

# Trust in AI

Many AI systems, particularly those built on deep learning neural networks, are fundamentally unexplainable and unpredictable. [4] [16]

These systems learn to classify data by adjusting parameters within interconnected “neurons”, similar to the human brain. However, due to the vast number of parameters (often in the trillions), the reasons behind the decisions made by these AI systems are often cloudy. [4] [16]

When humans can't comprehend something, they often trust it less.

Therefore, trust is grounded in predictability, and if AI systems behave in ways that are not expected, their perceived trustworthiness diminishes. [1] [5]

Having the capability to generate tremendous benefits for individuals and society, AI also gives rise to certain risks that should be properly managed.

Given that, overall, AI's benefits outweigh its risks, we must ensure to follow the road that maximizes the benefits of AI while minimizing its risks. To ensure that we stay on the right track, a human-centric approach to AI is needed, forcing us to keep in mind that the development and use of AI should not be seen as a means in itself, but as having the goal to increase human well-being. [4] [5] [7]

Trustworthy AI will be our north star, since human beings will only be able to confidently and fully reap the benefits of AI if they can trust the technology. [2]  
Trustworthy AI has two components:

(1) it should respect fundamental rights, applicable regulation and core principles and values, ensuring an “ethical purpose” and (2) it should be technically robust and reliable since, even with good intentions, a lack of technological mastery can cause unintentional harm. [2]

## Ensuring AI's ethical purpose

Ensuring AI's ethical purpose, by setting out the fundamental rights, principles and values that it should comply with guidance on the realization of Trustworthy AI, tackling both ethical purpose and technical robustness. [2] [3]

This is done by listing the requirements for Trustworthy AI and offering an overview of technical and non-technical methods that can be used for its implementation. [2]

Operationalizes the requirements by providing a concrete but non-exhaustive assessment list for Trustworthy AI. This list is then adapted to specific use cases. [2]

## The 7 Key requirements

The guidelines put forward a set of 7 key requirements that AI systems should meet in order to be deemed trustworthy:

- **Human agency and oversight:** AI systems should empower human beings, allowing them to make informed decisions and fostering their fundamental rights.
- **Technical Robustness and safety:** AI systems need to be resilient and secure.
- **Privacy and data governance:** Full respect for privacy and data protection, adequate data governance mechanisms must also be ensured.
- **Transparency:** The data, system and AI business models should be transparent.
- **Diversity, non-discrimination and fairness:** Unfair bias must be avoided.
- **Societal and environmental well-being:** AI systems should benefit all human beings, including future generations.
- **Accountability:** Mechanisms must be in place to ensure responsibility and accountability for AI systems and their outcomes. [10] [2] [3] [7]

## Lack of transparency

The lack of transparency in AI systems is a significant barrier to their trustworthiness. A suggestion is that to build trust in AI, we need to understand how it makes decisions. This involves not only understanding the algorithms and data used to train the AI, but also the values and biases that may be embedded in these algorithms. [2] [3] [12]

## Important notes

We need to ensure that AI does not discriminate against certain groups or individuals and that it respects privacy and confidentiality.

### Government

There are some AIs that are banned in the European atmosphere. Due to the nature of the "Fediverse" we should know what is banned in what parts of the world and ensure that they are not available globally or continent wide.

The American government also has rules (bill) that the AI needs to follow. Here below you will get the complete list of possible systems/AI that are banned for the Americans and the Europeans.

- Biometrics categorization systems that use sensitive characteristics (e.g. political, religious, philosophical beliefs, sexual orientation, race);
- Untargeted scraping of facial images from the internet or CCTV footage to create facial recognition databases;
- Emotion recognition in the workplace and educational institutions;
- Social scoring based on social behavior or personal characteristics;
- AI systems that manipulate human behavior to circumvent their free will;
- AI used to exploit the vulnerabilities of people (due to their age, disability, social or economic situation).

[6] [8] [11]

There are some Law enforcement exemptions in these, but we wouldn't come in contact with it.

Even if we (as the group) don't make these AIs, we need to think about how we can minimize the risk to the application. Because people will always have intentions, but they are sometimes things that we do not want.

## How to make the Data more secure

### Anonymization

Using data without taking care to protect the identity of the data owner can lead to lots of problems and potential lawsuits. Data anonymization easily put, is ensuring that we can't tell the actual data owner by looking at the data.[15]

### Key Ways to Anonymize a Data Set

#### Replacing the key

In some cases removing the key and replacing it with a random number is sufficient. However, care must be taken to accommodate the composition of the data.[15]

Taking into account the composition of the table, we may find that the ids are repeating but the repeats are valid data points that we want to keep. In this case, we want to replace the same key with the same random number each time.[15]

#### Hashing & Random Mixing

When replacing the key with random keys is insufficient we can take other approaches.[15]

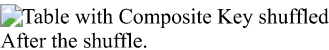
With Hashing converts the data into an alphanumeric or numeric code of fixed size, which cannot be easily reversed (if at all).[15]

Two key features of the hash:[15]

- It is technically feasible but practically impossible to reverse even the simplest hash i.e. you can go from A to hash(A) but not hash(A) to A.
- Hashing same value should provide same output each time and no two inputs should get same hash output (collision).

