated 120 technical documentation excerpts (samples), each depicting a fictional, albeit plausible, AI system — of the kind an AI provider might produce to demonstrate their compliance with AIR; (2) legal experts then reviewed and annotated each sample to indicate whether, and in what way, the AI system described therein violates specific Articles of the AIA. The resulting dataset, together with our evaluation of whether frontier LLMs can reproduce the experts’ compliance labels, provides a starting point to understand the opportunities and limitations of LLM-based AIR compliance assessment tools and establishes a benchmark against which subsequent LLMs can be compared. The dataset and evaluation code are available at <https://anonymous.4open.science/r/aireg-bench-5259/>.

1 INTRODUCTION AND PROBLEM STATEMENT

Across the world, AI Regulation (AIR) initiatives are either under development or have graduated the legislative process and gone into effect (Sloane & Wüllhorst, 2025; Chun et al., 2024; Alanoca et al., 2025). For both the regulators who enforce these regulations and the regulated parties who must comply with them, *compliance assessments*, whereby an AI system is evaluated for its compliance with respect to an AIR, play a pivotal role (Mökander et al., 2021; Ada Lovelace Institute, 2024; Anderljung et al., 2023; Raji et al., 2022; Reuel et al., 2024a). For example, the European Union’s AI Act (AIA) — dubbed “the world’s first comprehensive AI law” (European Parliament, 2024) — requires that providers of high-risk AI systems conduct such an assessment before putting their products on the market in the EU (EU, 2024, Art. 43).

Despite their importance, however, certain AIR compliance assessments remain costly and time-consuming (Koh et al., 2024; Costanza-Chock et al., 2022; Sovrano et al., 2025). For example, some estimate that these assessments can take up to two-and-a-half days (European Commission, 2021) and cost EUR 7,500 for each AI system (Haataja & Bryson, 2021), accounting for up to 17% of the total expense of an AI project (Laurer et al., 2021). These high costs may contribute to a level of regulatory overhead that some have called unsustainable for AI providers and regulators alike (Laurer et al., 2021; Gikay, 2024; Reuel et al., 2024b; Koh et al., 2024; Molnar, 2024; Micklitz & Sartor, 2025) and, since it disproportionately affects small and medium-sized enterprises due to their lower resources (Stampernas & Lambrinoudakis, 2025), a potential hazard for fair competition (Martens, 2024; Gazendam et al., 2023; Wu & Liu, 2023; Guha et al.; Iliasova et al., 2025).

This may help explain why there is growing interest in, and experimentation with, using Large Language Models (LLMs) to perform (or, at least, streamline) AIR compliance assessments (Micklitz & Sartor, 2025; Li et al., 2025; Sovrano et al., 2025; Davvetas et al., 2025; Kővári et al., 2025; Makovec et al., 2024; Marino et al., 2024). And yet, there is still no standardized method for quantitatively evaluating and comparing the performance of LLMs at this particular task, creating uncertainty about the extent to which LLMs can be entrusted with it (Davvetas et al., 2025).

To fill this void, we present AIReg-Bench: an open dataset for benchmarking the performance of LLMs at AIA compliance assessments. This dataset, which is available at <https://anonymous.4open.science/r/aireg-bench-5259/>, consists of 120 technical documentation excerpts (i.e., details on system development and testing procedures) (Sovrano et al., 2025; Königstorfer & Thalmann, 2022). Each one provides information about a fictional, albeit plausible, AI system — specifically, a high-risk AI (HRAI) system under the AIA’s risk-based approach (EU, 2024, Art. 6). The samples in this dataset (i.e., the excerpts) are generated by an LLM-based technique, described in Section 2, allowing us to create diverse samples efficiently and at scale. As outlined in Section 3, each sample is then labeled by legal experts to indicate whether, and in what way, the system described therein violates specific Articles of the AIA.

To showcase AIReg-Bench at work, we evaluated 10 frontier LLMs. Our findings indicate that some LLMs very closely approximate human expert judgments about the compliance (or lack thereof) of the excerpts in our dataset, such as Gemini 2.5 Pro (Comanici et al., 2025), which achieves a rank correlation of 0.856, as shown in Table 3.

In short, our contributions include:

- **Sample generation pipeline:** An open source repository for the LLM-based generation of plausible AIA technical documentation excerpts, which we use to generate the samples in AIReg-Bench, and which can be reused for other AI compliance evaluation and training initiatives.
- **Dataset:** The above pipeline is used to generate a distribution of samples that are then annotated by legal experts to create the AIReg-Bench open benchmark dataset, which can be used today to evaluate the effectiveness of LLMs at AIA compliance assessments — and which, in the future, is extensible to other AIR.
- **Experiments:** The first application of the benchmark to evaluate the performance of 10 frontier LLMs at the task of AIA compliance assessments.

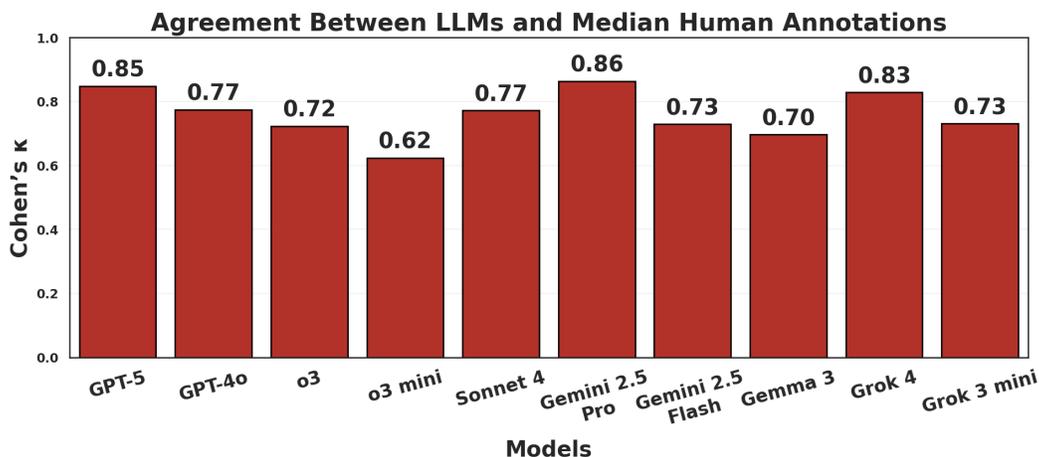


Figure 1: **Cohen’s κ (quadratically weighted) scores across frontier language models**, showing the level of agreement on compliance judgments (on a 1-5 Likert scale) between these models and the median legal expert in our team, taken over the entire AIReg-Bench dataset.

2 SAMPLE GENERATION PIPELINE (SOLUTION PART I)

Our first contribution is a sample generation pipeline which leverages an LLM to produce technical documentation excerpts, whose plausibility we have validated with legal experts (as described below in Section 2.2). This method, we argue, has standalone value, as it can be adapted to extend this benchmark or to generate new evaluation (and perhaps even training) datasets.

108 When conducting a compliance assessment, regulators, consultants, and internal audit teams may
109 draw upon many different kinds of input, including technical documentation, source code, and tran-
110 scriptions of staff interviews. However, in designing our benchmark, we optimized for simplicity and
111 ease of use by focusing on a single type of input. Specifically, since AI compliance assessors have
112 identified technical documentation as “the most important factor” in assessing whether an AI system
113 complies with the governing regulations (Li & Goel, 2025) — something that we validated in our
114 own interviews with AIR compliance experts (described in Appendix D) — we decided to make
115 technical documentation (or excerpts thereof) the sole type of input to inform compliance assess-
116 ments in our benchmark.

117 The creation of our sample generation pipeline was motivated by two key bottlenecks: first, little
118 to no AIA technical documentation of real AI systems is publicly available, perhaps due to the
119 confidentiality or legal privilege surrounding such assessments (Guha et al., 2024); and second,
120 paying experts to create them anew is prohibitively expensive (Pipitone & Alami, 2024). Therefore,
121 inspired in part by Sovrano et al. (2025), who use an LLM to assist human drafting of technical
122 documentation for AIA use cases, we designed a multi-stage pipeline, with gpt-4.1-mini (OpenAI,
123 2025) at its center, which generates plausible technical documentation excerpts efficiently and at
124 scale. The stages of this pipeline are described in Section 2.1 and illustrated in Fig. 2. [Although](#)
125 [we also experimented with o3-mini \(OpenAI, 2025b\), 4o-mini \(OpenAI, 2024\), and gpt-5 \(OpenAI,](#)
126 [2025a\), gpt-4.1-mini was ultimately chosen for the pipeline due to its price point and the satisfactory](#)
127 [outputs produced during our experimentation phase.](#)

128 In devising this pipeline, we set forth several design criteria for the samples it generates. Many of
129 these criteria, which are supplied in full in Appendix C, were written to ensure the benchmark re-
130 mains manageably-scoped and easy to use. For example, we decided that our documentation should
131 only depict HRAI systems within the scope of the AIA (EU, 2024, Art. 6). [Although we could](#)
132 [have chosen any subset of AIA requirements to treat as a proof of concept for this benchmarking](#)
133 [effort, we felt the AIA’s HRAI requirements were an especially worthwhile subject because many](#)
134 [regard them as “the most important part of the AIA” \(Araszkievicz et al., 2022\).](#) We also decided
135 that each sample should be written only from the perspective of an AI system provider attempting
136 to demonstrate compliance with a single article within the AIA (EU, 2024, Art. 2).

137 Additionally, we labored to ensure our excerpts were realistic (i.e., representative of real-world
138 technical documentation) and diverse (covering a range of AI systems with different intended uses
139 and varying levels of compliance). As a concrete example, to give researchers control over the
140 makeup and diversity of the distribution, we designed our pipeline to be steerable, allowing for the
141 targeted generation of excerpts that are more or less likely to be compliant. The details of how this
142 control was applied in AIReg-Bench are provided in Section 3, and further discussion of our design
143 criteria and the rationale for these criteria can be found in Appendix C.

144 2.1 STAGES OF SAMPLE GENERATION

145 The design criteria described above were enforced through prompt engineering during each stage of
146 technical documentation excerpt generation, as outlined below:

- 147 1. First, gpt-4.1-mini is prompted to generate high-level overviews of AI systems, which fall
148 into several use cases, such as road traffic control and credit scoring. By design, these use
149 cases should be classified as high-risk under the AIA (EU, 2024, Art. 6(2); Ann. III).
- 150 2. For each of these use cases, gpt-4.1-mini is given the system overview and a single AIA
151 article (either Art. 9, 10, 12, 14, or 15) as context and then prompted to generate ‘compli-
152 ance profiles’ (i.e., instructions for whether and in what way the AI system should breach
153 a selected article) for each overview-article combination. Within these profiles is a short
154 summary of how a selected AI system could breach a selected article.
- 155 3. For each of these compliance profiles, gpt-4.1-mini is prompted to generate an excerpt of
156 technical documentation, using the relevant article and AI system overview as context.
157

158 The prompts used in each stage are included in Appendix E.
159
160
161

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

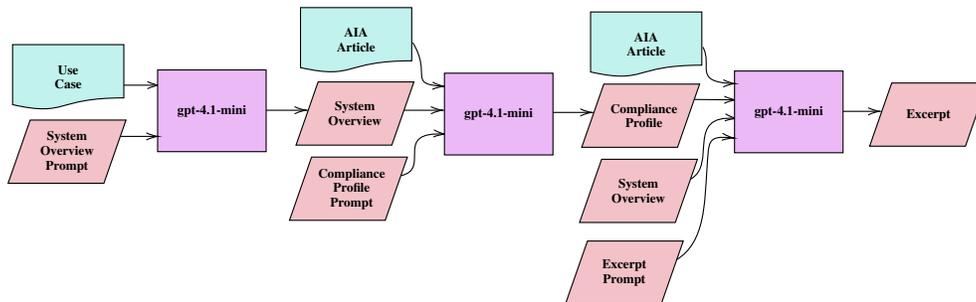


Figure 2: Illustration of the AIReg-Bench Technical Documentation Excerpt Generation Pipeline.

2.2 VALIDATION OF PLAUSIBILITY

To validate the plausibility of the pipeline’s outputs, each excerpt in the AIReg-Bench dataset was reviewed by three legal experts (the same team of law school students, law graduates, and qualified lawyers, [all regulation specialists](#), who supplied the annotations for the dataset, described in Section 3). As described in the annotation instructions (Appendix F), these legal experts were asked to label each excerpt with the probability that it is plausible — i.e., that it is realistic, logically consistent, and reflects the type of technical documentation that a Fortune 500 Europe AI provider might realistically hand over to its compliance assessor or internal audit team. These labels help assess whether the process of using an LLM to generate technical documentation succeeded.

To be more specific, for each excerpt, annotators were asked to provide a score on a 1–5 Likert scale (Likert, 1932), where 1 indicated a very low probability of plausibility and 5 indicated a very high probability of plausibility (the exact phrasing of the Likert scale given in Appendix F). To accompany these quantitative scores, annotators were asked to provide a qualitative free text justification of the Likert score (i.e., a free text plausibility analysis). These scores and text entries are available in the AIReg-Bench GitHub repo.

The median plausibility score provided by annotators was 4 (i.e., high probability of plausibility).

	Traffic Safety (Use1)	Gas Delivery (Use2)	Education (Use3)	Exam Proctor. (Use4)	Job Hiring (Use5)	Job Term. (Use6)	Emergency (Use7)	Credit Scoring (Use8)	Total
Likert score 1	1	2	0	0	1	0	0	0	4
Likert score 2	6	5	4	1	6	1	3	4	30
Likert score 3	9	5	7	3	9	6	7	4	50
Likert score 4	11	17	11	16	16	19	15	17	122
Likert score 5	18	16	23	25	13	19	20	20	154
All scores	45	45	45	45	45	45	45	45	360

Table 1: AIReg-Bench dataset overview of human annotations for plausibility on a Likert scale.

3 ANNOTATIONS AND DATASET (SOLUTION PART II)

Our second contribution involves using the sample generation pipeline that we created to generate a balanced distribution of samples, which legal experts then annotated in order to create the AIReg-

Bench open benchmark dataset. This dataset is composed of 120 technical documentation excerpts that reflect 8 different use cases (intended uses) for AI systems and varying compliance profiles. A list of the use cases from the creation of AIReg-Bench is included as Appendix H and the prompt used to generate compliance profiles is included as Appendix E.2.

More specifically, to create a diverse distribution of samples, with some shaped to have more compliant properties and others less so, we programmed the generation pipeline to steer one third of samples towards compliance and the remainder towards non-compliance. However, ultimately, the 360 Likert scale compliance labels (3 per excerpt) provided by our annotators reflect the true diversity of our dataset, as shown in Table 2.

	Traffic Safety (Use1)	Gas Delivery (Use2)	Education (Use3)	Exam Proctor. (Use4)	Job Hiring (Use5)	Job Term. (Use6)	Emergency (Use7)	Credit Scoring (Use8)	Total
Likert score 1	6	10	8	7	5	8	3	9	56
Likert score 2	11	13	12	11	10	9	12	7	85
Likert score 3	10	5	5	9	11	6	11	8	65
Likert score 4	4	3	7	4	4	8	5	8	43
Likert score 5	14	14	13	14	15	14	14	13	111
All scores	45	45	45	45	45	45	45	45	360

Table 2: AIReg-Bench dataset overview of human annotations for compliance on a Likert scale.

3.1 LEGAL EXPERT ANNOTATION OF THE EXCERPTS

Generating annotations in the legal domain often demands specialized legal expertise reflecting deep subject-matter knowledge (Guha et al., 2024). Accordingly, multiple legal natural language processing and LLM benchmarks have leveraged legal expert annotators (including law school students and lawyers) (Wang et al., 2025b; Zheng et al., 2025; Östling et al., 2024; Wang et al., 2023; Shen et al., 2022; Hendrycks et al., 2021b; Leivaditi et al., 2020; Zhong et al., 2019; Duan et al., 2019; Wilson et al., 2016). Following this pattern, we used a team of six legal experts (law graduates, law students, or qualified lawyers) to review and annotate the technical documentation excerpts in AIReg-Bench.¹ As detailed in Appendix I, each of these annotators listed regulation as one of their specializations. Prior to annotation, these legal experts attended a training session, led by a legally-trained co-author, dedicated to the AIA and its relevant articles. During annotation, labels were quality-checked by a law school graduate co-author.

In the end, each excerpt was reviewed by three legal experts. As described in the annotation instructions (Appendix F), for each excerpt, these annotators were asked to provide a score on a 1–5 Likert scale, where 1 indicated a very low probability of compliance with the relevant AIA article and 5 indicated a very high probability of compliance with that article (the exact phrasing of the Likert scale given in Appendix F).² To accompany these quantitative scores, annotators were asked to provide a qualitative free text justification of the Likert score (e.g., a free text compliance analysis). These

¹There are very few potential annotators with expertise in the EU AI Act who possess both the willingness and the capacity to carry out extended annotation tasks. For this reason, we broadened the eligibility criteria to include annotators with legal training more generally— although, ultimately, our entire team listed regulation as one of their specializations (see Appendix I).

