BigBear.ai Holdings, Inc.
6811 Benjamin Franklin Drive
Suite 200
Columbia, Maryland 21046

October 18, 2023

Filed Electronically

Suzanne V. Wilson
General Counsel and Associate Register of Copyrights
Maria Strong
Associate Register of Copyrights and Director of Policy and International Affairs
Library of Congress, Copyright Office
101 Independence Avenue, SE
Washington, DC 20559

Dear Ms. Strong and Ms. Wilson:

BigBear.ai Holdings, Inc. (BigBear.ai) appreciates the opportunity to submit the following comments to the request by the U.S. Copyright Office (the "Office") concerning copyright law and policy issues raised by artificial intelligence ("AI").

BigBear.ai recognizes and commends the Office's focus on striking the correct balance between promoting innovation by making training data available so that AI development may thrive and protecting the rights of copyright owners. As one of the few companies supporting the federal government in pioneering AI orchestration, BigBear.ai is dedicated to advancing trustworthy AI as it responsibly innovates and advances AI software and tools for the benefit of U.S. citizens and businesses. As the Office continues its work in this area, BigBear.ai strongly encourages the development of policies and regulations that support and encourage responsible innovation.

We appreciate your consideration of the following comments are available to answer any questions you may have.

Sincerely,

Mandy Long
CEO
BigBear.ai

**General Questions**

The Office has several general questions about generative AI in addition to the specific topics listed below. Commenters are encouraged to raise any positions or views that are not elicited by the more detailed questions further below.

1. **As described above, generative AI systems have the ability to produce material that would be copyrightable if it were created by a human author. What are your views on the potential benefits and risks of this technology? How is the use of this technology currently affecting or likely to affect creators, copyright owners, technology developers, researchers, and the public?**

Generative AI allows creation of myriad of works on a scale that could never be accomplished by humans alone. It also allows humans to create things (like artwork) that they might not be able to create unaided. There are also a number of risks: because generative AI necessarily requires training material, most of which is subject to copyright, it raises the question as to how or if the authors of this material should be consulted or compensated. This raises the existential question of what does it mean to create? In an age where art, music and culture often utilize appropriation, how is gen-AI different? Consumers may want (and have the right) to know whether a work was generated by AI, just as they have the right to know whether they are communicating with a human or a chatbot.

2. **Does the increasing use or distribution of AI-generated material raise any unique issues for your sector or industry as compared to other copyright stakeholders?**

Our issues are similar to those of others.

3. **Please identify any papers or studies that you believe are relevant to this Notice. These may address, for example, the economic effects of generative AI on the creative industries or how different licensing regimes do or could operate to remunerate copyright owners and/or creators for the use of their works in training AI models. The Office requests that commenters provide a hyperlink to the identified papers.**

N/A

4. **Are there any statutory or regulatory approaches that have been adopted or are under consideration in other countries that relate to copyright and AI that should be considered or avoided in the United States? [40] How important a factor is international consistency in this area across borders?**

Yes. Directive (EU) 2019/790 of 17 April 2019 on copyright and related rights in the Digital Single Market (Copyright Directive) introduced a text and data mining exception to copyright infringement under EU law.

Under the Directive, text and data mining in connection with research, educational, and cultural purposes is permitted without any caveats.  Text and data mining for any other use (including commercial use) is permitted except where the copyright owner "opts out."

In the UK, any content that can be accessed legally can be mined for non-commercial research purposes, including content subscribed to by a library or content that can be openly accessed online.  Draft legislation that would have extended TDM for commercial and other non-research purposes was introduced but was withdrawn after complaints from rightsholders.   BigBear.ai supports approaches to legislation that make content more readily available for AI innovation.

International consistency in copyright is desirable, particularly because the copyright issues raised by AI are likely to cross borders.

5.  **Is new legislation warranted to address copyright or related issues with generative AI? If so, what should it entail? Specific proposals and legislative text are not necessary, but the Office welcomes any proposals or text for review.**

Yes.  Although several AI-related court cases are working their way through the judicial system, rulings are likely to be very fact-specific, as IP cases often are.  Legislation would be helpful, assuming it strikes the right balance between protecting copyright owners' rights and allowing for innovation in an area that has tremendous potential.

**Training**

*If your comment applies only to a specific subset of AI technologies, please make that clear.*

6.  **What kinds of copyright-protected training materials are used to train AI models, and how are those materials collected and curated?**

 AI developers may use various types of copyright-protected training materials to train AI models, including text, images, audio, and more. These materials are typically collected and curated through a combination of methods, which can vary depending on the type of data and the goals of the AI project. Here's an overview of the kinds of copyrighted training materials and the collection and curation processes:

1.  **Text Data**:
    - **Books**: Text from books, including both public domain and copyrighted works, may be used. Public domain texts are freely available, while copyrighted texts may require permission or licensing agreements.
    - **Websites and Articles**: Text scraped from websites, news articles, blogs, and other online sources can be valuable training data.

2. **Images and Multimedia**:
   - **Images**: Copyrighted images, such as photographs, illustrations, and artwork, are used in computer vision tasks. Some AI developers purchase licenses for stock photos, while others create their own datasets through image capture or acquisition.
   - **Videos**: Video data may include clips from movies, TV shows, or user-generated content. Obtaining rights or licenses for copyrighted video material is essential.

3. **Audio Data**:
   - **Music and Audio Clips**: Copyrighted music tracks and audio clips are used in applications like music generation and speech recognition. Licensing agreements may be required for such data.
   - **Spoken Language**: Recorded spoken language data, including podcasts and audio books, can be used for speech recognition and natural language processing. Proper permissions are necessary for copyrighted audio content.

4. **Databases and Proprietary Sources**:
   - AI developers may have access to proprietary databases or datasets that are protected by copyright. These may include medical records, financial data, or other domain-specific information.
   - Licensing agreements, partnerships, or collaborations with data providers may be established to gain access to such datasets.

5. **Data Labeling and Annotation**:
   - Human annotators often label or annotate data for supervised learning tasks. This may involve marking objects in images, transcribing audio, or tagging text. Care must be taken to ensure that annotators do not introduce copyrighted content without appropriate authorization.

6. **Data Augmentation and Generation**:
   - AI developers sometimes create synthetic data or augment existing datasets to increase training data diversity. This can involve the use of generative algorithms to create new content or variations.

7. **Data Curation and Preprocessing**:
   - Data curation involves selecting, cleaning, and formatting the training data to make it suitable for AI model training. This may include removing duplicates, handling missing values, and ensuring data quality.

**6.1. How or where do developers of AI models acquire the materials or datasets that their models are trained on? To what extent is training material first collected by third-party entities (such as academic researchers or private companies)?**

There are several open-source datasets. Here is a very short list:

**BigBear.ai**

**Natural Language Processing**

- Amazon Reviews: A collection of over 35 million reviews from the last 18 years. It includes things like ratings, reviews in plain text, and user information. It also contains complete product information for reference.
- Wikipedia Links Data: The full power of Wikipedia including four million articles containing 1.9 billion words. Your search options are varied and include both word and phrase searches as well as pieces of paragraphs.