With random mixing is used when you might need data that has same structure and characteristics as the actual data set. In this you mix the table so the key is more random and the rows are not as it was before. An example below:[14] [15]

Table with Composite Key  
Before the shuffle.

Table with Composite Key shuffled  
After the shuffle.

If you also add some Fake data to it / delete some data it will also be more anonymized.

## Approaches

### Data Anonymization

[14] [15]

#	Approach	Description	Advantages	Disadvantages	Use Cases
1	Data Masking	Masks or disguises sensitive data by replacing characters with symbols or placeholders.	<ul style="list-style-type: none"><li>• Simplicity of implementation.</li><li>• Preservation of data structure.</li></ul>	<ul style="list-style-type: none"><li>• Limited protection against inference attacks.</li><li>• Potential negative impact on data analysis.</li></ul>	<ul style="list-style-type: none"><li>• Anonymizing email addresses in communication logs.</li><li>• Concealing rare names in datasets.</li><li>• Masking sensitive words in text documents.</li></ul>
2	Pseudonymization	Replaces sensitive data with pseudonyms or aliases or removes it altogether.	<ul style="list-style-type: none"><li>• Preservation of data structure.</li><li>• Data utility is generally preserved.</li><li>• Fine-grained control over pseudonymization rules.</li></ul>	<ul style="list-style-type: none"><li>• Pseudomized data is not anonymous data.</li><li>• Risk of re-identification is very high.</li><li>• Requires secure management of pseudonym mappings.</li><li>• Loss of fine-grained detail in the data.</li></ul>	<ul style="list-style-type: none"><li>• Protecting patient identities in medical research.</li><li>• Securing employee IDs in HR records.</li></ul>
3	Generalization/Aggregation	Aggregates or generalizes data to reduce granularity.	<ul style="list-style-type: none"><li>• Simple implementation.</li></ul>	<ul style="list-style-type: none"><li>• Risk of data distortion that affects analysis outcomes.</li><li>• Challenging to determine appropriate levels of generalization.</li></ul>	<ul style="list-style-type: none"><li>• Anonymizing age groups in demographic data.</li><li>• Concealing income brackets in economic research.</li></ul>
4	Data Swapping/Perturbation	Swaps or perturbs data values between records to break the link between individuals and their data.	<ul style="list-style-type: none"><li>• Flexibility in choosing perturbation methods.</li><li>• Potential for fine-grained control.</li></ul>	<ul style="list-style-type: none"><li>• Privacy-utility trade-off is challenging to balance.</li><li>• Risk of introducing bias in analyses.</li><li>• Selection of appropriate perturbation methods is crucial.</li></ul>	<ul style="list-style-type: none"><li>• E-commerce.</li><li>• Online user behavior analysis.</li></ul>
5	Randomization	Introduces randomness (noise) into the data to protect data subjects.	<ul style="list-style-type: none"><li>• Flexibility in applying to various data types.</li><li>• Reproducibility of results when using defined algorithms and seeds.</li></ul>	<ul style="list-style-type: none"><li>• Privacy-utility trade-off is challenging to balance.</li><li>• Risk of introducing bias in analyses.</li><li>• Selection of appropriate randomization methods is hard.</li></ul>	<ul style="list-style-type: none"><li>• Anonymizing survey responses in social science research.</li><li>• Online user behavior analysis.</li></ul>