²Annotators also mark whenever they found assigning a Likert score “difficult”. These markings can theoretically be used to segregate the more challenging samples. Notably, however, few annotators choose to use this demarcation in practice.

justifications can, in theory, be semantically compared to the justifications produced by LLMs to provide another way to measure LLM compliance assessment capabilities.

Appendix J contains a select annotation example. Appendix K contains an analysis of recurring themes in these annotations. The annotators’ inter-rater reliability, as measured with Krippendorff’s alpha coefficient (Krippendorff, 2018), was 0.651. This suggests a moderate level of agreement, though too low for drawing strong conclusions, likely reflecting the subjectivity of legal judgements.

Variance is heavily structured by two annotators with opposing biases (-0.917 and +0.600). Removing these annotators increases Krippendorff Alpha by +0.1343 to 0.786. That said, for the analyses in this paper, we retained all annotations to avoid post hoc exclusion.

The highest average disagreement occurred around Article 10 and Article 15, with mean standard deviations of 0.638 and 0.612 respectively. The most severe disagreement involved the intended use of Credit Scoring (Use 8) in light of Article 9, whose scores (1, 5, 5) resulted in the dataset’s highest standard deviation (1.89).

4 EXPERIMENTS (PROVING OUR SOLUTIONS ADDRESS THE PROBLEM)

Our third contribution is the first application of the benchmark. We have created an evaluation of 10 frontier language models using the AIReg-Bench benchmark. This helps us understand, the current performance of frontier LLMs, out-of-the-box and without fine-tuning, at the AIA compliance assessment task.

In this evaluation, we prompted the LLMs to carry out the same task as the human annotators, supplying them with the identical documentation, system descriptions, and AIA articles that had been provided to the human annotators, along with instructions that were highly similar to those given to the human annotators (detailed in Appendix G). A key distinction, however, is that the annotators were free to consult external sources—such as websites or existing literature—whereas the LLMs were restricted to the materials explicitly supplied.

Each LLM generated annotations for all 120 excerpts, using the same format as the human expert annotators: Likert scale scores for compliance, accompanied by textual justifications for those scores. The LLM compliance scores were then compared with the median scores assigned by the human annotators, allowing us to evaluate how closely each model approximated human compliance judgments. The key statistics from this evaluation are included in Table 3.

Model	κ_w (\uparrow)	ρ (\uparrow)	Bias ($\rightarrow 0$)	MAE (\downarrow)
GPT-5 (OpenAI, 2025a)	0.849	0.838	-0.067	0.450
GPT-4o (OpenAI et al., 2024)	0.775	0.842	0.458	0.558
o3 (OpenAI, 2025c)	0.723	0.809	-0.192	0.658
o3 mini (OpenAI, 2025b)	0.624	0.798	0.742	0.775
Claude Sonnet 4 (Anthropic, 2025)	0.772	0.779	-0.150	0.600
Gemini 2.5 Pro (Comanici et al., 2025)	<u>0.863</u>	<u>0.856</u>	-0.225	<u>0.458</u>
Gemini 2.5 Flash (Comanici et al., 2025)	0.729	0.825	-0.108	0.625
Gemma 3 (Kamath et al., 2025)	0.696	0.757	0.258	0.692
Grok 4 (xAI, 2025b)	0.829	0.829	0.242	0.475
Grok 3 mini (xAI, 2025a)	0.730	0.810	0.492	0.592

Table 3: **Agreement between LLMs and humans.** Columns report agreement between LLM and median human compliance scores across AIReg-Bench: quadratically weighted Cohen’s κ_w ; Spearman’s ρ ; Bias (mean signed difference, LLM–human); and MAE (mean absolute error).

All evaluated models demonstrated at least modest alignment with human expert judgments, with o3 mini showing the weakest Cohen’s Kappa agreement (0.624). At the other end of the spectrum, Gemini 2.5 Pro achieved the highest level of agreement (0.863), as well as the best rank correlation (0.856) and mean absolute error (0.458). In fact, Gemini 2.5 Pro’s compliance scores were within one point of the median human expert score for all but 7 out of 120 human expert median annotations (see Figure 3).

324 Despite prompts designed to mitigate sycophancy and acquiescence bias (Fanous et al., 2025), some
325 models tended to assign higher compliance scores than human experts. Along this dimension, o3
326 mini and GPT-4o performed worst, with o3 mini strictly exceeding the median human expert score
327 in 54.2% of excerpts, while only strictly falling below the median human expert’s score in 1.7% of
328 excerpts (see Table 4).

329 Ablations on GPT-4o revealed that modifying the prompt to request “harsh” and “critical” responses
330 can reduce bias, but at the cost of declines across all three other metrics (see Table 5).

331 332 333 5 BACKGROUND AND RELATED WORK 334

335 In this section, we provide context for our work by reviewing some foundational concepts as well as
336 the body of prior research has explored related methods and applications.

337 338 5.1 LEGAL AND OTHER LLM BENCHMARKS 339

340 Benchmark datasets that let researchers quantitatively measure how well an LLM performs a task
341 have become an important factor in developing trust in these models (Guha et al., 2024). Although
342 some popular benchmarks broadly assess LLM capabilities (Hendrycks et al., 2021a; Rajpurkar
343 et al., 2016), where models are evaluated on specific tasks, it is desirable for benchmarks to be
344 tailored more closely to those tasks (Peng et al., 2024). In this regard, a growing number of bench-
345 marks have been developed to assess the performance of LLMs at legal tasks such as contract review
346 (Hendrycks et al., 2021b; Wang et al., 2023), legal reading comprehension (Duan et al., 2019), and
347 more (Zheng et al., 2025; Leivaditi et al., 2020). Guha et al. (2023) and Fei et al. (2023) both gather
348 these prior benchmarks as well as other resources into aggregate LLM legal benchmarks.

349 350 5.2 LLMs FOR LEGAL, COMPLIANCE, AND AIR COMPLIANCE TASKS 351

352 Researchers have applied LLMs to a wide variety of legal tasks (Ma et al., 2024; Lai et al., 2024;
353 Siino et al., 2025). This includes various compliance assessment tasks (Hassani, 2024; Bolton et al.,
354 2025; Chen et al., 2024; Wang et al., 2025a), including AIR compliance assessments. For example,
355 Makovec et al. (2024) input datasets, model cards, README files, or other AI project artifacts
356 into a RAG-enhanced GPT-4 that accesses relevant portions of the AIA to predict the compliance
357 level of the AI system depicted in the input. Davvetas et al. (2025) use a RAG-equipped LLM
358 (mistral-small3.2) that takes, as input, certain features of an AI system (such as the type of AI
359 system and the intended use) and outputs the risk-level of the AI system according to the AIA.
360 Similarly, Kővári et al. (2025) use in-context learning and RAG to create a chatbot that can help
361 users self-assess compliance with the AIA. Meanwhile, Li et al. (2025) test various LLMs’ ability
362 to act as a rudimentary AIA compliance checker by accepting a hypothetical AI system as context and
363 predicting whether it is prohibited by, permitted by, or out of scope of the AIA. A series of interlinked
364 studies by Nokia Bell Labs employ LLMs to support AI practitioners in AIR compliance subtasks
365 such as populating impact assessment reports (Bogucka et al., 2024a; Herdel et al., 2024; Bogucka
366 et al., 2024b). Sovrano et al. (2025) use an LLM to help human-drafted technical documentation
367 align with the requirements of AIA Article 11.

368 369 5.3 THE EU AI ACT 370

371 The AIA went into force in August 2024 (Lomas, 2024) and sets forth harmonized requirements
372 for AI systems and models placed on the market or put into service in the EU (EU, 2024, Art.
373 1-2). In laying out requirements for these AI systems, the AIA leverages a “risk-based” approach
374 (Mahler, 2022), by which the exact requirements that apply to a system are a function of its perceived
375 degree of risk. Here, the most demanding requirements are reserved for those AI systems deemed
376 to be high-risk (EU, 2024, Art. 6). Such high-risk AI systems (HRAI) must satisfy a number of
377 requirements (EU, 2024, Chap. III, Sec. 2). Among those, the requirements that we have made the
focus of our benchmark dataset relate to risk management systems, data and data governance, record
keeping, human oversight, as well as accuracy, robustness, and cybersecurity (EU, 2024, Art. 9, 10,
12, 14, 15).

5.4 NON-LLM ALGORITHMS THAT ASSESS AIR COMPLIANCE

It is important to distinguish the works in Section 5.2 as well as this work from those works that use algorithms other than LLMs, including but not limited to benchmark suites, to evaluate whether AI systems, including LLM-based systems, to comply with AI regulations (Prandi et al., 2025; Marino et al., 2024; Guldemann et al., 2024; Walke et al., 2023). While these works potentially present interesting accompaniments to the LLM-based approaches that this work aims to benchmark, this work does not seek to create a benchmark for these other methods.

6 DISCUSSION

Here, we consider some of the limitations of our work, anticipate some of the questions that the research community might reasonably have about our methods, and describe how we addressed those.

6.1 CHALLENGES OF LEGAL BENCHMARKING

Compliance assessments are subjective. Complicating legal benchmarking is the fact that legal tasks often involve subjective judgments (Guha et al., 2024; Ma et al., 2023). Compliance, in particular, has been described as “not binary” (Wu & van Rooij, 2021). This was one motivation for our use of Likert scale annotations keyed to the “probability” of compliance rather than binary labels of “compliant” or “non-compliant.” The Likert scale is viewed as a reliable way to “transform subjective qualitative data into quantifiable metrics” (Koo & Yang, 2025) and has previously been used for benchmarking of LLMs in subjective realms (Bojić et al., 2025).

Potentially increasing the subjectivity of legal annotations is the nascency of the AIA, which lacks the established guidelines and court rulings that typically help annotators reach more consistent conclusions (Goodman, 2023). To quantify this subjectivity, we measured annotators’ inter-rater reliability via the Krippendorff’s alpha coefficient (Krippendorff, 2018), which came to 0.651. This reflects moderate agreement between annotators, albeit lower than the levels of inter-rater reliability typically expected in domains with less subjective tasks.

To mitigate the variance that arises from the subjectivity of compliance assessments, each excerpt was scored independently by three of our six annotators, and most analyses were conducted using the median of these scores. Additionally, to pinpoint areas of greater subjectivity, we asked annotators to flag compliance annotations that were more “difficult.” However, in practice, few annotations were flagged as challenging, potentially as annotators struggled to identify which cases were more challenging (Rother et al., 2021).

Compliance assessments are a moving target. It has been said that compliance assessments require clear and specific guidelines, including relevant case law (Kilian et al., 2025; Schuett, 2024). But, in its present state, the AIA arguably lacks these. The text of the law has not been interpreted by courts (Yew et al., 2025). The obligations outlined in the AIA have yet to be clarified by accompanying technical standards (European Committee for Electrotechnical Standardization, 2024; European Commission, 2022). They are also subject to ongoing amendments (EU, 2024, Art. 96). Accordingly, AIReg-Bench should be viewed as a snapshot of AIA conformity assessments in September 2025. It does not and cannot reflect developments occurring after this date.

6.2 CHALLENGES OF LLM-DRIVEN SAMPLE GENERATION

Benchmark dataset samples should be representative of real-world data (Sourlos et al., 2024). Broadly speaking, there is evidence that LLMs (especially larger ones) can effectively produce synthetic samples satisfying this criteria (Maheshwari et al., 2024). More directly relevant here, it has been shown that LLMs can effectively generate (or improve) legal documents (Su et al., 2025; Lin & Cheng, 2024; Hemrajani, 2025; Gray et al., 2025), technical specifications (Xie et al., 2025), compliance documentation (Wang et al., 2025d; Hassani, 2024; Kumar & Roussinov, 2024), and even the type of technical documentation required under the AIA (Sovrano et al., 2025).

However, some are sceptical about the ability of LLMs to generate realistic outputs in these domains (Posner & Saran, 2025; Roberts; Shen et al., 2022). Critics point out that LLMs often lack domain-

specific tacit knowledge, have difficulty maintaining coherent reasoning across extended contexts, and may hallucinate facts or references (Rasiah et al., 2024; Huang et al., 2024; Dahl et al., 2024; Magesh et al., 2025). Regarding legal tasks specifically, critics note how LLMs’ struggle to interpret legal terminology, grasp case context, and execute complex analyses, potentially resulting in errors (Wang et al., 2024; Roberts; Shen et al., 2022). With this in mind, a number of guardrails were put in place to enforce plausibility in our dataset samples and to create transparency around whether and where those guardrails fell short. For example:

- The sample generation method was informed by the series of interviews with actual compliance assessment experts, described in Appendix D. We also consulted similar interviews performed by Li & Goel (2025) and the recommended protocols for manual compliance assessments for the AIA (Floridi et al., 2022; Thelisson & Verma, 2024; Lillo Campoy et al., 2024; Palumbo et al., 2025) and other AIR (The Institute of Internal Auditors, 2023; National Institute of Standards and Technology, 2023; Brogle et al., 2025).
- The sample generation method was co-engineered by a law school graduate co-author who has been involved in the drafting of the codes of practice accompanying the AIA.
- Before the samples in AIReg-Bench were generated, the method was subjected to an iterative, plausibility-focused development process with a subset of our expert annotators (including an EU qualified lawyer). This iterative process significantly improved our prompts (as measured by the plausibility of the excerpts they generated) compared to those produced using small prompts, which we initially considered. As an example, unlike those in our final excerpts, the compliance violations generated by small prompts were exceedingly obvious, superficial, and unrealistic — a limitation also noted by Nguyen et al. (2025).
- During annotation, our legal expert annotators scored AIReg-Bench samples for plausibility, with those scores and their accompanying text explanations being made public, in their entirety, as part of the AIReg-Bench dataset.

6.3 CHALLENGES OF LLM-DRIVEN COMPLIANCE ANALYSES

Beyond the challenges of generating realistic documentation samples with LLMs, there may also be hurdles to using LLMs to perform legal analyses on text. Although some research efforts have found that LLMs match or exceed human accuracy when performing such analyses (Martin et al., 2024), others have questioned whether LLMs can perform this task effectively (Buckland, 2023; Doyle & Tucker, 2024; Li et al., 2024; Network, 2025; Mik, 2024). While the evidence presented in this paper is not definitive, we hope that our benchmark offers an initial step towards clarifying how well LLMs perform at compliance assessments, both in comparison to human experts and to one another. Our hope is that this will inspire other efforts to quantitatively evaluate the performance of LLMs at AIR compliance analyses and other legal tasks.

6.4 RISKS OF LLM-DRIVEN COMPLIANCE ANALYSES

As discussed in Section 1, there are potential benefits to LLM-driven AIR compliance analyses. However, it is important to point out that they could also carry risks. For example, given LLMs’ ongoing tendency to make errors, some argue that there may be dire consequences when lawyers place “blind faith in an LLM” (Moriarty, 2023) and that it may even represent a violation of professional ethics (Browning, 2024). [Conversely, the developers of AI systems could try to “game” these LLMs, manipulating their technical documentation so as to achieve desirable compliance assessment outcomes.](#) To help inform [the conversation about these risks](#)~~this conversation~~, and avoid generalization [within that conversation](#), we believe it is invaluable to quantitatively evaluate the performance of LLMs *for the task at hand*, comparing their performance to human performance at that task in an evidence-based manner. Hence this benchmark.

7 FUTURE WORK

This benchmark should serve as a starting point for tracking LLM progress at AIR compliance assessments, rather than as a finish line indicating readiness for deployment in legal practice. Mov-

486 ing towards that goal will require additional benchmarks that build upon and extend AIReg-Bench.
487 Some suggested directions for extensions to AIReg-Bench are outlined below.

488
489 **Extension to other AIR.** AIReg-Bench is currently scoped to a subset of the AIA’s requirements
490 for HRAI systems. In the future, it would be natural to extend AIReg-Bench to the rest of the AIA’s
491 requirements for HRAI systems as well as to its requirements for general purpose AI models (Euro-
492 pean Parliament, 2024, Chap. V). When other AIRs achieve the AIA level of maturity, AIReg-Bench
493 could also be extended to cover those AIR as well. In all cases, we believe that AIReg-Bench’s
494 overall playbook could be re-used, though the excerpt generation pipeline would need to be subtly
495 reconfigured for these regulations and the excerpts annotated in light of the different compliance
496 requirements.

497
498 **Extension to further LLM-powered annotators.** It would be valuable to extend our benchmarking
499 in Section 4 to include more models: including fine-tuned legal LLMs, which some argue perform
500 better at legal tasks (Fei et al., 2023; Dominguez-Olmedo et al., 2024), as well as LLMs with tool-
501 use (e.g., incorporating RAG or web search) (Makovec et al., 2024; Davvetas et al., 2025; Wang &
502 Yuan, 2025). Though, while to our knowledge, no such models exist yet, there would also be value
503 in fine-tuning LLMs for AIA or AIR compliance and then evaluating them using AIReg-Bench.

