



# Practical Guidelines on Ethical Machine Learning

## DE\_398\_EVIDEN - Ethical Machine Learning

Issued by: Eviden / Author names: Sebastian Schmidt, Lisa Fink, Djordje Benn-Maksimovic

Date: 30.01.2026

Ref: DE\_398\_ D398.6.3.1\_202601\_Practical\_Guidelines\_v1.0.docx

Official reference number service contract: 2024/DE\_398\_EVIDEN

Status: Public

*This document has been produced in the context of the Destination Earth Initiative and relates to tasks entrusted by the European Union to the European Centre for Medium-Range Weather Forecasts implementing part of this Initiative.*

*This document is funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the European Commission can be held responsible for them. The European Centre for Medium-Range Weather Forecasts is not liable in respect of this document and gives no warranty for the information provided.*

## Change Log

Version	Date	Description
0.1	30.06.2025	First draft
0.2	13.08.2025	Second draft version/Submitted version
0.3	11.09.2025	Third draft
0.4	25.11.2025	Fourth draft with added graphics for DestinE internal seminar
0.5	12.12.2025	Minor adaptations
1.0	30.01.2026	Version after internal consultation and as presented at workshop

## 1 Executive Summary

These practical guidelines support developers in building trustworthy AI-based weather and climate forecasting systems. Based on the “Guideline for Trustworthy Artificial Intelligence” by Fraunhofer IAIS [1] and prior research, they address key challenges such as the influence of training data on model performance, model reproducibility, and fairness. The guidelines provide actionable recommendations across the AI lifecycle, helping developers identify risks, define goals, and implement mitigation strategies. They emphasize consistent documentation and communication to support downstream providers in meeting legal obligations, including those under the EU AI Act [2].

## 2 Introduction and review of the AI Assessment Catalog

The integration of Artificial Intelligence (AI) into weather and climate forecasting is reshaping how predictions are generated and applied. While AI-based models offer improvements in accuracy, efficiency, and accessibility, they also introduce new ethical and legal challenges, especially in high-stakes domains such as disaster prevention and climate adaptation.

These **guidelines are primarily intended for developers** of AI-based forecasting systems. They aim to support responsible design and deployment, ensuring that ethical considerations and regulatory requirements are addressed. Since downstream providers may rely on documentation and transparency measures provided by ECMWF, the guidelines also promote consistent communication across the AI value chain by incorporating aspects of documentation and communication for downstream suppliers (while still being targeted at developers).

The structure of the guidelines is based on the **Guidelines for Designing Trustworthy Artificial Intelligence (AI Assessment Catalog)** [1], which provides a systematic framework for assessing trustworthiness across the AI lifecycle. This framework consists of three steps:

### 1. Protection Needs Analysis

The intended use of the system is defined, and its potential impact is assessed across trustworthiness dimensions. For AI-based weather forecasting, the most critical dimensions are **reliability, transparency, security, and fairness**.

### 2. Definition of Goals and Metrics

Based on the protection needs, measurable objectives are defined—e.g. standards for documentation, interpretability, availability, fairness, and accuracy.

### 3. Mitigation

Technical and organizational measures are derived to address identified risks. Each mitigation is documented to ensure traceability and support legal compliance. Comparing pre- and post-mitigation criteria enables a structured trustworthiness assessment.

To ensure usability for developers, the guidelines are presented as a list of lifecycle-specific considerations. While the full process is detailed in [1], this document focuses on lifecycle stages and trustworthiness dimensions, adapting key elements of goal setting and mitigation while omitting detailed risk area structures. Note that the Fraunhofer AI Assessment Catalog is a much broader document allowing AI assessments across domains, while the focus of these guidelines is not a structured assessment of an AI system, but rather guidance on ethical aspects to consider in the development process. As such, these guidelines are not suitable for a complete view of trustworthiness but are restricted to a concise range of important ethical focus areas.