**Sentiment Analysis**

- Standford Sentiment Treebank: Dataset containing sentiment notations for over 10,000 pieces of data from Rotten Tomatoes reviews rendered in HTML.
- Twitter US Airline Sentiment: Tweets collected about US Airlines with clear markers for positive, negative, and neutral tones, dated from 2015.

**Public Government Data**

- Data USA: A comprehensive overview of various sets of US public data in fun visualizations. It includes things like population, health, and jobs.
- EU Open Data Portal: Much like Data USA except with a concentration on countries belonging to the EU. It includes fields such as population, culture, energy, and health, among others.

**Finance and Economics**

- World Bank Open Data: Data concerning population demographics and key indicators for development.
- IMF Data: International Monetary Fund's collection of open data for things like debt rates, commodity pricing, international markets, and foreign exchange reserves.

**Facial Recognition**

- Labeled Faces In The Wild: Common dataset for facial recognition training. It includes 13,000 cropped faces plus a subset of people with two different pictures within the dataset.
- UMDFaces Dataset: Includes both still and video images. The dataset is annotated and features around 367,000 faces of over 8,000 subjects.

**Image Datasets**

- Imagenet: Dataset containing over 14 million images available for download in different formats. It also includes API integration and is organized according to the WordNet hierarchy.
- Google's Open Images: 9 million URLs to categorized public images in over 6,000 categories. Each image is licensed under creative commons.

**Health:**

- [Healthdata.gov](#): a resource from the US federal government providing data to improve health outcomes for the US population.
- [MIMIC Critical Care Database](#): Datasets for Computational Physiology with unidentified health data from 40,000 critical care patients (demographics, vital signs, medications, etc.)

**Media**

- [FiveThirtyEight Journalism](#): The numbers behind some of this journalism hub's stories. Useful for visualizations and data stories.
- [BuzzFeed Media](#): Open source data hub for everything in the realm of Buzzfeed. Everything their journalists used to produce the stories (the organization recommends reading the articles to get a better idea of how the data was used.

**Transportation**

- [US National Travel and Tourism Office](#): provides trustworthy datasets with big pictures of the tourism industry, including things like inbound and outbound travel and international visitor data.
- [Department of Transportation](#): datasets on each field that falls under the DOT including National Parks, driver registers, bridges and rail information, and port systems.

**Speech**

- [Flickr Audio Caption Corpus](#): 40,000 spoken captions from 8,000 images in a manageable size. It was initially designed for unsupervised speech pattern discovery.
- [Speech Commands Dataset](#): A continuously evolving collection of one second long utterances from thousands of different people. It's still receiving contributions and is useful for building basic voice interfaces.

**Sound**

- [FSD (Freesound)](#): A collection of every day sounds collected by contribution under an open source license.
- [Environmental Audio Datasets](#): It does contain some proprietary information, but a large portion is open source. It contains sound events tables and acoustic scenes tables.

**Dataset Aggregators**

- [OpenDataSoft](): 2600 data portals arranged in an interactive map formation or by country list. If you're looking for it, chances are, it's here.
- [Kaggle](): an online community of data scientists where users can work with and upload datasets. It's a community and a resource in one.
- [UCI Machine Learning Repository](): User contributed datasets in various levels of cleanliness. It's one of the originals, and you can download datasets without having to register anything.

**6.2. To what extent are copyrighted works licensed from copyright owners for use as training materials? To your knowledge, what licensing models are currently being offered and used?**

In many cases data is licensed under a strict business model license grant with indemnification from the licensee with 2$^{nd}$ derivative re-use rights thereof. The party responsible for granting data licenses, be it a vendor or a customer, must ensure that the agreement effectively outlines its ownership or other rights pertaining to the data through the following steps:

1. Obtain acknowledgements of their data rights from the licensee.
2. Include a carefully tailored definition of the data set being licensed within the agreement.

If the licensor is the rightful owner of the data, they should seek explicit acknowledgment from the licensee that the data provided under the agreement is solely and exclusively owned by the licensor. Moreover, for comprehensive data protection, the licensor should seek acknowledgments for the following:

- Include a carefully tailored definition of the data set being licensed within the agreement.

If the licensor is the rightful owner of the data, they should seek explicit acknowledgment from the licensee that the data provided under the agreement is solely and exclusively owned by the licensor. Moreover, for comprehensive data protection, the licensor should seek acknowledgments for the following:

- Emphasize that the licensor has invested substantial resources in gathering, assembling, and compiling the data, establishing it as a valuable asset.

- Specify the licensor's right to request additional fees for expanded data usage or alternative methods of utilization.

- Assert that the data is an original compilation protected by U.S. copyright laws and contains the trade secrets of the licensor.

  o In certain situations, a more narrowly defined licensed data set may be suitable. For example, in a data feed agreement where the licensee is restricted from generating any derived data, this approach accomplishes the following:

- Limits the scope of the licensed data.

- Allows the licensor to request additional fees for extended data usage or alternative methods of utilization.

  o Conversely, in a services agreement, the customer may opt for a broader data definition to encompass all data collected or received by the vendor, either directly or indirectly

from the customer, to facilitate the provision of services. This broader definition serves to prevent ambiguity regarding data ownership on the part of the service provider.

**6.3. To what extent is non-copyrighted material (such as public domain works) used for AI training? Alternatively, to what extent is training material created or commissioned by developers of AI models?**

The extent to which non-copyrighted material, such as public domain works, is used for AI training can vary depending on the specific AI model and the goals of the developers. Likewise, the use of training material created or commissioned by developers also varies. Here are some key points to consider:

1. Public Domain Works: Public domain works, which are not protected by copyright and are free for unrestricted use, are sometimes used as training data for AI models. For example, text from old books, historical documents, or publicly available scientific articles can be valuable sources of data for language models. The extent to which such data is used may depend on factors like the availability of relevant content and the objectives of the AI model.

2. Creative Commons and Open Access Content: Some AI developers may use content that is available under open licenses like Creative Commons, which allow for reuse and modification under certain conditions. Open access research papers and datasets are also commonly used as training material for AI models, especially in domains like natural language processing and computer vision.

3. Proprietary and Licensed Data: In many cases, AI developers may use proprietary or licensed datasets that they have acquired or created for training. This can include datasets containing text, images, audio, or other types of data. Developers often invest in these datasets to ensure they meet their specific needs and quality standards.

4. Data Curation: AI developers often curate and preprocess data to create suitable training sets. This process involves selecting, cleaning, and formatting data to make it usable for training AI models. Data curation can be a labor-intensive and resource-intensive task.

5. Synthetic Data: In some cases, developers may generate synthetic data to augment their training sets. This synthetic data is created using algorithms or simulations and can help address limitations in the availability of real-world data.

6. Domain-Specific Data: Depending on the AI model's intended application, developers may collect domain-specific data. For instance, developers working on medical AI may collect medical records, while those in autonomous driving may gather sensor data from vehicles.