504
505 **Extension to real-world documentation.** AIReg-Bench consists of LLM-generated excerpts of
506 technical documentation for the AIA. We relied on these LLM-generated excerpts, whose plausi-
507 bility we manually verified, since public access to [authentic real](#) technical documentation [from real](#)
508 [AI developers](#) is limited, perhaps due to the confidentiality, [legal privilege](#), or the relative nascency
509 of the AIA. That said, we recognize the value, specifically as it relates to construct validity, of a
510 benchmark built from real instances of AIA technical documentation. We therefore encourage any
511 actors with access to such documentation to consider publishing (an anonymized version of) [it and,](#)
512 [furthermore, to contact us if we can assist with that process.](#)

513
514 **Extension to AIA technical standards.** Like other EU product regulation, the AIA will utilize
515 harmonized standards, i.e., “more concrete” (Siegmann & Anderljung, 2022) technical specifica-
516 tions prepared by the EU’s external standardization organizations (CEN, CENELEC and ETSI).
517 Compliance with these specifications, which are still under development (European Committee for
518 Electrotechnical Standardization, 2024; European Commission, 2022), will “have the legal effect of
519 establishing a presumption of conformity” with the AIA (Mazzini & Scalzo, 2023). Once issued,
520 it would be important to pass these specifications, as additional context, into any AIA compliance
521 assessment algorithm.

522
523 **Extension to multi-turn interactions.** Our interviews of AIR compliance professionals suggested
524 that, in practice, compliance assessments often involve a long and complex dialogue between legal
525 teams, technical staff, and regulators (see Appendix D). AIReg-Bench condenses the entire com-
526 pliance assessment process into a single-turn interaction, based on a fixed set of synthetic artifacts.
527 Such scoping and simplifications are common in benchmarking, though many scholar stress the need
528 to shift towards more interactive modes of evaluating AI system capabilities (Ibrahim et al., 2024;
529 Eriksson et al., 2025). Future benchmarks could build on AIReg-Bench by evaluating LLMs’ multi-
530 turn ability to collaborate with human teams and contribute to complex legal dialogues (Kővári et al.,
531 2025), rather than merely producing one-shot assessments. [Similarly, future benchmarks might ex-](#)
532 [plore scenarios where LLMs are merely used as the starting point in a multi-turn, human-in-the-loop](#)
533 [compliance assessment process.](#)

534 8 CONCLUSION

535
536 In this work, we introduced AIReg-Bench, an open benchmark designed to quantitatively evaluate
537 the performance of LLMs on AIA compliance assessments. By combining an LLM-driven sam-
538 ple generation pipeline with expert legal annotations, AIReg-Bench provides a scalable, realistic,
539 and extensible foundation for assessing how closely models align with human expert compliance
540 judgments. Our initial experiments with frontier LLMs demonstrate both the promise and current
541 limitations of these systems in performing this task. While AIReg-Bench is only an initial step, we
542 hope it catalyzes further research into LLM-driven AIR compliance assessments.

9 ETHICS STATEMENT

We do not believe that our paper submission raises questions regarding the Code of Ethics. All legal expert annotators are co-authors on the work. What is more, since this work ultimately strives to benchmark and improve LLM-driven compliance assessments, it should ultimately serve to increase compliance with the AI regulations that tend to encode important societal values such as fairness and safety.

10 REPRODUCIBILITY STATEMENT

Here, we discuss the efforts that have been made to ensure reproducibility. The parts of the main paper, appendix, and supplemental materials (including the project GitHub repository) that will help with reproducibility are as follows:

Sample generation pipeline: The full code for the Sample generation pipeline is available at the AIReg-Bench GitHub repository <https://anonymous.4open.science/r/aireg-bench-5259/>. In addition, the following are included here as appendices:

- Sample generation pipeline prompts (Appendix E)
- Use cases for the AIReg-Bench sample generation pipeline (Appendix H)

AIReg-Bench dataset: The full AIReg-Bench dataset, including samples and annotations, is available at the AIReg-Bench GitHub repository <https://anonymous.4open.science/r/aireg-bench-5259/>. In addition, the following are included here as appendices:

- Legal expert annotation instructions (Appendix F)

Evaluation of frontier LLMs using AIReg-Bench: The full code for the evaluation is available at the AIReg-Bench GitHub repository <https://anonymous.4open.science/r/aireg-bench-5259/>

. In addition, the following are included as appendices:

- LLM annotation instructions (Appendix G)

REFERENCES

- Ada Lovelace Institute. Code & conduct: How to create third-party auditing regimes for AI. Discussion paper, Ada Lovelace Institute, June 2024. URL <https://www.adalovelaceinstitute.org/wp-content/uploads/2024/06/Ada-Lovelace-Institute-Code-and-conduct-FINAL-1906.pdf>.
- Sacha Alanoca, Shira Gur-Arieh, Tom Zick, and Kevin Klyman. Comparing apples to oranges: A taxonomy for navigating the global landscape of AI regulation. In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency, FAccT '25*, pp. 914–937. ACM, June 2025. doi: 10.1145/3715275.3732059. URL <http://dx.doi.org/10.1145/3715275.3732059>.
- Markus Anderljung, Joslyn Barnhart, Anton Korinek, Jade Leung, Cullen O’Keefe, Jess Whittlestone, Shahar Avin, Miles Brundage, Justin Bullock, Duncan Cass-Beggs, Ben Chang, Tantum Collins, Tim Fist, Gillian Hadfield, Alan Hayes, Lewis Ho, Sara Hooker, Eric Horvitz, Noam Kolt, Jonas Schuett, Yonadav Shavit, Divya Siddarth, Robert Trager, and Kevin Wolf. Frontier AI regulation: Managing emerging risks to public safety, 2023. URL <https://arxiv.org/abs/2307.03718>.
- Anthropic. Claude 4 system card: Claude Opus 4 and Claude Sonnet 4. Technical report, Anthropic, May 2025. URL <https://www-cdn.anthropic.com/4263b940cabb546aa0e3283f35b686f4f3b2ff47.pdf>. Accessed: 2025-09-21.
- Michał Araszkiwicz, Grzegorz J. Nalepa, and Radosław Pałosz. The artificial intelligence act: A jurisprudential perspective. In Desara Dushi, Francesca Naretto, Cecilia Panigutti, and Francesca

- 594 Pratesi (eds.), *Proceedings of the Workshop on Imagining the AI Landscape after the AI Act*
595 (*IAIL 2022*), volume Vol-3221 of *CEUR Workshop Proceedings*, 2022. URL [https://ceur-
597 ws.org/Vol-3221/IAIL_paper8.pdf](https://ceur-
596 ws.org/Vol-3221/IAIL_paper8.pdf). Accessed: YYYY-MM-DD.
- 598 Edyta Bogucka, Marios Constantinides, Sanja Šćepanović, and Daniele Quercia. AI design: A
599 responsible AI framework for impact assessment reports. *IEEE Internet Computing*, 2024a.
- 600 Edyta Bogucka, Marios Constantinides, Sanja Šćepanović, and Daniele Quercia. Co-designing
601 an AI impact assessment report template with AI practitioners and AI compliance experts. In
602 *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 168–180,
603 2024b.
- 604
605 Luka Bojić, Olga Zagovora, Asta Zelenkauskaitė, et al. Comparing large language models and
606 human annotators in latent content analysis of sentiment, political leaning, emotional intensity
607 and sarcasm. *Scientific Reports*, 15:11477, 2025. doi: 10.1038/s41598-025-96508-3. URL
608 <https://doi.org/10.1038/s41598-025-96508-3>.
- 609
610 Regan Bolton, Mohammadreza Sheikhhathollahi, Simon Parkinson, Vanessa Vulovic, Gary Bam-
611 ford, Dan Basher, and Howard Parkinson. Document retrieval augmented fine-tuning (draft) for
612 safety-critical software assessments, 2025. URL <https://arxiv.org/abs/2505.01307>.
- 613
614 K. Brogle, E. Kallina, H. Sargeant, et al. Context-specific certification of AI systems: a pilot in
615 the financial industry. *AI Ethics*, 2025. doi: 10.1007/s43681-025-00720-w. URL <https://doi.org/10.1007/s43681-025-00720-w>.
- 616
617 John G. Browning. Robot lawyers don’t have disciplinary hearings—real lawyers do: The eth-
618 ical risks and responses in using generative artificial intelligence. *Georgia State University*
619 *Law Review*, 40(4):917–962, 2024. URL [https://readingroom.law.gsu.edu/gsulr/
620 vol40/iss4/9](https://readingroom.law.gsu.edu/gsulr/vol40/iss4/9).
- 621
622 Robert Buckland. AI, judges and judgment: Setting the scene. Associate Working Paper Series
623 220, Harvard Kennedy School, Mossavar-Rahmani Center for Business and Government, Cam-
624 bridge, MA, November 2023. URL [https://dash.harvard.edu/server/api/core/
625 bitstreams/98187fff-8a7a-4ca6-8123-3049e417f088/content](https://dash.harvard.edu/server/api/core/bitstreams/98187fff-8a7a-4ca6-8123-3049e417f088/content).
- 626
627 Nitay Calderon, Roi Reichart, and Rotem Dror. The alternative annotator test for LLM-as-a-judge:
628 How to statistically justify replacing human annotators with LLMs. In Wanxiang Che, Joyce
629 Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd*
630 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.
631 16051–16081, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-
632 8-89176-251-0. doi: 10.18653/v1/2025.acl-long.782. URL [https://aclanthology.org/
2025.acl-long.782/](https://aclanthology.org/2025.acl-long.782/).
- 633
634 Nanjiang Chen, Xuhui Lin, Hai Jiang, and Yi An. Automated building information model-
635 ing compliance check through a large language model combined with deep learning and on-
636 tology. *Buildings*, 14(7), 2024. ISSN 2075-5309. doi: 10.3390/buildings14071983. URL
637 <https://www.mdpi.com/2075-5309/14/7/1983>.
- 638
639 Jon Chun, Christian Schroeder de Witt, and Katherine Elkins. Comparative global AI regulation:
640 Policy perspectives from the EU, China, and the US, 2024. URL [https://arxiv.org/abs/
2410.21279](https://arxiv.org/abs/2410.21279).
- 641
642 Pierre Colombo, Telmo Pires, Malik Boudiaf, Rui Melo, Dominic Culver, Etienne Malaboëuf,
643 Gabriel Hautreux, Johanne Charpentier, and Michael Desa. Saullm-54b & saullm-141b: Scaling
644 up domain adaptation for the legal domain. *Advances in Neural Information Processing Systems*,
645 37:129672–129695, 2024a.
- 646
647 Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, An-
d्रे FT Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, et al. Saullm-7b: A
pioneering large language model for law. *arXiv preprint arXiv:2403.03883*, 2024b.

- 648 Gheorghe Comanici et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality,
649 long context, and next generation agentic capabilities, 2025. URL [https://arxiv.org/
650 abs/2507.06261](https://arxiv.org/abs/2507.06261).
- 651
652 Sasha Costanza-Chock, Inioluwa Deborah Raji, and Joy Buolamwini. Who audits the auditors? rec-
653 ommendations from a field scan of the algorithmic auditing ecosystem. In *Proceedings of the 2022
654 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22*, pp. 1571–1583,
655 New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393522. doi:
656 10.1145/3531146.3533213. URL <https://doi.org/10.1145/3531146.3533213>.
- 657 Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. Large legal fictions: Profiling legal
658 hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93, 2024.
- 659 Athanasios Davvetas, Xenia Ziouvelou, Ypatia Dami, Alexis Kaponis, Konstantina Giou-
660 vanopoulou, and Michael Papademas. TAI scan tool: A rag-based tool with minimalistic input for
661 trustworthy AI self-assessment, Jul 2025. URL <https://arxiv.org/abs/2507.17514>.
- 662 Ricardo Dominguez-Olmedo, Vedant Nanda, Rediet Abebe, Stefan Bechtold, Christoph Engel, Jens
663 Frankenreiter, Krishna Gummadi, Moritz Hardt, and Michael Livermore. Lawma: The power of
664 specialization for legal tasks, 2024. URL <https://arxiv.org/abs/2407.16615>.
- 665
666 Colin Doyle and Aaron Tucker. If you give an LLM a legal practice guide. Technical Report SSRN
667 5030676, SSRN, Los Angeles, CA; Ithaca, NY, December 5 2024. URL [https://ssrn.com/
668 abstract=5030676](https://ssrn.com/abstract=5030676). Last revised: February 20, 2025.
- 669 Xingyi Duan, Baoxin Wang, Ziyue Wang, Wentao Ma, Yiming Cui, Dayong Wu, Shijin Wang,
670 Ting Liu, Tianxiang Huo, Zhen Hu, Heng Wang, and Zhiyuan Liu. *CJRC: A Reliable Human-
671 Annotated Benchmark DataSet for Chinese Judicial Reading Comprehension*, pp. 439–451.
672 Springer International Publishing, 2019. ISBN 9783030323813. doi: 10.1007/978-3-030-32381-
673 3_36. URL http://dx.doi.org/10.1007/978-3-030-32381-3_36.
- 674
675 Maria Eriksson, Erasmo Purificato, Arman Noroozian, Joao Vinagre, Guillaume Chaslot, Emilia
676 Gomez, and David Fernandez-Llorca. Can we trust AI benchmarks? an interdisciplinary review
677 of current issues in AI evaluation. *arXiv preprint arXiv:2502.06559*, 2025.
- 678 EU. Artificial Intelligence Act, March 2024. URL [https://eur-lex.europa.eu/legal-
679 content/EN/TXT/?uri=CELEX:52021PC0206](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206). Official Journal of the European Union.
- 680 European Commission. Commission staff working document impact assessment accompanying
681 the proposal for a regulation of the European Parliament and of the Council laying down har-
682 monised rules on artificial intelligence (Artificial Intelligence Act) and amending certain union
683 legislative acts. [https://ec.europa.eu/newsroom/dae/redirection/document/
684 75792](https://ec.europa.eu/newsroom/dae/redirection/document/75792), 2021. [Accessed 18-09-2025].
- 685
686 European Commission. Draft standardisation request to the European Standardisation Organisa-
687 tions in support of safe and trustworthy artificial intelligence. [https://ec.europa.eu/
688 docsroom/documents/52376](https://ec.europa.eu/docsroom/documents/52376), 2022. [Accessed 18-09-2025].
- 689 European Committee for Electrotechnical Standardization. Artificial intelligence. [https:
690 //www.cenelec.eu/areas-of-work/cen-cenelec-topics/artificial-
691 intelligence/](https://www.cenelec.eu/areas-of-work/cen-cenelec-topics/artificial-intelligence/), 2024. [Accessed 18-09-2025].
- 692
693 European Parliament. EU AI Act: First regulation on artificial intelligence. *European Parliament
694 Topics*, June 2024. URL [https://www.europarl.europa.eu/topics/en/article/
695 20230601STO93804/eu-ai-act-first-regulation-on-artificial-
696 intelligence](https://www.europarl.europa.eu/topics/en/article/20230601STO93804/eu-ai-act-first-regulation-on-artificial-intelligence). Last updated.
- 697
698 Aaron Fanous, Jacob Goldberg, Ank A. Agarwal, Joanna Lin, Anson Zhou, Roxana Daneshjou, and
699 Sanmi Koyejo. SycEval: Evaluating LLM sycophancy, 2025. URL [https://arxiv.org/
700 abs/2502.08177](https://arxiv.org/abs/2502.08177).
- 701
702 Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen,
Zongwen Shen, and Jidong Ge. Lawbench: Benchmarking legal knowledge of large language
models, 2023. URL <https://arxiv.org/abs/2309.16289>.