12-03-2024 12:21

TrustInAISummary.md

#	Approach	Description	Advantages	Disadvantages	Use Cases
6	Data Redaction	Removes or obscures specific parts of the dataset containing sensitive information.	<ul style="list-style-type: none"> <li>• Simplicity of implementation.</li> </ul>	<ul style="list-style-type: none"> <li>• Loss of data utility, potentially significant.</li> <li>• Risk of removing contextual information.</li> <li>• Data integrity challenges.</li> </ul>	<ul style="list-style-type: none"> <li>• Concealing personal information in legal documents.</li> <li>• Removing private data in text documents.</li> </ul>
7	Homomorphic Encryption	Encrypts data in such a way that computations can be performed on the encrypted data without decrypting it, preserving privacy.	<ul style="list-style-type: none"> <li>• Strong privacy protection for computations on encrypted data.</li> <li>• Supports secure data processing in untrusted environments.</li> <li>• Cryptographically provable privacy guarantees.</li> </ul>	<ul style="list-style-type: none"> <li>• Encrypted data cannot be easily worked with if previously unknown to the user.</li> <li>• Complexity of encryption and decryption operations.</li> <li>• Performance overhead for cryptographic operations.</li> <li>• May require specialized libraries and expertise.</li> </ul>	<ul style="list-style-type: none"> <li>• Basic data analytics in cloud computing environments.</li> <li>• Privacy-preserving machine learning on sensitive data.</li> </ul>
8	Federated Learning	Trains machine learning models across decentralized edge devices or servers holding local data samples, avoiding centralized data sharing.	<ul style="list-style-type: none"> <li>• Preserves data locality and privacy, reducing data transfer.</li> <li>• Supports collaborative model training on distributed data.</li> <li>• Suitable for privacy-sensitive applications.</li> </ul>	<ul style="list-style-type: none"> <li>• Complexity of coordination among edge devices or servers.</li> <li>• Potential communication overhead.</li> <li>• Ensuring model convergence can be challenging.</li> <li>• Shared models can still leak privacy.</li> </ul>	<ul style="list-style-type: none"> <li>• Healthcare institutions collaboratively training disease prediction models.</li> <li>• Federated learning for mobile applications preserving user data privacy.</li> <li>• Privacy-preserving AI in smart cities.</li> </ul>
9	Synthetic Data Generation	Creates artificial data that mimics the statistical properties of the original data while protecting privacy.	<ul style="list-style-type: none"> <li>• Strong privacy protection with high data utility.</li> <li>• Preserves data structure and relationships.</li> <li>• Scalable for generating large datasets.</li> </ul>	<ul style="list-style-type: none"> <li>• Accuracy and representativeness of synthetic data may vary depending on the generator.</li> <li>• May require specialized algorithms and expertise.</li> </ul>	<ul style="list-style-type: none"> <li>• Sharing synthetic healthcare data for research purposes.</li> <li>• Synthetic data for machine learning model training.</li> <li>• Privacy-preserving data sharing in financial analysis.</li> </ul>
10	Secure Multiparty Computation (SMPC)	Enables multiple parties to jointly compute functions on their private inputs without revealing those inputs to each other, preserving privacy.	<ul style="list-style-type: none"> <li>• Strong privacy protection for collaborative computations.</li> <li>• Suitable for multi-party data analysis while maintaining privacy.</li> <li>• Offers security against collusion.</li> <li>• Complexity of protocol design and setup.</li> </ul>	<ul style="list-style-type: none"> <li>• Performance overhead, especially for large-scale computations.</li> <li>• Requires trust in the security of the computation protocol.</li> </ul>	<ul style="list-style-type: none"> <li>• Privacy-preserving data aggregation across organizations.</li> <li>• Collaborative analytics involving sensitive data from multiple sources.</li> <li>• Secure voting systems.</li> </ul>

References

[1] <https://www.theguardian.com/commentisfree/2021/oct/02/the-truth-about-artificial-intelligence-it-isnt-that-honest>

[2] [https://www.academia.edu/38205904/ETHICS\\_GUIDELINES\\_FOR\\_TRUSTWORTHY\\_AI](https://www.academia.edu/38205904/ETHICS_GUIDELINES_FOR_TRUSTWORTHY_AI)

[3] [https://www.academia.edu/67787782/Trustworthy\\_AI\\_From\\_Principles\\_to\\_Practices](https://www.academia.edu/67787782/Trustworthy_AI_From_Principles_to_Practices)

[4] <https://theconversation.com/why-humans-cant-trust-ai-you-dont-know-how-it-works-what-its-going-to-do-or-whether-itll-serve-your-interests-213115>

[5] <https://www.nist.gov/publications/trust-and-artificial-intelligence>

[6] <https://www.europarl.europa.eu/news/en/press-room/20231206IPR15699/artificial-intelligence-act-deal-on-comprehensive-rules-for-trustworthy-ai>

[7] <https://ieeexplore.ieee.org/abstract/document/10188681>

[8] <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>

[9] <https://scienceexchange.caltech.edu/topics/artificial-intelligence-research/trustworthy-ai>

[10] <https://www.hindawi.com/journals/ijis/2023/4459198/>

[11] <https://onlinelibrary.wiley.com/doi/full/10.1111/rego.12512>

[12] <https://ieeexplore.ieee.org/abstract/document/10086944>

[13] <https://www.tandfonline.com/doi/full/10.1080/07421222.2023.2196773?needAccess=true> (not used because it's paid)

[14] <https://mostly.ai/blog/data-anonymization-tools>

[15] <https://towardsdatascience.com/anonymizing-data-sets-c4602e581a35>

[16] [https://openaccess.thecvf.com/content/WACV2023/papers/Ciftci\\_My\\_Face\\_My\\_Choice\\_Privacy\\_Enhancing\\_Deepfakes\\_for\\_Social\\_Media\\_WACV\\_2023\\_paper.pdf](https://openaccess.thecvf.com/content/WACV2023/papers/Ciftci_My_Face_My_Choice_Privacy_Enhancing_Deepfakes_for_Social_Media_WACV_2023_paper.pdf)