The specific approach taken depends on factors like the model's purpose, ethical considerations, data availability, and the resources at the disposal of the developers. It's essential for AI developers to be mindful of legal and ethical considerations when selecting and using training data, particularly regarding copyright and privacy issues. Additionally, transparency and disclosure about the sources of training data are becoming increasingly important in the AI community.

**6.4.** **Are some or all training materials retained by developers of AI models after training is complete, and for what purpose(s)? Please describe any relevant storage and retention practices.**

The retention of training materials by developers of AI models after training is complete can vary based on several factors, including the type of data, project goals, legal requirements, and ethical considerations. Here are some common scenarios and purposes for retaining training materials, along with relevant storage and retention practices:

1. **Retention of Raw Data**:
   - **Purpose**: Developers may retain the raw training data to maintain a reference for model evaluation, future retraining, or to address potential issues or questions that may arise after deployment.
   - **Storage and Retention Practices**: Raw data can be stored in secure, organized databases or data repositories. Retention policies should consider data privacy, compliance with legal regulations, and the ongoing value of the data.

2. **Retention of Model Checkpoints and Parameters**:
   - **Purpose**: Developers often keep model checkpoints and parameter settings to track the model's training history, enable fine-tuning, or facilitate further research.
   - **Storage and Retention Practices**: Model checkpoints and parameters can be stored in version control systems, cloud-based repositories, or on-premises servers. Security measures should be in place to protect this information.

3. **Data Retention for Compliance and Auditing**:
   - **Purpose**: In regulated industries (e.g., healthcare or finance), data retention is necessary to comply with legal and regulatory requirements. Auditing purposes may also necessitate data retention.
   - **Storage and Retention Practices**: Sensitive data is often stored securely with access controls and encryption. Retention policies should align with relevant regulations, specifying the duration data must be retained.

4. **Continued Data Collection and Augmentation**:
   - **Purpose**: Developers may continue to collect new data or augment existing datasets to improve model performance or adapt to changing circumstances.
   - **Storage and Retention Practices**: New data should be integrated into the existing dataset or kept separately, depending on the use case. Data augmentation methods should be documented.

5. **Research and Benchmarking**:
   - **Purpose**: Retaining training materials can be valuable for research, benchmarking, and comparison studies with new models or techniques.
   - **Storage and Retention Practices**: Data can be organized in repositories, and detailed records of data sources and preprocessing steps should be maintained for research reproducibility.

6. **Ethical Considerations and Data Deletion**:
   - **Purpose**: In some cases, ethical considerations or data privacy regulations may require the deletion of certain training materials to protect individuals' privacy.
   - **Storage and Retention Practices**: Data deletion should be carried out securely, and records of data destruction should be maintained to demonstrate compliance.

7. **Data Security and Access Controls**:
   - **Purpose**: Retained data and materials must be secured to prevent unauthorized access or breaches.
   - **Storage and Retention Practices**: Implement robust security measures, access controls, encryption, and regular security audits to protect retained materials.

   The same practices come into play with Personal Identifiable Information and any sensitive information thereof.

7. **To the extent that it informs your views, please briefly describe your personal knowledge of the process by which AI models are trained. The Office is particularly interested in:**

7.1. **How are training materials used and/or reproduced when training an AI model? Please include your understanding of the nature and duration of any reproduction of works that occur during the training process, as well as your views on the extent to which these activities implicate the exclusive rights of copyright owners.**

It's important to note that while AI training may involve the reproduction of copyrighted materials, it typically doesn't involve the distribution, public display, or other activities that might infringe on copyright owners' exclusive rights. Moreover, the primary purpose of AI model training is to develop models with improved capabilities, not to create copies of copyrighted materials for commercial distribution. Whether training an AI model on copyrighted materials constitutes infringement is a complex and evolving issue of fact and law that is discussed elsewhere in this response.

Training materials are used and potentially reproduced during the training of AI models in ways that can implicate the exclusive rights of copyright owners. The nature and duration of reproduction can vary depending on the type of AI model and the training process. Here's an overview:

1. **Nature of Reproduction**:
   - **Data Ingestion**: During training, AI models ingest and process training data, which often involves making copies of the data in memory or storage for computational purposes. This copying is typically transient and occurs for the duration of the training process.
   - **Feature Extraction**: For text or image-based models, features are extracted from the training data. This may involve creating derivative representations of the data for model input, which is also a form of reproduction.
   - **Model Updates**: As the model learns, it updates its internal parameters. These updates can be seen as a form of transformation or adaptation of the original training data.

2. **Duration of Reproduction**:
   - **Temporary Reproduction**: The reproductions made during training are often temporary and exist solely to facilitate the learning process. They are not typically used for any permanent storage or distribution.
   - **Training Timeframe**: The duration of reproduction is tied to the training process, which varies depending on the complexity of the model and the volume of training data. Once training is complete, many of the intermediate copies are no longer retained.

3. **Implication for Copyright Owners**:
   - **Transformative Use**: In many cases, AI models are considered transformative in their use of copyrighted materials. They don't use the materials for their original purpose (e.g., reading a book) but rather for statistical analysis and pattern recognition.
   - **Fair Use**: In some jurisdictions, the "fair use" doctrine or similar exceptions to copyright law may apply to the use of copyrighted materials for AI training. This would depend on factors such as the purpose of use, the nature of the copyrighted work, the amount used, and the potential market impact.

### 7.2.  How are inferences gained from the training process stored or represented within an AI model?

Inferences gained from the training process are stored or represented within an AI model through its internal parameters and architecture. These inferences, often referred to as the knowledge or learned patterns, are encoded in the model's weights, biases, and network connections. Here's how this process works:

1. **Weight Matrices and Neural Network Architecture**:
   - In many AI models, especially neural networks, the model's architecture consists of layers of interconnected nodes (neurons) organized in a specific topology.
   - Each connection between nodes is associated with a weight value, and each node typically has a bias value. These weights and biases are the core components that store the learned information.

2. **Training Process**:
   - During the training process, the model learns to make predictions or classifications by adjusting its weights and biases in response to the training data.
   - The model computes predictions using the weighted sum of inputs, passes this through activation functions, and compares the results to the actual target values in the training data.
   - Errors are calculated, and optimization algorithms (e.g., gradient descent) update the weights and biases iteratively to minimize these errors.

3. **Encoded Knowledge**:
   - As the training progresses, the model's weights and biases are adjusted to capture patterns, correlations, and statistical relationships in the training data.
   - These weight values represent the model's knowledge or learned features. For example, in an image recognition model, certain weights may correspond to recognizing edges, textures, or specific object shapes.

4. **Representation and Storage**:
   - The knowledge or inferences gained from training are stored implicitly in the model's weight matrices and biases.
   - These weights and biases are stored as numerical values in memory, either on local devices, cloud servers, or specialized hardware accelerators.
   - The model's architecture, along with these numerical values, represents the trained knowledge.

5. **Inference Phase**:
   - During the inference phase (i.e., when the model is used to make predictions or classifications), input data is processed through the model's architecture using the stored weights and biases.
   - The model computes outputs based on the learned knowledge encoded in these parameters.
   - The inferences made during inference are a result of the knowledge acquired during training.