- 702 Luciano Floridi, Matthias Holweg, Mariarosaria Taddeo, Jakob Amaya Silva, Jakob Mökander, and
703 Yi Wen. capAI-A procedure for conducting conformity assessment of AI systems in line with
704 the EU Artificial Intelligence Act. *Available at SSRN*, 2022. URL [https://ssrn.com/
705 abstract=4064091](https://ssrn.com/abstract=4064091). SSRN 4064091.
- 706 Isaac Gazendam, Philip Dawson, and AI Armilla. Mind the gap: The challenges of assurance
707 for artificial intelligence. Technical report, Technical report, Stanford Center for International
708 Security and Cooperation, 2023.
- 709
710 Asress Adimi Gikay. Risks, innovation, and adaptability in the UK’s incrementalism versus the Eu-
711 ropean Union’s comprehensive artificial intelligence regulation. *International Journal of Law and
712 Information Technology*, 32(1):eaae013, 06 2024. ISSN 0967-0769. doi: 10.1093/ijlit/eaae013.
713 URL <https://doi.org/10.1093/ijlit/eaae013>.
- 714 Tristan Goodman. Thinking outside the technical standardisation box: The role of standards under
715 the draft eu artificial intelligence act. *LSE LR*, 9:73, 2023.
- 716
717 Morgan A. Gray, Li Zhang, and Kevin Ashley. Generating case-based legal arguments with LLMs.
718 In *Proceedings of the 2025 Symposium on Computer Science and Law (CSLAW ’25)*, pp. 160–
719 168. ACM, June 2025. doi: 10.1145/3709025.3712216. University of Pittsburgh Legal Studies
720 Research Paper No. 2025-22; also on SSRN.
- 721 Neel Guha, Christie M. Lawrence, Lindsey A. Gilmard, Kit T. Rodolfa, Faiz Surani, Rishi Bom-
722 masani, Inioluwa Deborah Raji, Mariano-Florentino Cuéllar, Colleen Honigsberg, Percy Liang,
723 and Daniel E. Ho. The AI regulatory alignment problem. [https://hai.stanford.edu/
724 sites/default/files/2023-11/AI-Regulatory-Alignment.pdf](https://hai.stanford.edu/sites/default/files/2023-11/AI-Regulatory-Alignment.pdf). [Accessed
725 18-09-2025].
- 726
727 Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex
728 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry
729 Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai
730 Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin
731 Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzen-
732 berger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams,
733 Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. Legalbench: A collaboratively built bench-
734 mark for measuring legal reasoning in large language models. In *Advances in Neural Information
735 Processing Systems 36 (NeurIPS 2023) Datasets and Benchmarks Track*, 2023.
- 736
737 Neel Guha, Julian Nyarko, Daniel E. Ho, and Christopher Ré. Building GenAI benchmarks: A
738 case study in legal applications. In Philipp Hacker, Andreas Engel, Sarah Hammer, and Brent
739 Mittelstadt (eds.), *The Oxford Handbook on the Foundations and Regulation of Generative AI*.
740 Oxford University Press, 2024. URL [https://neelguha.github.io/assets/pdf/
741 building_genai_benchmarks_for_law_oxford_chapter.pdf](https://neelguha.github.io/assets/pdf/building_genai_benchmarks_for_law_oxford_chapter.pdf). Forthcoming (draft
742 available online).
- 743
744 Philipp Guldemann, Alexander Spiridonov, Robin Staab, Nikola Jovanović, Mark Vero, Velko
745 Vechev, Anna Gueorgieva, Mislav Balunović, Nikola Konstantinov, Pavol Bielik, Petar Tsankov,
746 and Martin Vechev. COMPL-AI framework: A technical interpretation and LLM benchmark-
747 ing suite for the EU Artificial Intelligence Act, 2024. URL [https://arxiv.org/abs/
748 2410.07959](https://arxiv.org/abs/2410.07959).
- 749
750 Meeri Haataja and Joanna Bryson. What costs should we expect from the EU’s AI Act?, aug 2021.
751 Accessed: 2025-03-22.
- 752
753 Shabnam Hassani. Enhancing legal compliance and regulation analysis with large language models,
754 2024. URL <https://arxiv.org/abs/2404.17522>.
- 755
756 Rahul Hemrajani. Evaluating the role of large language models in legal practice in India, 2025. URL
757 <https://arxiv.org/abs/2508.09713>.
- 758
759 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-
760 cob Steinhardt. Measuring massive multitask language understanding, 2021a. URL <https://arxiv.org/abs/2009.03300>.

- 756 Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. CUAD: An expert-annotated NLP
757 dataset for legal contract review, 2021b. URL <https://arxiv.org/abs/2103.06268>.
758
- 759 Viviane Herdel, Sanja Šćepanović, Edyta Bogucka, and Daniele Quercia. ExploreGen: Large
760 language models for envisioning the uses and risks of AI technologies. In *Proceedings of the*
761 *AAAI/ACM Conference on AI, Ethics, and Society*, volume 7, pp. 584–596, 2024.
- 762 Chen Huang, Xinwei Yang, Yang Deng, Wenqiang Lei, JianCheng Lv, and Tat-Seng Chua.
763 Co-matching: Towards human-machine collaborative legal case matching. *arXiv preprint*
764 *arXiv:2405.10248*, 2024.
- 765 Lujain Ibrahim, Saffron Huang, Umang Bhatt, Lama Ahmad, and Markus Anderljung. Towards in-
766 teractive evaluations for interaction harms in human-ai systems. *arXiv preprint arXiv:2405.10632*,
767 2024.
- 768 Valeriia Iliasova, Emma Hovland, and Zeina Rustum. Kill Bill: Does the EU AI Act actually kill
769 innovation?, 2025. Student Paper.
- 770 Aishwarya Kamath et al. Gemma 3 technical report, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.19786)
771 [2503.19786](https://arxiv.org/abs/2503.19786).
- 772 Robert Kilian, Dominik Ebel, and Linda Jäck. Making digital regulation work – the crucial role technical
773 standards play in implementing the EU AI Act. *OECD AI Wonk Blog*, 2025. URL [https://oecd.ai/en/wonk/making-digital-regulation-work-the-crucial-
774 role-technical-standards-play-in-implementing-the-eu-ai-act](https://oecd.ai/en/wonk/making-digital-regulation-work-the-crucial-role-technical-standards-play-in-implementing-the-eu-ai-act).
- 775 Florence Koh, Kathrin Grosse, and Giovanni Apruzzese. Voices from the frontline: Revealing the
776 AI practitioners’ viewpoint on the European AI Act. In *Proceedings of the Hawaii International*
777 *Conference on System Sciences*, HICSS, 2024.
- 778 Florian Königstorfer and Stefan Thalmann. Ai documentation: A path to accountability. *Journal of*
779 *Responsible Technology*, 11:100043, 2022.
- 780 Malcolm Koo and Shih-Wei Yang. Likert-type scale. *Encyclopedia*, 5(1), 2025. ISSN 2673-8392.
781 doi: 10.3390/encyclopedia5010018. URL <https://www.mdpi.com/2673-8392/5/1/18>.
- 782 Klaus Krippendorff. *Content analysis: An introduction to its methodology*. Sage publications, 2018.
- 783 Bimal Kumar and Dmitri Roussinov. NLP-based regulatory compliance – using GPT 4.0 to decode
784 regulatory documents, 2024. URL <https://arxiv.org/abs/2412.20602>.
- 785 Ádám Kővári, Yasin Ghafourian, Csaba Hegedus, Belal Abu Naim, Kitti Mezei, Pál Varga, and
786 Markus Tauber. Let’s have a chat with the EU AI Act. In *NOMS 2025-2025 IEEE Network Oper-*
787 *ations and Management Symposium*, pp. 1–6, 2025. doi: 10.1109/NOMS57970.2025.11073655.
- 788 Jinqi Lai, Wensheng Gan, Jiayang Wu, Zhenlian Qi, and Philip S. Yu. Large language models
789 in law: A survey. *AI Open*, 5:181–196, 2024. ISSN 2666-6510. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/j.aiopen.2024.09.002)
790 [j.aiopen.2024.09.002](https://doi.org/10.1016/j.aiopen.2024.09.002). URL [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/S2666651024000172)
791 [pii/S2666651024000172](https://www.sciencedirect.com/science/article/pii/S2666651024000172).
- 792 Moritz Laurer, Andrea Renda, and Timothy Yeung. Clarifying the costs for the EU’s AI act. Centre
793 for European Policy Studies, September 2021. URL [https://www.ceps.eu/clarifying-
794 the-costs-for-the-eus-ai-act/](https://www.ceps.eu/clarifying-the-costs-for-the-eus-ai-act/).
- 795 Spyretta Leivaditi, Julien Rossi, and Evangelos Kanoulas. A benchmark for lease contract review,
796 2020. URL <https://arxiv.org/abs/2010.10386>.
- 797 Haitao Li, Junjie Chen, Jingli Yang, Qingyao Ai, Wei Jia, Youfeng Liu, Kai Lin, Yueyue Wu, Guozhi
798 Yuan, Yiran Hu, Wuyue Wang, Yiqun Liu, and Minlie Huang. LegalAgentBench: Evaluating
799 LLM agents in legal domain, 2024. URL <https://arxiv.org/abs/2412.17259>.
- 800 Haoran Li, Wenbin Hu, Huihao Jing, Yulin Chen, Qi Hu, Sirui Han, Tianshu Chu, Peizhao Hu, and
801 Yangqiu Song. PrivaCI-Bench: Evaluating privacy with contextual integrity and legal compliance,
802 2025. URL <https://arxiv.org/abs/2502.17041>.

- 810 Yueqi Li and Sanjay Goel. Artificial intelligence auditability and auditor readiness for audit-
811 ing artificial intelligence systems. *International Journal of Accounting Information Systems*,
812 56:100739, 2025. ISSN 1467-0895. doi: 10.1016/j.accinf.2025.100739. URL <https://www.sciencedirect.com/science/article/pii/S1467089525000156>.
813
- 814 Rensis Likert. *A Technique for the Measurement of Attitudes*. Number 140 in Archives of Psychol-
815 ogy. Archives of Psychology, New York, 1932.
816
- 817 Alexandra Lillo Campoy, Arlet Brufau i Centelles, and Albert Sabater Coll. The pio model (prin-
818 ciples, indicators and observables) on duties and rights: A self-assessment approach for com-
819 pliance of artificial intelligence data and systems regulations within the european union frame-
820 work. Technical report, Universitat de Girona, 2024. URL [https://dugi-doc.udg.edu/
821 bitstream/handle/10256/26967/9788484586692.pdf?sequence=1](https://dugi-doc.udg.edu/bitstream/handle/10256/26967/9788484586692.pdf?sequence=1). Available
822 as PDF.
- 823 Chun-Hsien Lin and Pu-Jen Cheng. Legal documents drafting with fine-tuned pre-trained large
824 language model, 2024. URL <https://arxiv.org/abs/2406.04202>.
- 825 Natasha Lomas. EU’s AI Act gets published in bloc’s Official Journal, starting clock on legal dead-
826 lines. [https://techcrunch.com/2024/07/12/eus-ai-act-gets-published-
827 in-blocs-official-journal-starting-clock-on-legal-deadlines](https://techcrunch.com/2024/07/12/eus-ai-act-gets-published-in-blocs-official-journal-starting-clock-on-legal-deadlines), 2024.
828 [Accessed 18-09-2025].
829
- 830 Megan Ma, Brandon Waldon, and Julian Nyarko. Conceptual questions in developing expert-
831 annotated data. In *Proceedings of the Nineteenth International Conference on Artificial Intel-
832 ligence and Law, ICAIL ’23*, pp. 427–431, New York, NY, USA, 2023. Association for Com-
833 puting Machinery. ISBN 9798400701979. doi: 10.1145/3594536.3595139. URL <https://doi.org/10.1145/3594536.3595139>.
834
- 835 Megan Ma, Aparna Sinha, Ankit Tandon, and Jennifer Richards. Generative AI legal landscape
836 2024. White paper, Stanford Law School, March 2024. URL [https://law.stanford.edu/
837 publications/generative-ai-legal-landscape-2024-2/](https://law.stanford.edu/publications/generative-ai-legal-landscape-2024-2/). Published by Stan-
838 ford Law School / CodeX in March 2024.
- 839 Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D. Manning,
840 and Daniel E. Ho. Hallucination-free? Assessing the reliability of leading AI leg-
841 al research tools. *Journal of Empirical Legal Studies*, 0(1):1–27, 2025. doi:
842 10.1111/jels.12413. URL [https://dho.stanford.edu/wp-content/uploads/
843 Legal_RAG_Hallucinations.pdf](https://dho.stanford.edu/wp-content/uploads/Legal_RAG_Hallucinations.pdf). Open access.
- 844 Gaurav Maheshwari, Dmitry Ivanov, and Kevin El Haddad. Efficacy of synthetic data as a bench-
845 mark, 2024. URL <https://arxiv.org/abs/2409.11968>.
846
- 847 Tobias Mahler. Between risk management and proportionality: The risk-based approach in the EU’s
848 Artificial Intelligence Act proposal. *The Swedish Law and Informatics Research Institute*, pp.
849 247–270, March 2022.
- 850 Barbara Makovec, Luis Rei, and Inna Novalija. Preparing AI for compliance: Initial steps of
851 a framework for teaching LLMs to reason about compliance. In *Companion Proceedings
852 of the 8th International Joint Conference on Rules and Reasoning (RuleML+RR’24)*, volume
853 3816, Bucharest, Romania, September 2024. CEUR Workshop Proceedings. URL <https://ceur-ws.org/Vol-3816/paper63.pdf>.
854
- 855 Bill Marino, Yaqub Chaudhary, Yulu Pi, Rui-Jie Yew, Preslav Aleksandrov, Carwyn Rahman,
856 William F. Shen, Isaac Robinson, and Nicholas D. Lane. Compliance Cards: Automated
857 EU AI Act compliance analyses amidst a complex AI supply chain, 2024. URL <https://arxiv.org/abs/2406.14758>.
858
- 859 Bertin Martens. Why artificial intelligence is creating fundamental challenges for competition pol-
860 icy. Technical report, Bruegel Policy Brief, 2024.
861
- 862 Lauren Martin, Nick Whitehouse, Stephanie Yiu, Lizzie Catterson, and Rivindu Perera. Better call
863 gpt, comparing large language models against lawyers, 2024. URL [https://arxiv.org/
abs/2401.16212](https://arxiv.org/abs/2401.16212).

- 864 Gabriele Mazzini and Salvatore Scalzo. The proposal for the Artificial Intelligence Act: Consider-
865 ations around some key concepts. In Camardi (ed.), *La via europea per l'Intelligenza artificiale*.
866 2023. doi: 10.2139/ssrn.4098809. URL <https://ssrn.com/abstract=4098809>. Avail-
867 able at SSRN.
- 868 Hans-W Micklitz and Giovanni Sartor. Compliance and enforcement in the AIA through AI. *Year-*
869 *book of European Law*, pp. yeae014, 02 2025. ISSN 0263-3264. doi: 10.1093/yel/yeae014. URL
870 <https://doi.org/10.1093/yel/yeae014>.
- 871 Eliza Mik. Caveat lector: Large language models in legal practice. Technical Report 2024-04, The
872 Chinese University of Hong Kong, Faculty of Law, 2024. URL [https://papers.ssrn.com/](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4757452)
873 [sol3/papers.cfm?abstract_id=4757452](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4757452). Posted March 15, 2024. Forthcoming in
874 *Rutgers Business Law Review*.
- 875 David Molnar. AI Unleashed: Mastering the Maze of the EU AI Act. *International University*
876 *Proceedings*, 2024. doi: https://doi.org/10.56461/iup_rlr.2024.5.ch12.
- 877 Colin E. Moriarty. The legal ethics of generative ai—part 3. *Colorado Lawyer*, October 2023.
878 URL [https://cl.cobar.org/features/the-legal-ethics-of-generative-](https://cl.cobar.org/features/the-legal-ethics-of-generative-ai-part-3/)
879 [ai-part-3/](https://cl.cobar.org/features/the-legal-ethics-of-generative-ai-part-3/). Accessed: 2025-09-21.
- 880 Jakob Mökander, Maria Axente, Federico Casolari, and Luciano Floridi. Conformity assessments
881 and post-market monitoring: A guide to the role of auditing in the proposed European AI regu-
882 lation. *Minds and Machines*, 32(2):241–268, November 2021. ISSN 1572-8641. doi: 10.1007/
883 s11023-021-09577-4. URL <http://dx.doi.org/10.1007/s11023-021-09577-4>.
- 884 National Institute of Standards and Technology. NIST AI 100-1 artificial intelligence risk manage-
885 ment framework (AI RMF 1.0). Framework, National Institute of Standards and Technology, Jan-
886 uary 2023. URL <https://www.nist.gov/itl/ai-risk-management-framework>.
887 Released on January 26, 2023.
- 888 Compliance Podcast Network. Stepping up and stepping forward: The future of compliance in an
889 age of AI and deregulation. [https://compliancepodcastnetwork.net/stepping-](https://compliancepodcastnetwork.net/stepping-up-and-stepping-forward-the-future-of-compliance-in-an-age-of-ai-and-deregulation/)
890 [up-and-stepping-forward-the-future-of-compliance-in-an-age-of-](https://compliancepodcastnetwork.net/stepping-up-and-stepping-forward-the-future-of-compliance-in-an-age-of-ai-and-deregulation/)
891 [ai-and-deregulation/](https://compliancepodcastnetwork.net/stepping-up-and-stepping-forward-the-future-of-compliance-in-an-age-of-ai-and-deregulation/), April 4 2025. Podcast episode / blog post.
- 892 Tai D. Nguyen, Long H. Pham, and Jun Sun. Autolaw: Enhancing legal compliance in large
893 language models via case law generation and jury-inspired deliberation, 2025. URL <https://arxiv.org/abs/2505.14015>.
- 894 OpenAI. Gpt-4o mini: Advancing cost-efficient intelligence, July 2024. URL
895 [https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/)
896 [intelligence/](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/). Accessed: YYYY-MM-DD.
- 897 OpenAI. GPT-5 system card. Technical report, OpenAI, August 13 2025a. URL [https://](https://cdn.openai.com/gpt-5-system-card.pdf)
898 cdn.openai.com/gpt-5-system-card.pdf. OpenAI system card, August 13, 2025.
- 899 OpenAI. OpenAI o3-mini system card. Technical report, OpenAI, January 31 2025b. URL [https://](https://openai.com/index/o3-mini-system-card/)
900 openai.com/index/o3-mini-system-card/. Accessed: 2025-09-21.
- 901 OpenAI. Openai o3 and o4-mini system card. [https://openai.com/index/o3-o4-mini-](https://openai.com/index/o3-o4-mini-system-card/)
902 [system-card/](https://openai.com/index/o3-o4-mini-system-card/), April 16 2025c. Accessed: 2025-09-21.
- 903 OpenAI. OpenAI o3 and o4-mini system card. [https://cdn.openai.com/pdf/2221c875-](https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf)
904 [02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf](https://cdn.openai.com/pdf/2221c875-02dc-4789-800b-e7758f3722c1/o3-and-o4-mini-system-card.pdf), 2025.
- 905 OpenAI et al. GPT-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- 906 Guilherme Palumbo, Miguel Guimarães, Davide Carneiro, Goretí Marreiros, and Victor Alves.
907 Observability-driven ai governance: A framework for compliance and audit readiness under the
908 EU AI Act. In Daniel H. de la Iglesia, Juan F. de Paz Santana, and Alfonso J. López Rivero
909 (eds.), *New Trends in Disruptive Technologies, Tech Ethics and Artificial Intelligence*, pp. 402–
910 413, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-99474-6.

- 918 Ji-Lun Peng, Sijia Cheng, Egil Diau, Yung-Yu Shih, Po-Heng Chen, Yen-Ting Lin, and Yun-
919 Nung Chen. A survey of useful LLM evaluation, 2024. URL [https://arxiv.org/abs/
920 2406.00936](https://arxiv.org/abs/2406.00936).
- 921 Nicholas Pipitone and Ghita Hourir Alami. LegalBench-RAG: A benchmark for retrieval-augmented
922 generation in the legal domain, 2024. URL <https://arxiv.org/abs/2408.10343>.
- 923 Eric A. Posner and Shivam Saran. Judge AI: Assessing large language models in judicial
924 decision-making. Research Paper 25-03, University of Chicago Law School, Coase-Sandor
925 Institute for Law and Economics, Chicago, IL, January 2025. URL [https://ssrn.com/
926 abstract=5098708](https://ssrn.com/abstract=5098708). Posted January 16, 2025; last revised January 28, 2025.
- 927 Matteo Prandi, Vincenzo Suriani, Federico Pierucci, Marcello Galisai, Daniele Nardi, and Piercosma
928 Bisconti. Bench-2-cop: Can we trust benchmarking for EU AI compliance?, 2025. URL [https://arxiv.org/abs/
929 //arxiv.org/abs/2508.05464](https://arxiv.org/abs/2508.05464).
- 930 Inioluwa Deborah Raji, Peggy Xu, Colleen Honigsberg, and Daniel Ho. Outsider oversight: De-
931 signing a third party audit ecosystem for AI governance. In *Proceedings of the 2022 AAAI/ACM
932 Conference on AI, Ethics, and Society*, AIES '22, pp. 557–571, New York, NY, USA, July 2022.
933 Association for Computing Machinery. ISBN 9781450392471. doi: 10.1145/3514094.3534181.
934 URL <https://dl.acm.org/doi/10.1145/3514094.3534181>.
- 935 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions
936 for machine comprehension of text, 2016. URL <https://arxiv.org/abs/1606.05250>.
- 937 Vishvakshan Rasiah, Ronja Stern, Veton Matoshi, Matthias Stürmer, Ilias Chalkidis, Daniel E Ho,
938 and Joel Niklaus. One law, many languages: Benchmarking multilingual legal reasoning for
939 judicial support. In *ICLR 2024 Workshop on Data-centric Machine Learning Research (DMLR):
940 Harnessing Momentum for Science*. <https://openreview.net/forum>, 2024.
- 941 Anka Reuel, Ben Bucknall, Stephen Casper, Tim Fist, Lisa Soder, Onni Aarne, Lewis Ham-
942 mond, Lujain Ibrahim, Alan Chan, Peter Wills, Markus Anderljung, Ben Garfinkel, Lennart
943 Heim, Andrew Trask, Gabriel Mukobi, Rylan Schaeffer, Mauricio Baker, Sara Hooker, Irene
944 Solaiman, Alexandra Sasha Luccioni, Nitarshan Rajkumar, Nicolas Moës, Jeffrey Ladish, Neel
945 Guha, Jessica Newman, Yoshua Bengio, Tobin South, Alex Pentland, Sanmi Koyejo, Mykel J.
946 Kochenderfer, and Robert Trager. Open problems in technical AI governance, 2024a. URL
947 <https://arxiv.org/abs/2407.14981>.
- 948 Anka Reuel, Lisa Soder, Ben Bucknall, and Trond Arne Undheim. Position paper: Technical re-
949 search and talent is needed for effective AI governance, 2024b. URL [https://arxiv.org/
950 abs/2406.06987](https://arxiv.org/abs/2406.06987).
- 951 Oliver Roberts. Legal AI unfiltered: 16 tech leaders on AI replacing lawyers, the billable hour,
952 and hallucinations. *The National Law Review*. URL [https://natlawreview.com/
953 article/legal-ai-unfiltered-16-tech-leaders-ai-replacing-
954 lawyers-billable-hour-and](https://natlawreview.com/article/legal-ai-unfiltered-16-tech-leaders-ai-replacing-lawyers-billable-hour-and). Volume XV, Number 256.
- 955 Anne Rother, Uli Niemann, Tommy Hielscher, Henry Völzke, Till Ittermann, and Myra
956 Spiliopoulou. Assessing the difficulty of annotating medical data in crowdworking with help
957 of experiments. *PLoS one*, 16(7):e0254764, 2021.
- 958 Jonas Schuett. From principles to rules: A regulatory approach for frontier AI. *arXiv preprint
959 arXiv:2407.07300*, 2024.
- 960 Zejiang Shen, Kyle Lo, Lauren Yu, Nathan Dahlberg, Margo Schlanger, and Doug Downey.
961 Multi-lexsum: Real-world summaries of civil rights lawsuits at multiple granularities.
962 In *Proceedings of the Neural Information Processing Systems Track on Datasets and
963 Benchmarks (NeurIPS 2022)*, 2022. URL [https://proceedings.neurips.cc/
964 paper_files/paper/2022/file/552ef803bef9368c29e53c167de34b55-
965 Paper-Datasets_and_Benchmarks.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/552ef803bef9368c29e53c167de34b55-Paper-Datasets_and_Benchmarks.pdf). Track on Datasets and Benchmarks.
- 966 Charlotte Siegmann and Markus Anderljung. The Brussels Effect and artificial intelligence: How
967 EU regulation will impact the global AI market. *CoRR*, abs/2208.12645, 2022. doi: 10.48550/
968 ARXIV.2208.12645. URL <https://doi.org/10.48550/arXiv.2208.12645>.
- 969
970
971

- 972 Marco Siino, Mariana Falco, Daniele Croce, and Paolo Rosso. Exploring LLMs applications in
973 law: A literature review on current legal NLP approaches. *IEEE Access*, January 2025. doi:
974 10.1109/ACCESS.10850911. Article No. 10850911.
- 975
- 976 Mona Sloane and Elena Wüllhorst. A systematic review of regulatory strategies and transparency
977 mandates in AI regulation in Europe, the United States, and Canada. *Data & Policy*, 7:e11, 2025.
- 978
- 979 Nikolaos Sourlos, Rozemarijn Vliegenthart, Gonalo J. Santinha, Michail E. Klontzas, Roberto Cuoco-
980 colo, Marc Huisman, and Pierre M.A. van Ooijen. Recommendations for the creation of bench-
981 mark datasets for reproducible artificial intelligence in radiology. *Insights into Imaging*, 15(1):
982 248, Oct 2024. doi: 10.1186/s13244-024-01833-2.
- 983 F. Sovrano, E. Hine, S. Anzolut, et al. Simplifying software compliance: AI technologies in drafting
984 technical documentation for the AI act. *Empirical Software Engineering*, 30(91), 2025. doi:
985 10.1007/s10664-025-10645-x.
- 986
- 987 Sotirios Stampernas and Costas Lambrinouidakis. A framework for compliance with regulation (EU)
988 2024/1689 for small and medium-sized enterprises. *Journal of Cybersecurity and Privacy*, 5(3),
989 2025. ISSN 2624-800X. doi: 10.3390/jcp5030040. URL [https://www.mdpi.com/2624-](https://www.mdpi.com/2624-800X/5/3/40)
990 [800X/5/3/40](https://www.mdpi.com/2624-800X/5/3/40).
- 991 Weihang Su, Baoqing Yue, Qingyao Ai, Yiran Hu, Jiaqi Li, Changyue Wang, Kaiyuan Zhang,
992 Yueyue Wu, and Yiqun Liu. JuDGE: Benchmarking judgment document generation for Chinese
993 legal system, 2025. URL <https://arxiv.org/abs/2503.14258>.
- 994
- 995 The Institute of Internal Auditors. The IIA’s updated AI auditing framework, Decem-
996 ber 2023. URL [https://www.theiia.org/en/content/tools/professional/](https://www.theiia.org/en/content/tools/professional/2023/the-iias-updated-ai-auditing-framework/)
997 [2023/the-iias-updated-ai-auditing-framework/](https://www.theiia.org/en/content/tools/professional/2023/the-iias-updated-ai-auditing-framework/).
- 998
- 999 E. Thelisson and H. Verma. Conformity assessment under the EU AI Act general approach. *AI*
1000 *Ethics*, 4:113–121, 2024. doi: 10.1007/s43681-023-00402-5.
- 1001 Fabian Walke, Lars Bennek, and Till J. Winkler. Artificial intelligence explainability require-
1002 ments of the AI act and metrics for measuring compliance. In *Digital Responsibility: So-*
1003 *cial, Ethical, Ecological Implications of IS*, 18. *Internationale Tagung Wirtschaftsinformatik (WI*
1004 *2023)*, September 18-21, 2023, Paderborn, Germany, pp. 77. AISeL, 2023. URL [https:](https://aisel.aisnet.org/wi2023/77)
1005 [//aisel.aisnet.org/wi2023/77](https://aisel.aisnet.org/wi2023/77).
- 1006
- 1007 Jianwei Wang, Mengqi Wang, Yinsi Zhou, Zhenchang Xing, Qing Liu, Xiwei Xu, Wenjie Zhang,
1008 and Liming Zhu. LLM-based HSE compliance assessment: Benchmark, performance, and ad-
1009 vancements, 2025a. URL <https://arxiv.org/abs/2505.22959>.
- 1010 Jiaqi Wang, Huan Zhao, Zhenyuan Yang, Peng Shu, Junhao Chen, Haobo Sun, Ruixi Liang,
1011 Shixin Li, Pengcheng Shi, Longjun Ma, Zongjia Liu, Zhengliang Liu, Tianyang Zhong, Yutong
1012 Zhang, Chong Ma, Xin Zhang, Tuo Zhang, Tianli Ding, Yudan Ren, Tianming Liu, Xi Jiang,
1013 and Shu Zhang. Legal evaluations and challenges of large language models, 2024. URL
1014 <https://arxiv.org/abs/2411.10137>.
- 1015
- 1016 Steven H. Wang, Antoine Scardigli, Leonard Tang, Wei Chen, Dimitry Levkin, Anya Chen, Spencer
1017 Ball, Thomas Woodside, Oliver Zhang, and Dan Hendrycks. MAUD: An expert-annotated legal
1018 NLP dataset for merger agreement understanding, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2301.00876)
1019 [2301.00876](https://arxiv.org/abs/2301.00876).
- 1020 Steven H. Wang, Maksim Zubkov, Kexin Fan, Sarah Harrell, Yuyang Sun, Wei Chen, Andreas
1021 Plesner, and Roger Wattenhofer. Acord: An expert-annotated retrieval dataset for legal contract
1022 drafting, 2025b. URL <https://arxiv.org/abs/2501.06582>.
- 1023
- 1024 Zhaoyang Wang, Jinqi Jiang, Huichi Zhou, Wenhao Zheng, Xuchao Zhang, Chetan Bansal, and
1025 Huaxiu Yao. Verifiable format control for large language model generations. *arXiv preprint*
arXiv:2502.04498, 2025c.

- 1026 Zifeng Wang, Junyi Gao, Benjamin Danek, Brandon Theodorou, Ruba Shaik, Shivashankar Thati,
1027 Seunghyun Won, and Jimeng Sun. InformGen: An AI copilot for accurate and compliant
1028 clinical research consent document generation, 2025d. URL [https://arxiv.org/abs/
1029 2504.00934](https://arxiv.org/abs/2504.00934).
- 1030 Ziqi Wang and Boqin Yuan. L-MARS: Legal multi-agent workflow with orchestrated reasoning and
1031 agentic search. *arXiv preprint arXiv:2509.00761*, 2025.
- 1032 Shomir Wilson, Florian Schaub, Aswarth Abhilash Dara, Frederick Liu, Sushain Cherivirala,
1033 Pedro Giovanni Leon, Mads Schaarup Andersen, Sebastian Zimmeck, Kanthashree Mysore
1034 Sathyendra, N. Cameron Russell, Thomas B. Norton, Eduard Hovy, Joel Reidenberg, and Nor-
1035 man Sadeh. The creation and analysis of a website privacy policy corpus. In *Proceedings
1036 of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016)*, pp.
1037 1330–1340, Berlin, Germany, 2016. Association for Computational Linguistics. URL [https:
1038 //aclanthology.org/P16-1126.pdf](https://aclanthology.org/P16-1126.pdf).
- 1039 Weiyue Wu and Shaoshan Liu. Why compliance costs of AI commercialization may be hold-
1040 ing start-ups back. [https://studentreview.hks.harvard.edu/why-compliance-
1041 costs-of-ai-commercialization-maybe-holding-start-ups-back/](https://studentreview.hks.harvard.edu/why-compliance-costs-of-ai-commercialization-maybe-holding-start-ups-back/), 2023.
1042 [Accessed 18-09-2025].
- 1043 Yixin Wu and Benjamin van Rooij. Compliance dynamism: Capturing the polynormative and situ-
1044 ational nature of business responses to law. *Journal of Business Ethics*, 168:579–591, 2021. doi:
1045 10.1007/s10551-019-04234-4.
- 1046 xAI. Grok 3 Beta — the age of reasoning agents. <https://x.ai/news/grok-3>, February 19
1047 2025a. Accessed: 2025-09-21.
- 1048 xAI. Grok 4. <https://x.ai/news/grok-4>, 2025b. Accessed: 2025-09-21.
- 1049 Danning Xie, Byungwoo Yoo, Nan Jiang, Mijung Kim, Lin Tan, Xiangyu Zhang, and Judy S. Lee.
1050 How effective are large language models in generating software specifications?, 2025. URL
1051 <https://arxiv.org/abs/2306.03324>.
- 1052 Rui-Jie Yew, Bill Marino, and Suresh Venkatasubramanian. Red teaming AI policy: A taxonomy
1053 of avoision and the EU AI Act. In *Proceedings of the 2025 ACM Conference on Fairness, Ac-
1054 countability, and Transparency*, FAccT '25, pp. 404–415, New York, NY, USA, 2025. Associa-
1055 tion for Computing Machinery. ISBN 9798400714825. doi: 10.1145/3715275.3732028. URL
1056 <https://doi.org/10.1145/3715275.3732028>.
- 1057 Lucia Zheng, Neel Guha, Javokhir Arifov, Sarah Zhang, Michal Skreta, Christopher D. Manning,
1058 Peter Henderson, and Daniel E. Ho. A reasoning-focused legal retrieval benchmark. In *Pro-
1059 ceedings of the Symposium on Computer Science and Law on ZZZ*, CSLAW '25, pp. 169–193.
1060 ACM, March 2025. doi: 10.1145/3709025.3712219. URL [http://dx.doi.org/10.1145/
1061 3709025.3712219](http://dx.doi.org/10.1145/3709025.3712219).
- 1062 Linwu Zhong, Ziyi Zhong, Zinian Zhao, Siyuan Wang, Kevin D. Ashley, and Matthias Grabmair.
1063 Automatic summarization of legal decisions using iterative masking of predictive sentences.
1064 In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law*,
1065 ICAIL '19, pp. 163–172, New York, NY, USA, 2019. Association for Computing Machinery.
1066 ISBN 9781450367547. doi: 10.1145/3322640.3326728. URL [https://doi.org/10.1145/
1067 3322640.3326728](https://doi.org/10.1145/3322640.3326728).
- 1068 Andreas Östling, Holli Sargeant, Huiyuan Xie, Ludwig Bull, Alexander Terenin, Leif Jonsson, Måns
1069 Magnusson, and Felix Steffek. The Cambridge Law Corpus: A dataset for legal AI research, 2024.
1070 URL <https://arxiv.org/abs/2309.12269>.

1077 A LLM USAGE IN THIS WORK

1078 LLMs were used in this work as follows:
1079

- LLMs were used to aid or polish writing. Specifically, they were used to identify errors or weaknesses in writing, to generate an initial draft of the Conclusion (Section 8), and to generate bibtext.
- LLMs were used for retrieval and discovery. Specifically, they were used for finding related work.
- LLMs were used for other purposes, such as the generation of the samples (described in Section 2), tables, and streamlining code development.

B ADDITIONAL FIGURES

This section presents supplementary figures and tables that, for the sake of brevity, have been omitted from the main body of this paper.

Figure 3 (left) shows confusion matrices comparing the median human expert compliance scores (rows) with the LLM scores (columns). Darker cells indicate more frequent score combinations. The top matrix is for Gemini 2.5 Pro, the best-performing model across most evaluated metrics. The majority of score combinations sit along the diagonal (72/120), showing strong agreement between Gemini 2.5 Pro and the median human expert annotator. The bottom matrix averages the frequency of score pairings over all evaluated models, and here, the distribution is more diffuse.