6. **Transfer Learning and Fine-Tuning**:
   - In some cases, pre-trained models with learned knowledge are used as starting points for new tasks (transfer learning). In these cases, the model's weights may be fine-tuned with a smaller dataset specific to the new task.

It's important to note that the exact representation and storage of knowledge within an AI model can vary based on the model type. For example, recurrent neural networks (RNNs), convolutional neural networks (CNNs), and transformers all have different architectures and mechanisms for storing knowledge. Additionally, knowledge representation can involve millions or even billions of parameters in large-scale models like BERT, LLAMA, GPT-3, and GPT-4 to name a few large language models.

The ability to efficiently represent and store learned knowledge is what enables AI models to make predictions, recognize patterns, and perform tasks based on their training. These stored inferences are the essence of the AI model's capabilities.

**7.3.** **Is it possible for an AI model to "unlearn" inferences it gained from training on a particular piece of training material? If so, is it economically feasible? In addition to retraining a model, are there other ways to "unlearn" inferences from training?**

The most foolproof way to unlearn would be to start from scratch and remove the training material from the dataset.  Training from scratch is usually the best mode and considered the baseline.  In most cases this is very cost prohibitive, however current research in the following areas are underway (this is in no way comprehensive):

1.  *Differential Privacy:* Proposed by Dwork et al where differential privacy offers probabilistic guarantees about the privacy of individual records in a database. The bounds changes in model parameters that may be induced by any single training point. While several efforts make it possible to learn with differential privacy, this guarantee is different from what most wish to provide.
2.  *Statistical Query Learning*: Cao et al. model unlearning in the statistical query learning framework. By doing so, they are able to unlearn a point in time when the learning algorithm queries data in an order decided prior to the start of learning. In this setting, it is possible to know exactly how individual training points contributed to model parameter updates. However, their approach is not generalizable and does not easily scale to more complex models These models are trained using adaptive statistical query algorithms which make queries that depend on all queries previously made.
3.  *Decremental Learning:* Ginart et al. consider the problem from a data-protection regulation standpoint. They present a formal definition of complete data erasure which can be relaxed into a distance-bounded definition. Deletion time complexity bounds are provided. They note that the deletion and privacy problems are orthogonal, which means deletion capability does not imply privacy nor vice versa.
4.  *SISA (Sharded, Isolated, Sliced, Aggregated)* Bourtoule et al training replicates the model being learned several times where each replica receives a disjoint shard (or subset) of the dataset—similar to current distributed training strategies. This refers to each replica as a constituent model. However, SISA training deviates from current strategies in the way incremental model updates are propagated or shared—there is no flow of information between constituent models.

**7.4. Absent access to the underlying dataset, is it possible to identify whether an AI model was trained on a particular piece of training material?**

In practice, it's challenging to definitively identify whether an AI model was trained on a specific piece of training material without access to the dataset. Many AI developers take steps to obfuscate or aggregate their training data sources to protect proprietary information and data privacy. Additionally, the process of training a model often involves complex transformations and generalizations, making it challenging to trace model outputs back to specific input data. However, there are some *possible* methods:

1.  **Model Behavior and Output**: One indirect way to infer the presence of specific training materials is to analyze the behavior and output of the AI model. If the model consistently generates or classifies content related to the specific training material, it could suggest the influence of that material. However, this is not a definitive method as the model may have been generalized from the training data.
2.  **Model Fine-Tuning**: If a model is fine-tuned on a specific dataset or task, it may exhibit characteristics that align with the fine-tuning data. Researchers and practitioners sometimes disclose information about fine-tuning, which can provide clues about the training material.
3.  **Textual Artifacts**: In natural language processing (NLP) models, certain textual artifacts or biases may emerge from the training data. Researchers have developed methods to detect and analyze these biases, which could hint at the presence of specific training materials. However, this doesn't pinpoint the exact sources.
4.  **Metadata and Attribution**: In some cases, metadata or attribution information may be embedded in the model or disclosed by the developers. This can provide information about the sources of training data.
5.  **Comparison with Known Data**: If you have access to known data that closely matches the training material in question, you could compare the model's behavior or outputs with that data to look for similarities. This method can provide some insights but is not foolproof.
6.  **Reverse Engineering**: In some cases, researchers have attempted to reverse engineer AI models to gain insights into their training data. However, this can be a complex and resource-intensive process and may not yield complete information.

**8. Under what circumstances would the unauthorized use of copyrighted works to train AI models constitute fair use? Please discuss any case law you believe relevant to this question.**

The fair use analysis should be the same for AI training as it would be for any other use, employing the following factors: the purpose and character of the use, including whether such use is of a commercial nature or is for non-profit educational purposes; the nature of the copyrighted work; the amount and substantiality of the portion used in relation to the copyrighted work as a whole; and the effect of the use upon the potential market for or value of the copyrighted work. While not directly on point the *Google v. Oracle* and *Warhol Foundation v. Goldsmith* Supreme Court cases, discussed below, begin to shed some light on fair use in general, albeit not in the specific context of AI.

**8.1.** In light of the Supreme Court's recent decisions in *Google* v. *Oracle America* [41] and *Andy Warhol Foundation* v. *Goldsmith,*[42] how should the "purpose and character" of the use of copyrighted works to train an AI model be evaluated? What is the relevant use to be analyzed? Do different stages of training, such as pre-training and fine-tuning,[43] raise different considerations under the first fair use factor?

Purpose and character should always be considered as part of a fair use analysis, but the more likely arguments will be made under a transformative fair use defense. In *Google v. Oracle*, the Supreme Court found that Google's unlicensed use of the declarative code of Oracle's APIs was not an unlawful exploitation of Oracle's work but rather a necessary step in the creation of a transformative new and original work—Google's Android-based phones. In the *Warhol* case, the Supreme Court found that the Warhol Foundation's licensing of an orange silkscreen made by Andy Warhol based on Goldsmith's photograph was not a fair use. Reconciling these two cases is challenging and applying them to the use of copyrighted materials for AI training even more so. The ultimate determination as to whether the use of copyrighted material to train an AI algorithm is fair use is likely to hinge on whether the trained algorithm is sufficiently different than the original works to constitute a transformative use. As with most IP cases, the facts of any individual instance will be outcome-determinative.

**8.2.** How should the analysis apply to entities that collect and distribute copyrighted material for training but may not themselves engage in the training?

At a high level, the analysis should be the same, regardless of the activities of the accused infringer, i.e., whether the use make is fair use. As mentioned above, this will likely turn on whether the accused infringer's use is transformative. That analysis may be different for an aggregator reselling the copyrighted material as opposed to the end-user of the data who is training an algorithm. The transformative use argument is much stronger for the latter, assuming that the output of that training is transformative as compared to the input.

**8.3.** The use of copyrighted materials in a training dataset or to train generative AI models may be done for noncommercial or research purposes. How should the fair use analysis apply if AI models or datasets are later adapted for use of a commercial nature? Does it make a difference if funding for these noncommercial or research uses is provided by for-profit developers of AI systems?

The fair use analysis should remain the same, but each of these fact patterns is likely to generate a difference outcome. Non-commercial or research purposes are generally regarded as fair use. That should not change simply because AI has been inserted into the analysis. Adaptation for commercial use will change this analysis. Again, IP cases are very fact-specific, but fair use analysis should remain the same, at least at a high level.

**8.4.** What quantity of training materials do developers of generative AI models use for training? Does the volume of material used to train an AI model affect the fair use analysis? If so, how?

The quantity of training materials used by developers of generative AI models can vary widely based on several factors, including the specific AI model, complexity of the task, available resources, and desired level of performance. Large-scale generative models like BERT, LLAM, and GPT-3 and GPT-4, which are designed to generate human-like text, typically use massive datasets consisting of hundreds of gigabytes, or even terabytes, of text data. In short, the volume of data required to train an AI model varies from model to model and depends on its ultimate use/purpose. The volume of training material used to train an AI model can potentially affect the fair use analysis in the context of copyright law, but it is just one of several factors considered in such an analysis. Here's how the volume of training data affects fair use determinations:

1. **Amount Used in Relation to the Whole**: Fair use analysis often considers whether the amount of copyrighted material used is reasonable in relation to the purpose of use. Using a substantial portion of a large dataset for training might be considered reasonable if it's necessary to achieve the model's intended functionality.
2. **Transformative Purpose**: One of the key factors in fair use analysis is whether the use of copyrighted material is transformative. If the AI model's training process results in a highly transformative output (e.g., generating entirely new text), it may weigh in favor of fair use.
3. **Effect on Market Value**: Fair use analysis also considers whether the use of copyrighted material negatively affects the market value of the original work. The larger the dataset used, the more potential there is for overlap with the original work, which could impact this factor.
4. **Commercial vs. Non-Commercial Use**: Whether the use of training materials is commercial or non-commercial can also influence fair use analysis. Commercial use is generally subject to stricter scrutiny.

It's important to note that fair use is a complex legal doctrine, and courts evaluate each case on its individual merits. The volume of training data is just one factor among many considered in a fair use analysis, and the outcome depends on the specific circumstances and legal jurisdiction.

Developers of generative AI models often strive to comply with copyright laws and ethical standards. They may obtain licenses for copyrighted data when necessary, curate training data to minimize copyright concerns, and take steps to ensure that the use of training materials is justifiable based on the model's transformative purpose and other factors.

**8.5.** **Under the fourth factor of the fair use analysis, how should the effect on the potential market for or value of a copyrighted work used to train an AI model be measured?** [46] **Should the inquiry be whether the outputs of the AI system incorporating the model compete with a particular copyrighted work, the body of works of the same author, or the market for that general class of works?**

As mentioned earlier, each of these cases will turn on the individual facts. For example, training an algorithm on the works of an author to create similar works will be a more compelling plaintiff's case than using that training data to create a piece of art. The fourth factor is the effect of the use upon the potential market for or value of the copyrighted work. Each individual case must be decided on the merits of that

case, so certainly the inquiry should focus on the copyrighted work itself and perhaps the body of works of the same author.

9. **Should copyright owners have to affirmatively consent (opt in) to the use of their works for training materials, or should they be provided with the means to object (opt out)?**

Both "opt-in" and "opt-out" pose challenges. It may be worth considering either a voluntary or compulsory licensing scheme for certain content used as AI training data.

9.1. **Should consent of the copyright owner be required for all uses of copyrighted works to train AI models or only commercial uses? [47]**

The consent of the copyright owner is not required for all uses of copyrighted works to train AI model. The doctrine of fair use should be applied here in the first instance, as with other uses of copyrighted works. The basic premise in a fair use analysis is that there are certain uses of copyrighted works that should not be considered infringement of the copyright owner's rights. As the federal government issues new regulations, those regulations should be harmonized with the fair use doctrine rather than supplanting it.

9.2. **If an "opt out" approach were adopted, how would that process work for a copyright owner who objected to the use of their works for training? Are there technical tools that might facilitate this process, such as a technical flag or metadata indicating that an automated service should not collect and store a work for AI training uses? [48]**

Copyright owners could embed metadata or watermarks in their content to indicate their preferences regarding usage in AI training. AI models, data providers, and training datasets could honor these indicators. For example, an image could include metadata specifying whether it can be used for AI training. Several tools for watermarking have been available for quite some time. However, Google recently released SynthID, a tool designed to watermark AI-generated images in a way that is imperceptible to humans but discernible by an AI-detection program.

API (Application Programmer Interface)-Based Controls allow Platform providers to offer APIs that allow copyright owners to register their works and set usage preferences. When AI models access content through such platforms, they could respect these preferences. For instance, a copyright owner could set their content to be excluded from AI training datasets.

9.3. **What legal, technical, or practical obstacles are there to establishing or using such a process? Given the volume of works used in training, is it feasible to get consent in advance from copyright owners?**

The obstacles here appear to be primarily practical in nature. Given the volume of works used in training (and the fact that more are being generated every minute), it does not seem feasible to obtain advance consent from copyright owners.

**9.4.** **If an objection is not honored, what remedies should be available? Are existing remedies for infringement appropriate or should there be a separate cause of action?**

Existing remedies for infringement should be adequate. Creating a separate cause of action is a slippery slope, as every new technology could present the opportunity to amend the law, which would be unwise.

**9.5.** **In cases where the human creator does not own the copyright—for example, because they have assigned it or because the work was made for hire—should they have a right to object to an AI model being trained on their work? If so, how would such a system work?**

If a human creator does not own the copyright for any reason, they do not have the right to license the work or prevent it from being used. We do not have a "moral rights" system in the US, and it does not make sense to create one for AI. If someone has alienated their IP rights, presumably they have received whatever value they expected from it. If they were cheated, that is a different cause of action/remedy.

10. **If copyright owners' consent is required to train generative AI models, how can or should licenses be obtained?**

Different licenses are being developed for this purpose, similar to open source and creative commons licenses for software and source code. These licenses could be very helpful. The other option is a compulsory licensing scheme, which would certainly facilitate AI innovation, but would likely be met with objection from some copyright holders as an unlawful taking.

**10.1.** **Is direct voluntary licensing feasible in some or all creative sectors?**

As simplified licenses become developed and available, this approach could become customary, similar to open-source licensing of software.

**10.2.** **Is a voluntary collective licensing scheme a feasible or desirable approach?** [49] **Are there existing collective management organizations that are well-suited to provide those licenses, and are there legal or other impediments that would prevent those organizations from performing this role? Should Congress consider statutory or other changes, such as an antitrust exception, to facilitate negotiation of collective licenses?**

A voluntary collective licensing scheme could work but comes with certain tradeoffs. It has the advantage of being voluntary, but that also means that there may be many works that will not be available under this approach. This approach has existed in the music world for many years, with organizations like ASCAP and BMI as two of the more prominent players. Diligence would need to be performed (perhaps in the form of an FRP) as to which organization (if any) would be suitable for this role.