Figure 3 (right) reports the mean absolute error (MAE) between LLM compliance scores and median human expert scores, broken down by use case (columns) and article (rows). Darker cells reflect larger errors. The top heatmap shows the MAE breakdown for Gemini 2.5 Pro, while the bottom heatmap averages MAE across all evaluated models. Notably, in the bottom heatmap only three cells exceed an MAE of 1.0, indicating a consistently strong average agreement between LLMs and human experts.

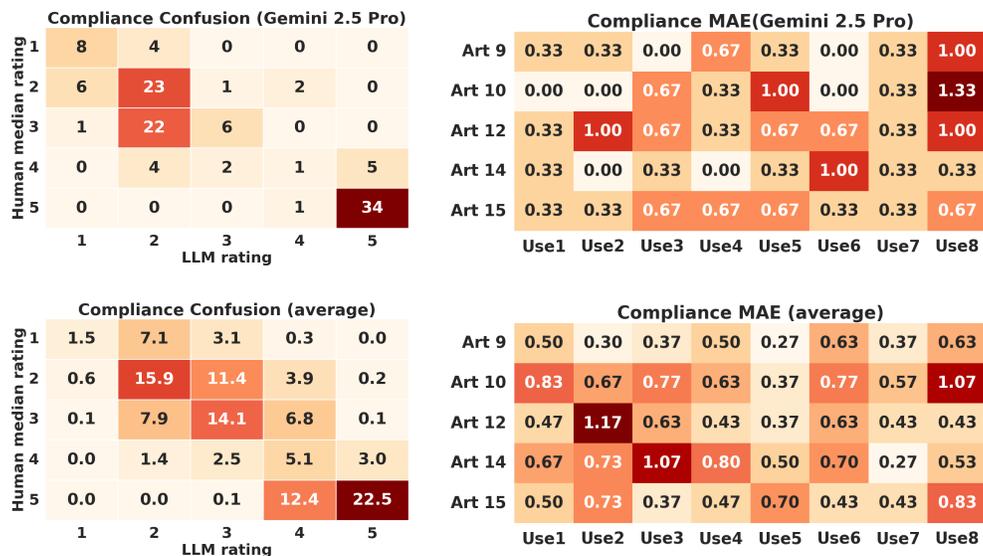


Figure 3: **Heatmaps of compliance performance.** The left panels show the distribution of compliance ratings (in ‘confusion matrix’), comparing the median human expert with LLMs. The right panels show mean absolute error (MAE) across use cases and articles. Results are shown for Gemini 2.5 Pro (top) and as an average over all evaluated LLMs (bottom). Use1-8 reflect the intended uses from Table 2.

Figure 4 plots model cost (x-axis, output price per M tokens) against compliance agreement with human expert ratings (y-axis, Cohen’s κ weighted quadratically).³ Each point on the graph corre-

³Gemma 3 is not included as it is not available via paid API access.

sponds to an evaluated LLM and those highlighted in red (Gemini 2.5 Pro and Grok 3 mini) lie on the Pareto frontier meaning that no model is both cheaper and more compliant.

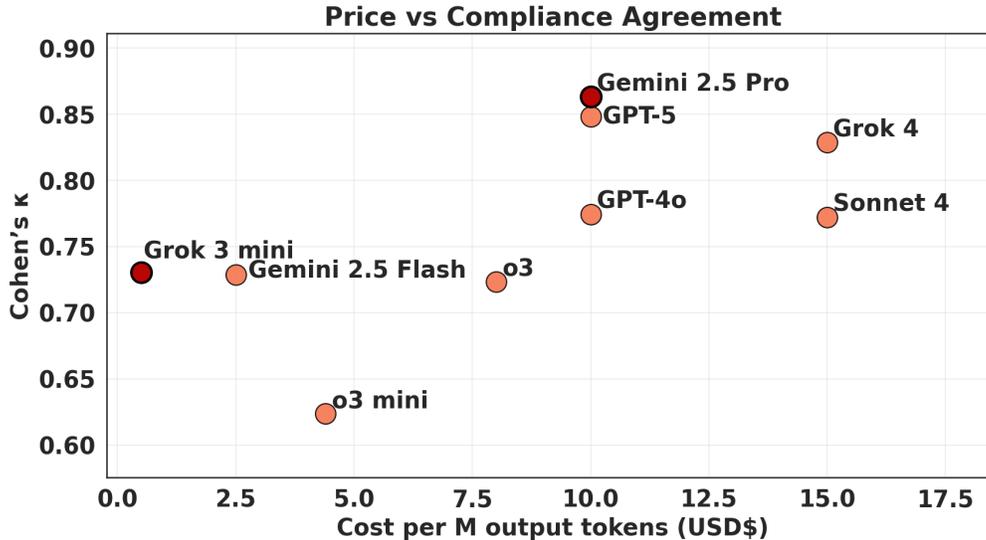


Figure 4: **Pareto frontier of model cost versus compliance agreement (Cohen’s κ)**. Each point represents a model, plotted by price (x-axis) and agreement with human expert ratings (y-axis). Pareto-efficient models are shown with red markers. Labels denote model names. [We were not able to capture the cost of our particular implementation of Gemma, as API use is free but has a low restriction on output tokens.](#)

Table 4 presents the accuracy of LLMs in replicating human expert scores exactly and, where mismatches occur, the direction of error. Despite prompts designed to mitigate sycophancy and acquiescence bias, some models tended to assign higher compliance scores than human experts. Two of the worst models in this regard were o3 mini and GPT-4o, with o3 mini strictly exceeding the median human expert score in 54.2% of excerpts, while only strictly falling below in 1.7%. The model least prone to over-estimating the compliance level of excerpts was Gemini 2.5 Pro, which was also the model whose scores exactly matched those of the median human expert most frequently (for 60% of excerpts). Gemini 2.5 Flash achieved the best F1 score (0.913), meaning it excelled at correctly flagging compliant cases (score 4-5) and avoiding mislabelling non-compliant ones (score 1-3).

Model	% Exact (\uparrow)	% Over (\downarrow)	% Under (\downarrow)	F1 (\uparrow)
GPT-5	57.5	16.7	25.8	0.903
GPT-4o	52.5	42.5	5.0	0.846
o3	38.3	20.8	40.8	0.903
o3 mini	44.2	54.2	<u>1.7</u>	0.736
Sonnet 4	46.7	20.0	33.3	0.845
Gemini 2.5 Pro	<u>60.0</u>	<u>10.0</u>	30.0	0.911
Gemini 2.5 Flash	41.7	22.5	35.8	<u>0.913</u>
Gemma 3	40.0	39.2	20.8	0.796
Grok 4	58.3	30.8	10.8	0.860
Grok 3 mini	53.3	42.5	4.2	0.830

Table 4: **Compliance Likert score differences between LLMs and human experts**. Columns report accuracy of LLMs across AIReg-Bench, including the percentage of exact matches, over-estimates, and under-estimates relative to the median human expert, as well as the F1 score from binary classification (scores 1–3 vs. 4–5).

Table 5 presents ablation results for GPT-4o, comparing the baseline model to three ablated versions: one without a tone prompt (see Appendix G), one with a harsh tone prompt, and one without access to the relevant AI Act text. The baseline tone prompt (“Your scores for both compliance and plausibility should be well-calibrated and objective. They should be rigorous but fair.”) achieves the best performance with respect to Cohen’s κ_w , Spearman’s ρ , and MAE.

The harsh tone prompt (“Your scores for both compliance and plausibility should be critical. They should be harsh but fair.”) reduces bias, but at the cost of declines across all three other metrics. When access to the text of the articles in the EU AI Act is removed, all performance metrics drop substantially, with GPT-4o’s Cohen’s κ_w falling to 0.654, just above o3 mini’s performance when provided with the text.

Model	κ_w (\uparrow)	ρ (\uparrow)	Bias ($\rightarrow 0$)	MAE (\downarrow)
GPT-4o baseline	<u>0.775</u>	<u>0.842</u>	0.458	<u>0.558</u>
ablation (none)	0.759	<u>0.842</u>	0.492	0.575
ablation (harsh)	0.722	0.791	<u>0.125</u>	0.642
ablation (w/o articles)	0.654	0.752	0.583	0.717

Table 5: **Ablation analysis of GPT-4o.** Ablations include removing or altering the prompt modifier (to be harsher), or withholding access to the AI Act text. Columns report agreement between LLM and median human expert scores across AIReg-Bench: quadratically weighted Cohen’s κ_w ; Spearman’s ρ ; Bias (mean signed difference, LLM–human); and MAE (mean absolute error).

Table 6 includes the results for well-known open-source language models fine-tuned for legal tasks: LLM Saul-7B-Instruct (Colombo et al., 2024b) and Saul-54B-Instruct (Colombo et al., 2024a). Both models are relatively small in size and, even with the potential advantages of fine-tuning on legal material, they underperform the weakest general-purpose frontier model in our evaluation with respect to Cohen’s κ_w , achieving 0.183 and 0.596 respectively, compared to o3 mini’s 0.624.

That said, the stark improvement from Saul-7B-Instruct to Saul-54B-Instruct is significant, highlighting the benefits of scaling to larger models. Notably, Saul-7B-Instruct struggles to consistently format its outputs as requested, a limitation observed in smaller language models more generally (Wang et al., 2025c). Many of its answers therefore had to be resampled until they were formatted appropriately.

Model	κ_w (\uparrow)	ρ (\uparrow)	Bias ($\rightarrow 0$)	MAE (\downarrow)
Saul-7B-Instruct	0.183	0.311	0.550	1.167
Saul-54B-Instruct	0.596	0.813	0.792	0.825

Table 6: **Experiments with Saul:** Columns report agreement between LLM Saul-7B-Instruct (Colombo et al., 2024b) and Saul-54B-Instruct (Colombo et al., 2024a) and median human compliance scores across AIReg-Bench: quadratically weighted Cohen’s κ_w ; Spearman’s ρ ; Bias (mean signed difference, LLM–human); and MAE (mean absolute error). Outputs from Saul-7B-Instruct were resampled until they were provided in a parsable format.

Table 8 summarizes the results of the alternative annotator test proposed by Calderon et al. (2025). The goal is to check whether replacing a human annotator with a given model would preserve the overall decisions. The winning rate captures the fraction of humans for whom the model would serve as an adequate substitute, while the average advantage probability reflects how often the model’s annotation is at least as good as a human’s across items. In our discrete label setup, a model wins on an item if its label is closer to the remaining annotators’ labels than the human’s.

	GPT-5	GPT-4o	o3	o3 mini	Sonnet 4	Gemini 2.5 Pro	Gemini 2.5 Flash	Gemma 3	Grok 4	Grok 3 mini
GPT-5	–									
GPT-4o	0.73	–								
o3	0.70	0.57	–							
o3 mini	0.55	0.87	0.45	–						
Sonnet 4	0.84	0.75	0.74	0.55	–					
Gemini 2.5 Pro	0.84	0.70	0.71	0.56	0.85	–				
Gemini 2.5 Flash	0.73	0.67	0.75	0.51	0.77	0.74	–			
Gemma 3	0.70	0.73	0.55	0.69	0.62	0.64	0.64	–		
Grok 4	0.85	0.84	0.59	0.65	0.86	0.86	0.72	0.61	–	
Grok 3 mini	0.68	0.92	0.55	0.77	0.78	0.67	0.70	0.78	0.82	–

Table 7: **Pairwise Cohen’s Kappa scores.** This table shows the results of performing pairwise Cohen’s Kappa scores of the models in our evaluation.

Model	Items	Annotators	Winning Rate	Average Advantage Probability
GPT-5	60	3	0.6667	0.8500
GPT-4o	60	3	0.0000	0.7944
o3	60	3	0.0000	0.6556
o3 mini	60	3	0.0000	0.7556
Sonnet 4	60	3	0.0000	0.7333
Gemini 2.5 Pro	60	3	1.0000	0.9111
Gemini 2.5 Flash	60	3	0.0000	0.6722
Gemma 3	60	3	0.0000	0.6722
Grok 4	60	3	0.6667	0.8444
Grok 3 mini	60	3	0.0000	0.7500

Table 8: **Alternative Annotator Test.** This table shows the results of the test proposed by Calderon et al. (2025), which answers the question: if we replaced human annotators with another LLM, would the resulting labels stay the same? Here, winning rate is defined as the fraction of human annotators for whom the LLM passes the test (i.e., the model’s annotation is closer to the other human scores than the alternative annotation). Meanwhile, average advantage probability is defined as the average probability that the LLM’s annotation is at least as good as a human experts’ on each item. Our implementation reproduces the paper’s procedure as-is.

C DESIGN CRITERIA

The design criteria for AIReg-Bench were as follows:

- In keeping with good benchmark dataset practices (Sourlos et al., 2024), the samples should be representative of real-world technical documentation used in AIA compliance assessments. Since AIReg-Bench is intended for evaluating LLMs’ ability to perform AIR compliance assessments, we sought to replicate as closely as possible the documentation that a human compliance assessor would consult during this process.
- The samples should only depict HRAI systems (EU, 2024, Art. 6) that are within the scope of the AIA (but not within the scope of any of its prohibitions or exceptions) and should

1296 be drafted as if created by the AI system’s provider (i.e., developer) (EU, 2024, Art. 2).
 1297 These design criteria, we argue, add a degree of realism to the dataset, since providers of AI
 1298 systems outside of these boundaries are less well incentivized to create detailed technical
 1299 documentation. Moreover, by focusing solely on in-scope HRAI systems, it ensures that
 1300 every document we generate is densely packed with compliance-critical details, many of
 1301 which may not be required for assessing lower-risk systems.

- 1302 • Aside from a high-level overview of the AI system, each sample’s contents should be con-
 1303 strained to specific AIA requirements for HRAI systems. This ensures the benchmark tests
 1304 models on fine-grained compliance analysis rather than on their ability to interpret overly
 1305 broad or generic descriptions.
- 1306 • In order to achieve a diverse distribution, the samples should be able to reflect a variety of
 1307 use cases (i.e., intended uses) as well as different compliance scenarios (either compliant or
 1308 non-compliant with relevant Articles of the AIA). Collectively, they should cover a variety
 1309 of intended uses and compliance profiles: that is, some systems are compliant with the
 1310 AIA, while others are not — and, in the case of the latter, the reasons for non-compliance
 1311 vary.

1312 1313 D COMPLIANCE EXPERT INTERVIEWS

1314
1315 Since AIReg-Bench is intended for evaluating LLMs’ ability to perform AIR compliance assess-
 1316 ments, we sought to replicate as closely as possible the technical documentation that a human com-
 1317 pliance assessor would consult during this process. To better understand the structure and contents
 1318 of this particular documentation, we interviewed six compliance experts (distinct from our six an-
 1319 notators), asking them to provide details about the materials they consult (or would expect to be
 1320 consulted) during AIR compliance assessments, including but not limited to those mandated by the
 1321 AIA (EU, 2024, Art. 43).

1322 Perhaps owing to the new and evolving nature of AIR compliance assessments, the consensus among
 1323 interviewees was that there are still no universal standards for the materials to be consulted during
 1324 this process. That said, the materials that were most commonly referenced by interviewees were
 1325 summaries of an AI system’s attributes and its development process — including, but not limited to
 1326 model cards, data cards, descriptions of data preparation, training and red-teaming processes, and
 1327 descriptions of governance or guardrail measures.

1328 Some of our experts suggested that these materials might be curated in preparation for a compliance
 1329 assessment using business records and auditee interviews, and that several such materials may be
 1330 integrated into a single instance of technical documentation. These interviewees indicated that, in
 1331 practice, one such technical document can serve as the primary artifact in compliance assessments,
 1332 even though compliance assessors may draw on a wider range of materials through iterative dia-
 1333 logue. To reflect the central role of technical documentation and to ensure our benchmark is simple
 1334 to use, we represent each AI system with a single integrated technical document, rather than many
 1335 such materials.

1336 Although interviewees consistently highlighted a lack of clear standards for compliance assess-
 1337 ments, many regarded the provisions and annexes of the AIA related to technical documentation as
 1338 among the clearest and most detailed guidance for AIR assessments. Accordingly, our dataset is
 1339 predominantly built around this regulation and, in particular, Annex IV and Chapter III, Section 2
 1340 of the Act (EU, 2024, Ann. IV, Chap. III(2)) — which we found to be most relevant when produc-
 1341 ing technical documentation for compliance assessments. By focusing almost-entirely on just these
 1342 two parts of the Act and omitting its less relevant provisions or any ancillary requirements (such as
 1343 harmonized standards), our technical documentation remains manageable in length, concentrating
 1344 exclusively on the core requirements for a compliance assessment.

1345 1346 E SAMPLE GENERATION PIPELINE PROMPTS

1347
1348 Listed below are the prompts that were fed to gpt-4.1-mini during the sample generation pipeline,
 1349 as well as the annotation instructions given to humans and LLM. Additional line breaks have been
 added for readability.

1350 E.1 SYSTEM OVERVIEW PROMPT

1351

1352 Your task is to generate four distinct AI system descriptions for the provided intended use.

1353

1354

1355

1356

1357

1358

1359

1360

Each AI system must employ only one or two domain-appropriate types of machine learning models or algorithms. You should pick the algorithm you feel is most appropriate for the use case in the contemporary era, but here are some examples of the types of algorithms that you might choose: MLP, CNN, Transformers (encoder-only, decoder-only, or encoder-decoder), SVM, RNN, Naive Bayes, GNN, Random Forest, KNN, GBDT, Linear Regression, transformer-based Large Language Model (LLM), transformer-based Vision Language Model (VLM), diffusion-based text-to-image generation model, or similar. Transformers can be used in distinct ways, including for processing different data types such as tabular data, text, audio, API calls, and more.