**10.3. Should Congress consider establishing a compulsory licensing regime?** [50] **If so, what should such a regime look like? What activities should the license cover, what works would be subject to the**

**license, and would copyright owners have the ability to opt out? How should royalty rates and terms be set, allocated, reported and distributed?**

It is worthy of consideration. The license would need to cover all activities that would otherwise constitute infringement. An opt-out process would seem to defeat the purpose of establishing this regime. The royalty rates could be set by a panel of experts for each industry. The revenues generated could be distributed to individuals in order to encourage innovation—whether by way of grants, scholarships or other similar means.

**10.4. Is an extended collective licensing scheme [51] a feasible or desirable approach?**

This could also work. Each of these options present trade-offs: the more encompassing the regime, the more effective it will be in promoting innovation through the use of AI by preventing litigation relating to the use of training material for this purpose. On the other hand, content creators may complain that a compulsory licensing regime constitutes an unlawful taking of their intellectual property. The more permissive and voluntary the regime, the less effective it will be for AI innovators, but will presumably generate less complaints from content creators who have a choice as to whether to license their content in this way.

**10.5. Should licensing regimes vary based on the type of work at issue?**

No, they should not because all works will be used in the same manner—to train AI. The aim should be to use the license revenues to encourage human creativity.

**11. What legal, technical or practical issues might there be with respect to obtaining appropriate licenses for training? Who, if anyone, should be responsible for securing them (for example when the curator of a training dataset, the developer who trains an AI model, and the company employing that model in an AI system are different entities and may have different commercial or noncommercial roles)?**

Practically speaking, it could be challenging to obtain all the correct licenses for training material one wishes to use. In some instance there may be a simple license available, such as a Community Data License Agreement, but in other instances there may be no license available, either because the owner does not wish to license their work, or they simply have not thought to do so. As with any copyrighted material, anyone who could be a potential infringer would be well advised to obtain a license, regardless of where they fall in the workflow.

**12. Is it possible or feasible to identify the degree to which a particular work contributes to a particular output from a generative AI system? Please explain.**

Identifying the degree to which a particular work contributes to a specific output from a generative AI system can be challenging and is often not feasible in practice. Several factors contribute to the difficulty of making such determinations:

1. **Complexity of Model Learning**: Generative AI systems, especially large-scale models like GPT-3 or GPT-4, are highly complex. They learn from vast amounts of training data and develop intricate internal representations. It's challenging to pinpoint the direct influence of a single work or input in the final output because the model's responses are the result of the collective knowledge learned from the entire training dataset.
2. **Feature Attribution**: Attribution techniques, such as feature attribution or saliency methods, attempt to highlight the input features that contribute most to a model's output. However, these methods provide a high-level understanding of feature importance and may not precisely identify the contribution of a specific work.
3. **Data Aggregation**: Training data for generative models typically consists of a wide variety of sources, making it challenging to isolate the impact of one particular work. The model may have seen countless variations of similar content from different sources.
4. **Temporal Sequence**: In many cases, generative models work with sequences of data, such as text or speech. The model's responses are influenced not only by the immediate input but also by the context of the entire sequence, including previous inputs and internal states.
5. **Internal Representation**: Generative models transform input data into a complex internal representation, which is difficult to interpret directly. Deciphering how a specific work contributes to this representation is a non-trivial task.
6. **Latent Factors and Generalization**: Generative models often learn "latent factors" and generalize from their training data. They may generate outputs that resemble multiple sources simultaneously, and it can be challenging to attribute specific content to a single source.

It is worth noting here the complexity of these latent factors or variables defined by these algorithms and methodologies:

In a generative AI model, latent variables are hidden or unobservable variables that represent underlying factors or features of the data. These variables are used to capture the structure and variability of the data and are fundamental to the generative process of the model. Latent variables are essential for generative models to learn and generate new data samples that resemble the training data. Here's a closer look at latent variables in generative AI models:

1. **Generative Models and Latent Variables**:
   - Generative AI models aim to learn a probabilistic representation of a dataset. They model the probability distribution of the data and can generate new data samples that are statistically like the training data.
   - Latent variables are introduced to capture the underlying patterns, features, or factors that explain the observed data. They act as hidden representations within the model.
2. **Variational Autoencoders (VAEs)**:
   - VAEs are a type of generative model that use latent variables. In a VAE, an encoder network maps observed data into a latent space, and a decoder network map points in the latent space back to data space.
   - The latent variables in VAEs represent continuous, probabilistic distributions and are used to generate new data points by sampling from these distributions.

3. **Generative Adversarial Networks (GANs)**:
   - GANs consist of two networks: a generator and a discriminator. The generator generates data samples from random noise (latent variables), while the discriminator distinguishes between real and generated data.
   - The latent variables in GANs represent random noise vectors that are used as input to the generator to produce synthetic data.
4. **Topic Modeling**:
   - In natural language processing (NLP), latent Dirichlet allocation (LDA) is an example of a topic modeling technique that uses latent variables. LDA aims to discover latent topics within a collection of documents.
   - The latent variables in LDA represent the mixture proportions of topics for each document and the distribution of words for each topic.
5. **Dimensionality Reduction**:
   - In some cases, latent variables are introduced for dimensionality reduction purposes. For example, principal component analysis (PCA) is a dimensionality reduction technique that identifies latent variables (principal components) that capture the most variance in the data.
6. **Applications**:
   - Latent variables can be used for various generative tasks, such as image synthesis, text generation, recommendation systems, and anomaly detection.
   - They are also used in unsupervised learning tasks, where the model learns to represent data in a more compact and informative way.
7. **Interpretability**:
   - The interpretation of latent variables can be challenging, as they are often abstract and may not have a direct, human-interpretable meaning.
   - Understanding the semantics of latent variables is an active research area in AI and machine learning.

**13. What would be the economic impacts of a licensing requirement on the development and adoption of generative AI systems?**

Implementing a licensing requirement for the development and adoption of generative AI systems would have several economic impacts, both positive and negative. The specific effects would depend on the scope, terms, and enforcement of such licensing requirements. Here are some potential economic impacts to consider:

**Positive Economic Impacts**:

1. **Revenue Generation**: Licensing requirements can create a new revenue stream for content creators and copyright holders. Developers and organizations using generative AI systems may need to obtain licenses, leading to licensing fees and royalties paid to content owners.

2. **Incentive for Content Creation**: Licensing requirements may incentivize content creators to produce high-quality, valuable training data and works for generative AI systems. They may see increased demand for their content due to its potential use in AI models.
3. **Legal and Regulatory Compliance**: Licensing can provide clarity and legal protection for AI developers and users. It can help organizations avoid legal disputes and potential liability for copyright infringement.
4. **Market for Licensing Services**: Licensing requirements can create a market for licensing services and intermediaries that facilitate licensing agreements, negotiations, and royalty payments. This can stimulate economic activity in the licensing industry.

**Negative Economic Impacts**:

1. **Increased Costs**: Licensing requirements can lead to increased costs for AI developers and organizations. They may need to pay for licenses, royalties, and legal services, potentially raising the barrier to entry for smaller players and startups.
2. **Stifled Innovation**: Stringent licensing requirements may stifle innovation in the AI field, particularly for researchers and organizations with limited resources. High licensing fees and complex negotiations could discourage experimentation and development.
3. **Reduced Accessibility**: Licensing fees and requirements could limit access to generative AI technologies. Smaller businesses, educational institutions, and non-profit organizations may find it more challenging to afford licenses, hindering widespread adoption.
4. **Potential for Monopolies**: Licensing requirements may favor large content owners and established players, potentially leading to the consolidation of power and monopolistic behavior in the AI industry.
5. **Legal and Compliance Costs**: Organizations may incur significant legal and compliance costs to navigate complex licensing agreements and ensure they are adhering to licensing terms, potentially diverting resources from other activities.
6. **Impact on Open Source**: Licensing requirements could affect the development and adoption of open-source AI projects. The open-source community might face challenges in incorporating licensed data and complying with licensing terms.
7. **International Variability**: Licensing requirements may vary by jurisdiction, adding complexity and uncertainty for AI developers and users operating internationally.
8. **Market Uncertainty**: Licensing requirements could introduce uncertainty into the AI market, as developers may be unsure about the availability and terms of licensing for specific datasets or works.

14. **Please describe any other factors you believe are relevant with respect to potential copyright liability for training AI models.**

N/A.

**Transparency & Recordkeeping**

**15. In order to allow copyright owners to determine whether their works have been used, should developers of AI models be required to collect, retain, and disclose records regarding the materials used to train their models? Should creators of training datasets have a similar obligation?**

While it is a best practice to keep records of the data used for any endeavor for a variety of reasons, including intellectual property/licensing compliance as well privacy and security, it is not clear that mandatory disclosure would provide a public benefit in all cases while introducing additional burden on the AI developer or creator of datasets.

There may be other reasons for recordkeeping; for example, if the creator of a dataset is charging for the use of that dataset, the licensee community is likely to require some sort of indemnity, which would make the imperative for knowing what is in the dataset all the more compelling for its creator, but the level of recordkeeping should be left to the individual organization.

**15.1. What level of specificity should be required?**

It is not clear that this should be a statutory requirement, as this could be very burdensome on AI innovators.  This is a decision that every individual organization would need to make based on their usage of data and content.

**15.2. To whom should disclosures be made?**

Disclosures may need to be made in the course of discovery in a litigation.  It would not serve public policy to create a new, statutory obligation to do so in a general context.

**15.3. What obligations, if any, should be placed on developers of AI systems that incorporate models from third parties?**

These obligations already exist by virtue of current intellectual property laws.  It is preferable to work within the existing legal framework rather than create new laws for new technology.

**15.4. What would be the cost or other impact of such a recordkeeping system for developers of AI models or systems, creators, consumers, or other relevant parties?**

This requirement could be very burdensome and could effectively put smaller innovators out of the running in AI.  That would be a very negative result here.

**16. What obligations, if any, should there be to notify copyright owners that their works have been used to train an AI model?**

N/A

**17. Outside of copyright law, are there existing U.S. laws that could require developers of AI models or systems to retain or disclose records about the materials they used for training?**

Yes. Data privacy and security laws may require such disclosure and while not law per se, there are certain standards-setting bodies (such as IEEE) that encourage transparency in AI. The Department of Defense has adopted ethical principles for artificial intelligence (https://www.defense.gov/News/News-Stories/article/article/2094085/dod-adopts-5-principles-of-artificial-intelligence-ethics/). The Defense Intelligence Unit has issued guidelines for AI design, development and deployment (https://www.diu.mil/responsible-ai-guidelines). In addition, organizations are increasingly setting their own internal standards for AI development.

**Generative AI Outputs**

*If your comment applies only to a particular subset of generative AI technologies, please make that clear.*

**Copyrightability**

**18. Under copyright law, are there circumstances when a human using a generative AI system should be considered the "author" of material produced by the system? If so, what factors are relevant to that determination? For example, is selecting what material an AI model is trained on and/or providing an iterative series of text commands or prompts sufficient to claim authorship of the resulting output?**

Here it is important to remember that photography was initially not considered "art" because the thought was that the camera took the photo, not the human. In the instance where the human selects the inputs for an AI-generated output, a case can be made that this is analogous to a human arranging the composition of a photograph and this output should be protected by copyright. The less input from a human, the more tenuous the argument is that the human created the output.

**19. Are any revisions to the Copyright Act necessary to clarify the human authorship requirement or to provide additional standards to determine when content including AI-generated material is subject to copyright protection?**

It is generally ill-advised to adjust legislation to accommodate new technologies. While the concept of copyright substantially predates modern technology, it has generally worked. Modifying copyright law to accommodate AI would not only be difficult and time-consuming, itis likely to create unintended consequences that are undesirable.

**20. Is legal protection for AI-generated material desirable as a policy matter? Is legal protection for AI-generated material necessary to encourage development of generative AI technologies and systems? Does existing copyright protection for computer code that operates a generative AI system provide sufficient incentives?**

Assuming that the materials meet all the other requirements for copyrightability, legal protection of human creation serves as an incentive for innovation. Against that must be balanced the concern that over-protecting intellectual property can stifle innovation. These considerations should be taken into account. Existing copyright protection for computer code that operated a generative AI system is fine for protecting the system but goes no further.

**20.1. If you believe protection is desirable, should it be a form of copyright or a separate *sui generis* right? If the latter, in what respects should protection for AI-generated material differ from copyright?**

It would be optimal to determine a way to use existing copyright law to protect works that meet all of the requirements of those laws. Given the pace at which technology is progressing, to change the copyright law for every new technology that comes along seems inadvisable.

**21. Does the Copyright Clause in the U.S. Constitution permit copyright protection for AI-generated material? Would such protection "promote the progress of science and useful arts"? [52] If so, how?**

The Copyright Clause certainly does not prohibit protection of AI-generated material. The authors of the Constitution highly valued innovation and part of the beauty of that document is that it has served this country very well for almost 250 years, partially because it has been able to accommodate changes and progress. This country continues to value innovation, and appropriate protection of copyrightable material should not depend on the method through which is was generated. Only then can we truly "promote the progress of science and useful arts."

**Infringement**

**22. Can AI-generated outputs implicate the exclusive rights of preexisting copyrighted works, such as the right of reproduction or the derivative work right? If so, in what circumstances?**

Yes. The question is likely to turn on whether the resulting work is a straight-up infringement, a derivative work, or a transformative use. This analysis will likely turn on the fair use analysis described above.

**23. Is the substantial similarity test adequate to address claims of infringement based on outputs from a generative AI system, or is some other standard appropriate or necessary?**

The substantial similarity test may not yield a determination of copyright infringement under these circumstances. It may be more logical to compare the training data to the copyrighted work rather than the output of the generative AI to the copyrighted work. For example, in the Getty Images v. Stability AI case, the example of the output was an AI-generated image of a cat wearing a scarf. Getty is not alleging that any of their photos bear substantial similarity to that image, but rather, that the Getty images used to train the algorithm were identical to their copyrighted works, without which the algorithm would not have been able to generate the cat-in-scarf image.