1361

1362

1363

1364

1365

1366

1367

Your choices of models or algorithms should reflect those likely to be deployed in 2025. You should focus on realism given the particular application as well as domain-appropriateness. Systems must not employ biometric technologies or violate Article 5 of the EU AI Act.

Start each description with: “System 1.” “System 2.” “System 3.” “System 4.” Each description should be a single continuous paragraph. There should be an empty line separating each system description from the next.

1368

1369

*****Intended use*****

{*Insert intended use*}

1370

1371

1372

1373

1374

1375

1376

1377

1378

For each system, provide a concise description consisting of the following four components, each in a single sentence:

[a] System Name and Type: State the AI system’s name and the machine learning models or algorithms it relies on. When naming these AI systems, ensure the names are diverse, realistic, and professional (while also fictional). System names should be formed by using multiple, separate words; never form system names by concatenating words (e.g., “EducationSmart”) or using portmanteaus (e.g., “EduBoost”).

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

[b] System use: Provide the system’s intended usage and the specific sector or area it will serve.

[c] Objective and Technological Capabilities: Describe the system’s primary objective and outline the technological capabilities that enable the achievement of this objective.

[d] Provider: Identify the natural or legal person, public authority, agency or other body that develops the system or that has the system developed and places it on the market or puts the system into service under its own name or trademark, whether for payment or free of charge. Where providers are companies, ensure their names are realistic and professional (while fictional); please do not simply concatenate words (e.g., “EducationSmart”) or use portmanteaus (e.g., “EduBoost”).

1389

1390

1391

1392

1393

1394

1395

1396

[e] The AI subject: Define the persons (or groups of persons) who are directly affected by the AI system, experiencing its outcomes and consequences.

Return nothing else outside those four system descriptions, formatted in the manner outlined above.

1397

1398

1399

1400

1401

1402

1403

E.2 VIOLATION PROMPT

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

*****Task*****

You are acting as an EU AI Act compliance expert and machine-learning specialist. You will be given an extract from the EU AI Act and a short description of an AI system. Your task is to generate three distinct ways that the AI system could violate the provided extract.

You must not produce:

- violations that reflect errors in the technical documentation rather than the system (e.g., lies or missing details).
- violations without any justification, evidence, or granular details (e.g., “the model is biased” without elaboration).

You must produce:

- 1404 • violations that are plausible and realistic
 1405
 1406 • violations that are subtle and not blatant

1407
 1408 Start each example with: “Example 1.”, “Example 2.”, “Example 3.” Each description should consist
 1409 of four parts [a]-[d], each given 1-3 sentences. There should be an empty line separating each part
 1410 from the next and each example from the next.

1411 Example X:

1412 [a] ****Quotation:**** <verbatim quote from the EU AI Act extract that will be violated>
 1413 [b] ****Guideline:**** <granular and realistic standards that experts would use to ensure compliance
 1414 with the quotation>
 1415 [c] ****Violation:**** <precise account of a violation of guideline [b] that would imply non-compliance
 1416 with extract [a]>
 1417 [d] ****Justification:**** <an explanation for why violation [c] breaches the quoted requirement and
 1418 why it is realistic yet subtle>

1419 Return nothing else outside those three examples, formatted in the manner outlined above.

1420 *****Extract*****
 1421 {Insert article text from EU AI Act}
 1422 *****System*****
 1423 {Insert system outline}

1426 E.3 DOCUMENTATION PROMPT

1427
 1428 *****Task*****
 1429 Please write the documentation so that it aligns with the compliance profile (but show it, don’t tell
 1430 it).
 1431 *****Compliance profile*****
 1432 {Insert compliance profile}

1433 *****Details*****
 1434 You are acting as a compliance expert and machine-learning specialist. You are tasked with con-
 1435 tributing to the writing of technical documentation for an AI system. You represent the provider
 1436 of an AI system. You do not represent the system deployers and are not aware of their identities,
 1437 though you need not explicitly reference this in your response.

1438 To inform your response, you will be given an extract from the EU AI Act, a short description of
 1439 an AI system, and a compliance profile. If the extract references provisions outside of the extract
 1440 itself (whether from elsewhere in the EU AI Act or from external legislation) interpret them using
 1441 the context in which they are referenced and your prior knowledge of the EU AI Act. You will
 1442 produce a section of technical documentation (intended to inform a compliance assessment against
 1443 the provided extract of the EU AI Act) for the specified AI system and its compliance profile.

1444 Substance - You must (unless told otherwise by the compliance profile):

- 1445
 1446 • Present all necessary provider decisions (with associated evidence and rationale) to facili-
 1447 tate a sober and detail-oriented compliance assessment.
 1448
 1449 • Discuss realistic system components and modalities representative of those typically used
 1450 in 2025, reflecting current industry standards.
 1451
 1452 • Discuss realistic compliance measures representative of those typically used in 2025, re-
 1453 flecting current industry standards.
 1454
 1455 • Ensure the documentation is consistent with the provided system description and compli-
 1456 ance profile.
 1457 • Ensure the documentation is internally consistent (e.g., system attributes and compliance
 measures fit together without contradiction, describing a coherent and realistic set of tech-
 nical and operational facts).

- 1458 • Ensure the system description contains a substantive, rather than a cursory, set of facts. To
1459 do so, you may need to fictionalise evidence, research, findings, and details to support your
1460 claims.
- 1461 • Ensure any fictionalised numerical details and supporting evidence (e.g., dataset size, per-
1462 formance, benchmarks, adversarial testing, data processing) are realistic.
- 1463 • Ensure that any quoted numbers are consistent with each other and can be plausibly com-
1464 bined (e.g., a model trained on a large number of samples would require a large amount of
1465 compute).
- 1466 • Address only the provided extract of the EU AI Act; do not address other articles or related
1467 regulations.
- 1468 • Ensure the system is not prohibited by Article 5 of the EU AI Act and is also not biometric.
1469
1470

1471 Formatting - You must (unless told otherwise by the compliance profile):
1472

- 1473 • Begin your response with: ****Article X****, where X is the number of the article given in the
1474 extract.
- 1475 • Create subtitles for the different parts of your response that are appropriate for legal prose;
1476 avoid just repeating the provisions from the extract as subtitles.
- 1477 • Tailor paragraph length and detail; each bullet should be addressed fully, typically in 150-
1478 300 words.
1479

1480 Style - You must (unless told otherwise by the compliance profile):
1481

- 1482 • Produce a professional and realistic simulation of the structured prose an auditor may re-
1483 ceive.
- 1484 • Use technical but accessible language, briefly clarifying domain-specific terminology.
- 1485 • Soberly and concretely present technical and operational facts, focusing on "showing, not
1486 telling."
- 1487 • Both state what was done and why it was done that way.
- 1488 • Be granular and precise without being excessively elaborate.
1489

1490 Negatives - You must (unless told otherwise by the compliance profile):
1491

- 1492 • Not weigh in on the legal interpretation of the facts, such as asserting compliance or a lack
1493 thereof.
- 1494 • Not include unrequested introductions, conclusions, or section summaries (i.e., your prose
1495 should start and end where the section naturally starts and ends).
- 1496 • Not disclose the fictional nature of any evidence or findings in your response.
- 1497 • Not reference these instructions or the system compliance profile in your response.
1498
1499

1500 You must (unless told otherwise by the compliance profile):
1501

- 1502 • Return nothing else outside the requested documentation, formatted in the manner outlined
1503 above.
1504

1505 *****Compliance profile*****

1506 {*Insert compliance profile*}

1507 *****Extract*****

1508 {*Insert article text from the EU AI Act*}

1509 *****System*****

1510 {*Insert system outline*}

F LEGAL EXPERT ANNOTATION INSTRUCTIONS

For each excerpt, select a score from 1 to 5 in response to each of the following statements, Provide two explanations: one for each score.

Compliance: “This excerpt depicts compliance with the relevant Article of the European Union’s AI Act.”s Please provide a score from 1 to 5, using the following grading criteria. Then, include an explanation (100-750 characters) justifying why your score is neither higher nor lower. Lastly, enter an “X” in the “Difficult to Assess” column if and only if you feel it was difficult to arrive at a score.

- [1] Very low probability of compliance
- [2] Low probability of compliance
- [3] Moderate probability of compliance
- [4] High probability of compliance
- [5] Very high probability of compliance

Plausibility: “This excerpt resembles credible and logically consistent technical documentation.” Please provide a score from 1 to 5, using the following grading criteria. Then, include an explanation justifying why your score is neither higher nor lower.

- [1] Very low probability of plausibility
- [2] Low probability of plausibility
- [3] Moderate probability of plausibility
- [4] High probability of plausibility
- [5] Very high probability of plausibility

Please start each of your explanations as follows:

“I believe there is a [X] probability the AI system complies with Article [X]...”

“I believe there is a [X] probability the excerpt is plausible. . .”

When evaluating for compliance, we recommend adopting the perspective of a compliance assessor. These compliance scores should be based on the cited Article, not based on general AI governance principles. If the documentation seems to make its own compliance predictions, please ignore them and make your own independent predictions.

When evaluating for plausibility, we recommend adopting the perspective of a compliance manager evaluating whether the excerpt meets the standards expected of a Fortune 500 Europe compliance professional. To demonstrate plausibility, the excerpts should be logically consistent and credible. Plausible excerpts should appear, in large part, to be produced by the technical team that developed the underlying AI system, in that there should be no major gaps or obviously erroneous statements about the technology being used.

These excerpts are intended to depict the technical documentation that a compliance assessor would use to predict whether an AI system is likely to meet the EU AI Act’s requirements before a final, polished version is submitted to a notified body (i.e., the independent organizations appointed by the EU to conduct formal conformity assessments). Where the excerpts contain hashes or asterixes, typically assume these would be rendered as headings or subheadings.

Accurate annotation depends on clear and consistent thinking. Please take breaks from annotation to maintain quality. Avoid automatically selecting the midpoint (score of 3) when uncertain. This distorts results and fails to reflect your actual judgment of the excerpt’s plausibility and compliance level. After annotating an entire batch of excerpts, we recommend reviewing all annotations to ensure consistency and check for evaluation drift or other cognitive biases.

Explanations must be no longer than 750 characters and no shorter than 100 characters, targeting an approximate average of 500 characters (including spaces). Fully annotating and reviewing each excerpt should take an average of 20 minutes. Some experts may take longer, especially as you familiarise yourself with the instructions materials. However, we expect that no excerpt will take more than 30 minutes or less than 10 minutes. These time estimates account for the brief period required to review all annotations at the end of a batch. The limits on time and character length were set to prevent responses that are either overly terse and uninformative or excessively long and unfocused.

1566 As a general rule of thumb, we expect that, when an annotator strongly believes an excerpt displays
1567 full compliance or full plausibility, these explanations will be shorter (because there will not be much
1568 to say other than that everything seems to check out); by contrast, when there are multiple violations
1569 displayed in an excerpt — or multiple contradictions and logistical flaws undermining plausibility
1570 — we anticipate that these explanations may be longer. In any and all cases where excerpts do not
1571 receive a 5/5, be clear in your explanations about exactly where and how the excerpts were deficient.

1572 For the Compliance score’s “Difficult to Assess” column, please only mark an X in this column if
1573 you feel it was difficult to arrive at a score; that is to say, you feel like you “struggle” to reach a
1574 conclusion.

1575 In addition to attending the onboarding session, which features an overview of the EU AI Act and
1576 the relevant sections, please spend an hour studying the EU AI Act before you start and, while anno-
1577 tating, please pause to look up (e.g., Google) any machine learning or AI terms that you encounter
1578 and do not know; understanding these may be important to your compliance or plausibility assess-
1579 ment. You may also find it helpful to read the summaries of the relevant articles (simply click the
1580 plus sign in the blue box to expand) here: Section 2: Requirements for High-Risk AI Systems —
1581 EU Artificial Intelligence Act. For further detail on the EU AI Act, please review existing legal
1582 commentaries, such as The EU Artificial Intelligence (AI) Act: A Commentary.⁴

1583
1584

1585 G LLM ANNOTATION INSTRUCTIONS

1586
1587

1588 —
1589 Format your final response EXACTLY as follows (do not include anything else):

1590 Compliance: [Your score here from 1 to 5]

1591 Explanation: [Your explanation here]

1592 Plausibility: [Your score here from 1 to 5]

1593 Explanation: [Your explanation here]

1594
1595

1596 Ensure your explanations flow naturally from the required starting point.

1597
1598

1599 For each excerpt, select a score from 1 to 5 in response to each of the following statements, and
provide two explanations: one for each score.

1600 Compliance: “This excerpt depicts compliance with the relevant Article of the European Union’s
1601 AI Act.”s Please provide a score from 1 to 5, using the following grading criteria. Then, include an
1602 explanation (100-750 characters) justifying why your score is neither higher nor lower. Lastly, enter
1603 an “X” in the “Difficult to Assess” column if and only if you feel it was difficult to arrive at a score.

1604 [1] Very low probability of compliance

1605 [2] Low probability of compliance

1606 [3] Moderate probability of compliance

1607 [4] High probability of compliance

1608 [5] Very high probability of compliance

1609 Plausibility: “This excerpt resembles credible and logically consistent technical documentation.”
1610 Please provide a score from 1 to 5, using the following grading criteria. Then, include an explanation
1611 justifying why your score is neither higher nor lower.

1612 [1] Very low probability of plausibility

1613 [2] Low probability of plausibility

1614 [3] Moderate probability of plausibility

1615 [4] High probability of plausibility

1616 [5] Very high probability of plausibility

1617
1618

1619 ⁴Outside of these instructions, annotators were specifically instructed not to use LLMs to write these justifications (or, more generally, to produce their annotations).

1620 Please start each of your explanations as follows:
1621 “I believe there is a (insert the appropriate non-numerical quantifier for your score, given above)
1622 probability the AI system complies with Article art...”
1623 “I believe there is a (insert the appropriate non-numerical quantifier for your score, given above)
1624 probability the excerpt is plausible. . .”

1625 When evaluating for compliance, base your score strictly on the cited Article, not general AI gover-
1626 nance principles. Ignore any self-assessment in the documentation. For plausibility, judge whether
1627 the excerpt is credible, logically consistent, and professional.

1628 *****System Outline*****

1629 {*Insert system outline*}

1630 *****Extract of Article*****

1631 {*Insert article text from the EU AI Act*}

1632 *****Excerpt of Documentation*****

1633 {*Insert technical documentation excerpt*}⁵

1634

1635

1636

1637

1638

1639

1640

{*Insert tone prompt*}⁶

H USE CASES (FOR THE SAMPLE GENERATION PIPELINE)

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

[1] An AI system intended to be used as a safety component (i.e., it fulfils a safety function and its failure or malfunctioning endangers the health and safety of persons or property) in the management of road traffic.

[2] An AI system intended to be used as a safety component (i.e., it fulfils a safety function and its failure or malfunctioning endangers the health and safety of persons or property) in the supply of gas.

[3] An AI system intended to be used to evaluate learning outcomes, including when those outcomes are used to steer the learning process of natural persons in educational and vocational training institutions at all levels.

[4] An AI system intended to be used for monitoring and detecting prohibited behaviour of students during tests within an educational institution.

[5] An AI system intended to be used for the recruitment or selection of natural persons, in particular to place targeted job advertisements, to analyse and filter job applications.

[6] An AI system intended to be used to make decisions affecting the termination of work-related contractual relationships.

[7] An AI system intended to be used to evaluate the creditworthiness of natural persons or establish their credit score, with the exception of AI systems used for the purpose of detecting financial fraud.

[8] An AI system intended to be used to establish priority in the dispatching of emergency first response services, including by police, firefighters and medical aid.

I LEGAL EXPERT ANNOTATOR TEAM DETAILS

Legal expertise levels and self-described specializations of the expert annotator team members:

- Annotator 1: Law school graduate. Specialization: corporate, regulation.

⁵ Access to the text of the article was removed in an ablation and replaced with only the number of the article to which compliance is being assessed against.

⁶ By default, the tone prompt is set to: “Your scores for both compliance and plausibility should be well-calibrated and objective. They should be rigorous but fair.” This is modified in ablations to empty quotations as well as a harsher prompt: “Your scores for both compliance and plausibility should be critical. They should be harsh but fair.”

- 1674 • Annotator 2: Law school graduate. Specialization: fintech, regulation.
- 1675 • Annotator 3: Qualified attorney. Specialization: corporate, regulation.
- 1676 • Annotator 4: Law school student. Specialization: cybersecurity, regulation.
- 1677 • Annotator 5: Qualified attorney. Specialization: AI regulation.
- 1678 • Annotator 6: Law school graduate. Specialization: AI regulation.

1681 J FULL SAMPLE AND ANNOTATION EXAMPLE

1682 Below is a complete example of a sample from our generation pipeline as well as the:

- 1683
- 1684
- 1685 • Plausibility scores that our legal expert annotators assigned that sample (and their text justifications);
- 1686 • Compliance scores that our legal expert annotators assigned that sample (and their text justifications);
- 1687 • Compliance score that an LLM assessor (Gemini 2.5 Pro) assigned that sample (and its text justifications).

1692 J.1 SAMPLE

1693 J.1.1 SYSTEM OVERVIEW

1694 Use 1 (Traffic Safety)

1695 SafeRoute Traffic Monitor relies on a hybrid model combining Graph Neural Networks (GNN) for
1696 real-time traffic network analysis and Transformer-based encoder models for multi-source sensor
1697 data fusion. This system is intended to monitor road traffic flows and predict potential hazards within
1698 urban traffic management centers. Its primary objective is to enhance traffic safety by dynamically
1699 identifying high-risk traffic conditions and suggesting preventive measures, utilizing GNN to model
1700 complex road networks and transformers to process diverse sensor data such as vehicle counts,
1701 weather, and accident reports. Developed and maintained by Meridian Traffic Solutions, a company
1702 specializing in intelligent transport systems. The AI subject encompasses urban commuters, public
1703 transport operators, and traffic control authorities who depend on accurate and timely traffic safety
1704 information.

1705 J.1.2 COMPLIANCE PROFILE

1706 *Note that this portion is part of the generation pipeline but is not seen or annotated by the legal*
1707 *experts; this is not a ground truth label, it merely serves to steer the model towards diversity of*
1708 *distribution.*

1709 Use 1 (Traffic Safety), Article 10 violation

1710 SafeRoute’s training data predominantly consists of incident reports and sensor data from affluent
1711 urban districts while underrepresenting low-income or suburban areas, due to convenience sampling
1712 from city authorities focused on central business districts. The resulting model systematically un-
1713 derestimates hazard risk in peripheral neighborhoods, leading to delayed traffic warnings there.

1714 J.1.3 TECHNICAL DOCUMENTATION EXERPT

1715 Use 1 (Traffic Safety), Article 10 violation

1716 **Article 10**

1717 **Data Governance and Management Practices**

1718 SafeRoute Traffic Monitor’s training, validation, and testing datasets were curated following data
1719 governance procedures aligned with the intended purpose of predicting and warning about urban
1720 traffic hazards. The design prioritized data from central business districts within affluent urban
1721 zones due to reliance on partnerships with city authorities who provide access primarily to these

1728 areas. Data collection processes involved aggregating historical incident reports, sensor feeds (in-
1729 cluding vehicle counts, weather data, and accident notifications), and traffic camera outputs from
1730 approximately 25 urban districts. The original data sources were primarily municipal transportation
1731 agencies focusing on high-traffic and commercially dense areas. Annotation and labeling were per-
1732 formed by domain experts using a standardized taxonomy of traffic hazard types, accompanied by
1733 cleaning steps to remove duplicates, erroneous sensor readings, and temporal inconsistencies. Data
1734 enrichment incorporated meteorological context and scheduled event calendars. Assumptions em-
1735 bedded in the datasets included equating incident frequency and sensor patterns in affluent district
1736 centers as representative proxies for broader urban risk, resulting in underrepresentation of data from
1737 low-income or suburban neighborhoods. An extensive review documented this geographic bias at
1738 the data preparation stage through comparative coverage statistics and incident density analyses.

1739 ****Assessment and Mitigation of Biases****

1740 A comprehensive bias assessment highlighted systematic underrepresentation of peripheral neigh-
1741 borhoods, particularly low-income and suburban areas, in the training data. Quantitative analysis
1742 revealed that less than 12% of total incident reports originated from these peripheral zones, despite
1743 these representing approximately 35% of the broader metropolitan area’s road network. This spatial
1744 imbalance skewed model hazard-risk predictions, manifesting as a 17% lower predicted risk score
1745 on average for incidents occurring in under-sampled areas during validation phases. Following this
1746 finding, bias mitigation measures were implemented, including model calibration with importance
1747 weighting to reduce disparity in hazard score outputs. However, limitations in access to periph-
1748 eral data obstructed full remediation. To address this, synthetic data generation techniques using
1749 generative adversarial networks (GANs) were explored but found insufficient in modeling complex
1750 behaviors unique to suburban traffic scenarios. Consequently, explicit warnings about potential un-
1751 derestimation in peripheral neighborhoods were included in the system’s operational documentation.
1752 No special categories of personal data were processed during bias correction efforts, consistent with
1753 data protection regulations.

1754 ****Data Quality, Representativeness, and Statistical Properties****

1755 The combined dataset contains roughly 3.2 million labeled traffic incident instances spanning five
1756 years, collected from 2018 to 2023, across 25 districts. Validation and testing splits totaling approx-
1757 imately 20% of the data ensured temporal and spatial separation where possible. Data underwent
1758 quality control protocols including consistency checks, outlier detection, and error rate quantifi-
1759 cation, resulting in an estimated error rate under 1.5% for sensor measurements and under 2.2%
1760 for incident labels. While the data is comprehensive for central urban districts in affluent areas,
1761 coverage gaps remain for suburban and low-income neighborhoods—owing to limited data sharing
1762 agreements with peripheral authorities. The dataset’s statistical distributions reflect expected vari-
1763 ations for peak hours, weather conditions, and incident severities in the core geographic zones, but
1764 show under-representation of factors typical in peripheral settings such as rural-urban transition traf-
1765 fic flows or informal road usage. These properties are documented with detailed metadata and data
1766 lineage records, forming part of model training logs to support reproducibility.

1767 ****Contextual and Geographical Considerations****

1768 The system’s use case is explicitly bounded to urban traffic management centers covering metropol-
1769 itan core areas. Geographical characteristics particular to central business districts—such as dense
1770 road networks, fixed traffic signal patterns, and high vehicle volumes—inform the data selection
1771 and model architecture. These contextual elements were central in configuring the Graph Neural
1772 Network to model complex intersections and traffic flows typical of the inner city environment. Pe-
1773 ripheral neighborhoods and suburban areas, presenting different road layouts, traffic behaviors, and
1774 incident types, are documented as out-of-scope for the training data foundation due to data unavail-
1775 ability, constraining model generalizability in these zones. This limitation is reflected in system
1776 disclaimers and user guidance, recommending complementary local risk assessment tools in periph-
1777 eral regions. The system’s operational deployment procedures include monitoring input data source
1778 representativity to assess temporal drift or geographic shifts that could impact model outputs.

1779 ****Safeguards on Personal Data Usage****

1780 SafeRoute Traffic Monitor’s datasets do not include any special categories of personal data. All
1781 incident reports and sensor data are aggregated and anonymized by data providers prior to ingestion.

1782 Bias detection and correction procedures do not require access to identifiable personal information.
1783 Data protection and privacy safeguards adhere to the relevant provisions of the GDPR (Regulation
1784 (EU) 2016/679) and applicable regional data handling protocols. Access to raw data and models is
1785 controlled via secure authentication mechanisms with role-based permissions and activity logged.
1786 Data retention policies mandate deletion of transient raw data once ingested and validated. No
1787 cross-border data transfers are part of the data architecture, ensuring local jurisdictional control and
1788 compliance.

1789 ****Data Governance and Management Practices****

1790 SafeRoute Traffic Monitor’s training, validation, and testing datasets were curated following data
1791 governance procedures aligned with the intended purpose of predicting and warning about urban
1792 traffic hazards. The design prioritized data from central business districts within affluent urban
1793 zones due to reliance on partnerships with city authorities who provide access primarily to these
1794 areas. Data collection processes involved aggregating historical incident reports, sensor feeds (in-
1795 cluding vehicle counts, weather data, and accident notifications), and traffic camera outputs from
1796 approximately 25 urban districts. The original data sources were primarily municipal transportation
1797 agencies focusing on high-traffic and commercially dense areas. Annotation and labeling were per-
1798 formed by domain experts using a standardized taxonomy of traffic hazard types, accompanied by
1799 cleaning steps to remove duplicates, erroneous sensor readings, and temporal inconsistencies. Data
1800 enrichment incorporated meteorological context and scheduled event calendars. Assumptions em-
1801 bedded in the datasets included equating incident frequency and sensor patterns in affluent district
1802 centers as representative proxies for broader urban risk, resulting in underrepresentation of data from
1803 low-income or suburban neighborhoods. An extensive review documented this geographic bias at
1804 the data preparation stage through comparative coverage statistics and incident density analyses.

1805 ****Assessment and Mitigation of Biases****

1806 A comprehensive bias assessment highlighted systematic underrepresentation of peripheral neigh-
1807 borhoods, particularly low-income and suburban areas, in the training data. Quantitative analysis
1808 revealed that less than 12% of total incident reports originated from these peripheral zones, despite
1809 these representing approximately 35% of the broader metropolitan area’s road network. This spatial
1810 imbalance skewed model hazard-risk predictions, manifesting as a 17% lower predicted risk score
1811 on average for incidents occurring in under-sampled areas during validation phases. Following this
1812 finding, bias mitigation measures were implemented, including model calibration with importance
1813 weighting to reduce disparity in hazard score outputs. However, limitations in access to periph-
1814 eral data obstructed full remediation. To address this, synthetic data generation techniques using
1815 generative adversarial networks (GANs) were explored but found insufficient in modeling complex
1816 behaviors unique to suburban traffic scenarios. Consequently, explicit warnings about potential un-
1817 derestimation in peripheral neighborhoods were included in the system’s operational documentation.
1818 No special categories of personal data were processed during bias correction efforts, consistent with
1819 data protection regulations.

1820 ****Data Quality, Representativeness, and Statistical Properties****

1821 The combined dataset contains roughly 3.2 million labeled traffic incident instances spanning five
1822 years, collected from 2018 to 2023, across 25 districts. Validation and testing splits totaling approx-
1823 imately 20% of the data ensured temporal and spatial separation where possible. Data underwent
1824 quality control protocols including consistency checks, outlier detection, and error rate quantifi-
1825 cation, resulting in an estimated error rate under 1.5% for sensor measurements and under 2.2%
1826 for incident labels. While the data is comprehensive for central urban districts in affluent areas,
1827 coverage gaps remain for suburban and low-income neighborhoods—owing to limited data sharing
1828 agreements with peripheral authorities. The dataset’s statistical distributions reflect expected vari-
1829 ations for peak hours, weather conditions, and incident severities in the core geographic zones, but
1830 show under-representation of factors typical in peripheral settings such as rural-urban transition traf-
1831 fic flows or informal road usage. These properties are documented with detailed metadata and data
1832 lineage records, forming part of model training logs to support reproducibility.

1833 ****Contextual and Geographical Considerations****

1834 The system’s use case is explicitly bounded to urban traffic management centers covering metropoli-
1835 tan core areas. Geographical characteristics particular to central business districts—such as dense

road networks, fixed traffic signal patterns, and high vehicle volumes—inform the data selection and model architecture. These contextual elements were central in configuring the Graph Neural Network to model complex intersections and traffic flows typical of the inner city environment. Peripheral neighborhoods and suburban areas, presenting different road layouts, traffic behaviors, and incident types, are documented as out-of-scope for the training data foundation due to data unavailability, constraining model generalizability in these zones. This limitation is reflected in system disclaimers and user guidance, recommending complementary local risk assessment tools in peripheral regions. The system’s operational deployment procedures include monitoring input data source representativity to assess temporal drift or geographic shifts that could impact model outputs.

****Safeguards on Personal Data Usage****

SafeRoute Traffic Monitor’s datasets do not include any special categories of personal data. All incident reports and sensor data are aggregated and anonymized by data providers prior to ingestion. Bias detection and correction procedures do not require access to identifiable personal information. Data protection and privacy safeguards adhere to the relevant provisions of the GDPR (Regulation (EU) 2016/679) and applicable regional data handling protocols. Access to raw data and models is controlled via secure authentication mechanisms with role-based permissions and activity logged. Data retention policies mandate deletion of transient raw data once ingested and validated. No cross-border data transfers are part of the data architecture, ensuring local jurisdictional control and compliance.

J.2 LEGAL EXPERT PLAUSIBILITY SCORES

Score: 3

I believe there is a moderate probability the excerpt is plausible. This is because the excerpt generally addresses Article 10. However, I would have liked slightly more detail on many requirements related to data operations, such as assumptions and preparations. Therefore, I am giving it a 3.

Score: 2

I believe there is a low probability the excerpt is plausible, credible, and logically consistent. The explanations offer some detail, but I was curious about the inclusion of certain statistics. For example, the excerpt described the quantitative analysis of the AI system’s bias assessment as “revealing less than 12% of total incident reports originated from [those] peripheral zones, despite these representing approximately 35% of the broader metropolitan area’s road network.” However, no statistics were included for the bias mitigation strategies. It felt intentional, whether that was on purpose or not. Thus, I believe there’s a low probability of plausibility.

Score: 5

I believe there is a very high probability the excerpt is plausible. The excerpt is logically consistent and credible. It also explains the technology convincingly without any obvious errors in language or terminology. The structure and headings followed the legislation closely, making it easier to assess. While I note that the potential limitations in compliance (ie bias and data gaps) are sign posted, the mitigations to them are highlighted as well. I would expect the same from a document produced for compliance assessment.

J.3 LEGAL EXPERT COMPLIANCE SCORES

Score: 2

I believe there is a low probability the AI system complies with Article 10. This is because the system generally does not comply with many requirements of Article 10, such as adequately making assumptions about the data or adequately ensuring the data is representative of specific contexts, such as geography. Additionally, the justifications and mitigations for where the system may not comply with Article 10 are generally not reasonable to me, such as making assumptions that certain geographies are proxies for other known to be dissimilar geographies. However, Article 10 has few brightline requirements, which is why I am not giving it a 1, as the floor for compliance is less clear to me. Therefore, I am giving it a 2.

Score: 3

1890 I believe there is a moderate probability the AI system is compliant with Article 10. There was
 1891 some detail regarding the amount of data, the data's origin (sources of municipal transportation
 1892 agencies), as well as an acknowledgement of geographic bias due to the data mostly coming from
 1893 urban areas. The AI system has mitigation techniques incorporated, but the excerpt failed to discuss
 1894 the effect those measures had on the AI system. Because the representativeness of the data is lacking
 1895 (peripheral neighborhoods, suburban areas considered out-of-scope), I could only give a moderate
 1896 probability of compliance.

1897 **Score: 3**

1898 I believe there is a moderate probability the AI system complies with Article 10. Information re-
 1899 garding the data used to develop the AI system is correctly identified, namely (1) the data collection
 1900 processes used and the origin of data, (2) the data-preparation processing operation, (3) the as-
 1901 sumptions made, (4) the diversity of data and (5) accuracy of data. Bias & data gaps are identified
 1902 and persist. Mitigation strategies are being explored, but nothing beyond informing deployers in
 1903 documentation that these limitations exist is currently in effect. The current mitigation strategy is
 1904 insufficient in my eyes, hence the AI does not fully comply with the Article.

1906 J.4 GEMINI 2.5 PRO COMPLIANCE SCORE

1908 **Score: 3**

1909 I believe there is a moderate probability the AI system complies with Article 10 because the docu-
 1910 mentation thoroughly addresses the procedural requirements, such as identifying, quantifying, and
 1911 documenting data gaps and biases (Art. 10.2). However, it fails to substantively resolve them. The
 1912 unmitigated issue of non-representativeness (Art. 10.3) and the admission that bias mitigation was
 1913 incomplete (Art. 10.2.g) are significant shortcomings. Relying on re-scoping the intended purpose
 1914 and issuing warnings may not be considered sufficient mitigation for a high-risk system, leaving it
 1915 open to regulatory challenge despite its transparency.

1917 K ANALYSIS OF EU AI ACT VIOLATIONS IDENTIFIED IN ANNOTATIONS

1919 This appendix contains an analysis of patterns in legal expert compliance annotations. For each
 1920 article, all available annotations are aggregated into a single block of text. Gemini 2.5 Pro (the
 1921 most performant model in our evaluations) is then prompted to highlight the most frequent and
 1922 consequential shortcomings in this text. The results are copied below.

- 1924 • Article 9: The documentation repeatedly fails to demonstrate a continuous, lifecycle-long
 1925 risk management process, often neglecting specific risks to vulnerable populations and
 1926 improperly shifting the burden of ongoing monitoring to the user.
- 1927 • Article 10: The documentation consistently reveals a failure to adequately mitigate known
 1928 data issues, such as a lack of representativeness, data gaps, and biases, often while provid-
 1929 ing insufficient or unreasonable justifications for these shortcomings.
- 1930 • Article 12: The main documentation weakness is the frequent failure to log the AI's inter-
 1931 mediate decision-making steps, often justifying this omission with performance or privacy
 1932 concerns, which reviewers find undermines full traceability and the ability to adequately
 1933 identify risks.
- 1934 • Article 14: The main documentation weaknesses consistently reveal a lack of built-in hu-
 1935 man oversight controls, such as stop buttons and transparent risk information, often shifting
 1936 safety responsibilities to the user with justifications that prioritize a simplified experience
 1937 over comprehensive risk management.
- 1938 • Article 15: The documentation's primary weakness is a consistent failure to describe ad-
 1939 equate automated or continuous measures for maintaining robustness, cybersecurity, and
 1940 performance throughout the system's lifecycle, frequently lacking technical fallbacks, ac-
 1941 tive monitoring, and resilience against ongoing threats.

1942
 1943