**24. How can copyright owners prove the element of copying (such as by demonstrating access to a copyrighted work) if the developer of the AI model does not maintain or make available records of what training material it used? Are existing civil discovery rules sufficient to address this situation?**

This element is difficult, not only from a discovery standpoint, but from a technology perspective. It is unclear whether AI training data is copied when used for this purpose, even in an ephemeral manner. Existing discovery rules are adequate, although it will be more difficult, but not impossible, to establish access to the training material. If a lawsuit has been filed, there is presumably a Rule 11 basis for doing so, so one would imagine that the plaintiff has some reason to believe that its copyrighted works are being used as training material.

**25. If AI-generated material is found to infringe a copyrighted work, who should be directly or secondarily liable—the developer of a generative AI model, the developer of the system incorporating that model, end users of the system, or other parties?**

The existing caselaw and theories on contributory infringement and vicarious liability can be applied to any actors in the system. Both theories require some sort of knowledge and intent, i.e., for contributory infringement, that the accused either induces, causes, or materially contributes to the infringing conduct of another (with knowledge of the infringing activity). For vicarious liability, the accused must have both the right and ability to supervise the infringing activity and also must have a direct financial interest in the activity.

**25.1. Do "open-source" AI models raise unique considerations with respect to infringement based on their outputs? [53]**

It is possible that open-source models may raise unique considerations with respect to infringement. For example, in an Apache-type license, the licensee may find themselves unlicensed to certain software if they assert a claim. One could devise a similar type license here, whereby a licensee could give up certain rights in their IP by licensing open-source models and, like the GPL variety of licenses, they could be written as "copyleft."

**26. If a generative AI system is trained on copyrighted works containing copyright management information, how does 17 U.S.C. 1202(b) apply to the treatment of that information in outputs of the system?**

Because the outputs of the system are generally not copies of the inputs, 17 USC 1202(b) does not apply.

**27. Please describe any other issues that you believe policymakers should consider with respect to potential copyright liability based on AI-generated output.**

**N/A**

**Labeling or Identification**

**28. Should the law require AI-generated material to be labeled or otherwise publicly identified as being generated by AI? If so, in what context should the requirement apply and how should it work?**

The public has an interest in understanding the provenance of the material it consumes. AI-generated material can often be flat-out wrong (e.g., Chat GPT), so it may be in the public interest to label such material as having been AI-generated.

**28.1. Who should be responsible for identifying a work as AI-generated?**

The party who is in the best position to know whether a work is AI-generated is the author, but the party displaying the work may also bear some responsibility to alert the public that the materials they are consuming was AI-generated. As the use and display of AI- generated materials increases, it may be more common to see representations and warranties relating to how the work was generated, and websites may choose to include banners or pop-ups (like with privacy) that alert the user that there may be AI-generated materials on the site and give them an opportunity to exit if that is not acceptable to them.

**28.2. Are there technical or practical barriers to labeling or identification requirements?**

There are both technical and practical barriers to labeling and identification requirements. First, a decision must be made as to what constitutes AI-generated material. Once that decision has been made, technology such as digital rights management watermarking could be used, although these watermarks can be removed. There may also be a scale issue here, depending on the volume of the content.

**28.3. If a notification or labeling requirement is adopted, what should be the consequences of the failure to label a particular work or the removal of a label?**

The consequences should be similar to other false labeling or advertising-type infractions, i.e., civil liability with a private right of action to recover damages.

**29. What tools exist or are in development to identify AI-generated material, including by standard-setting bodies? How accurate are these tools? What are their limitations?**

Text-based tools include Content at Scale, WriterAI, and ZeroGPT.

For images, the following can be used, similar to DRM:

**Reverse Image Search Engines:** Tools like Google Image Search can help identify similar images on the web, potentially leading to the source of the AI-generated image.

**Metadata Analysis:** Many AI-generated images might lack metadata or have unusual metadata. Examining the metadata associated with an image can provide clues about its authenticity.

**Exif Data:** Some AI-generated images might not include Exif data, which is typically present in images captured by cameras or smartphones. Analyzing the Exif data can help in identifying anomalies.

**Error Level Analysis (ELA**): ELA is a method used to detect modifications in an image. It can highlight areas with significant alterations, which could indicate AI-generated or manipulated content.

**Digital Forensics Tools:** Tools like Forensically offer various features to analyze image authenticity, including ELA and other forensics techniques.

**AI Tools for Image Forgery Detection**: Some AI tools have been developed to detect AI-generated or manipulated images. These tools often use AI to identify anomalies and signs of image manipulation.

**Blockchain Verification:** Some platforms use blockchain technology to verify the authenticity of images, allowing users to check if the image has been tampered with.

**User Community:** Online communities often work together to identify fake or AI-generated images. Websites like Reddit and platforms like Twitter have users dedicated to fact-checking and verifying images.

The following article provides an analysis of current art:

Elkhatat, A.M., Elsaid, K. & Almeer, S. Evaluating the efficacy of AI content detection tools in differentiating between human and AI-generated text. *Int J Educ Integr* **19**, 17 (2023). https://doi.org/10.1007/s40979-023-00140-5

**Additional Questions About Issues Related to Copyright**

**30. What legal rights, if any, currently apply to AI-generated material that features the name or likeness, including vocal likeness, of a particular person?**

The right of publicity under various state laws prevent the commercial use of another's name, likeness or voice. In addition, Section 43(a) of the Lanham Act may create a cause of action where the particular circumstances concerning the use of another's name, likeness or voice constitute unfair competition.

In the European Union, the use of another's name, likeness or voice in a manner sufficient to identify the person is prohibited without the person's express written consent.

**31. Should Congress establish a new federal right, similar to state law rights of publicity, that would apply to AI-generated material? If so, should it preempt state laws or set a ceiling or floor for state law protections? What should be the contours of such a right?**

If possible, it is desirable to work within the framework of current laws and regulations.

**32. Are there or should there be protections against an AI system generating outputs that imitate the artistic style of a human creator (such as an AI system producing visual works "in the style of" a specific artist)? Who should be eligible for such protection? What form should it take?**

It seems as though this issue should be able to be handled under the current laws. If there is legislation or case law that prohibits using source material (like a comic's standup routines) to generate competing source material, that will protect the original author. It does not seem wise to create additional laws for new technologies. Copyright has stood the test of time. We should not abandon it too hastily.

**33. With respect to sound recordings, how does section 114(b) of the Copyright Act relate to state law, such as state right of publicity laws? [54] Does this issue require legislative attention in the context of generative AI?**

Section 114(b) serves a very different purpose to State publicity laws. Section 114(b) as currently drafted, would not apply to AI generated sound recordings.

**34. Please identify any issues not mentioned above that the Copyright Office should consider in conducting this study.**

**N/A**