Each item follows the naming convention **[Dimension abbreviation – Lifecycle stage abbreviation – number]** and is marked as either a **general requirement** by **[GENERAL]**, applicable to all models, including NWP or an **AI-specific requirement** by **[AI-SPECIFIC]**, relevant only or primarily to machine learning-based systems. Note that the guidelines are not intending to cover all ethical aspects of NWP-based forecasting systems but clearly focus on AI-based forecasting. The purpose of this categorization is instead to make immediately

visible, which ethical issues are new in data-driven forecasting or exhibit new aspects in data-driven forecasting and which have existed with NWP-based forecasting before<sup>1</sup>.

Additionally, requirements are categorized by priority into **MUST (mandatory)**, **SHOULD (recommended)** and **COULD (optional or context-dependent)**, to distinguish between essential and desirable criteria.

### 3 Guidance for trustworthiness of data

Trustworthy training and inference data is the foundation of building trustworthy models and AI applications. Many characteristics of forecasting models are a direct result of their training data. This stage includes data assimilation – if an AI-based data assimilation system is built from scratch however, it is advisable to also consider aspects of the model and operation stages.

#### Reliability of Data Sets and Data Assimilation

##### Validation of data sources

- [RE-D-1] [GENERAL] Data sources **must** be reviewed for trustworthiness, e.g. by considering documentation provided by weather stations, available metadata and information on whether alterations have been made to weather station data, i.e. outlier removal
- [RE-D-2] [AI-SPECIFIC] It **should** be verified and documented that the data sufficiently cover the intended application domain, e.g. with regards to geographic coverage and diversity

##### Validation of data sets

- [RE-D-3] [AI-SPECIFIC] The impact of errors in training data on model performance **should** be quantified, e.g. by explainable AI or UQ techniques
- [RE-D-4] [AI-SPECIFIC] Tests for physical plausibility of training and inference data **should** be conducted, see e.g. [3]
- [RE-D-5] [GENERAL] Common tests for data consistency, outlier/missing value/double value detection, data distributions and data coverage **must** be conducted [4]
- [RE-D-6] [AI-SPECIFIC] For rare event detection, synthetic data **could** be considered to increase occurrence in training data

#### Fairness of Data Sets and Data Assimilation

##### Measuring biases and subgroup testing

- [FA-D-1] [AI-SPECIFIC] Datasets **should** be evaluated to ensure representativeness of data, e.g. by considering proportions of training data among subgroups
- [FA-D-2] [AI-SPECIFIC] If inequalities are noticed, mitigations to make the dataset more representative **could** be considered, e.g. by adding data from further (reliable) sources or, if those sources do not exist, by adding synthetic data

---

<sup>1</sup> Examples of new ethical aspects arising in data-driven forecasting include drift monitoring or questions of fairness of data sets. While validation procedures for these cases might be similar to the ones already being performed in NWP-based forecasting, the aspect of biases and drifts directly influencing future model behavior is new.

- [FA-D-3] [AI-SPECIFIC] Known biases in the training data and their potential downstream effects, especially in applications involving risk-sensitive decisions (e.g. flood warnings, agricultural planning) **should** be documented

## Security of Data Sets and Data Assimilation

### Access management

- [SE-D-1] [GENERAL] The (editing) access to data sets used for training and infrastructure **must** be restricted to developing teams requiring it
- [SE-D-2] [GENERAL] Changes to data sets **must** be logged and traceable for developers and potentially auditors
- [SE-D-3] [GENERAL] Protocols for unauthorized access and/or changes to the datasets **should** be established including reporting of incidents
- [SE-D-4] [AI-SPECIFIC] Outlier detection **should** be conducted in the latest version of the dataset before training

### Security of infrastructure

- [SE-D-5] [GENERAL] Transmission and processing of data **must** be safeguarded and encrypted
- [SE-D-6] [GENERAL] It **must** be ensured that infrastructure components are protected through strict access control and continuous monitoring
- [SE-D-7] [GENERAL] Regular backups of critical datasets and configurations **must** be implemented to ensure data availability and recovery in case of system failure or security incidents
- [SE-D-7] [GENERAL] Security certifications (i.e. ISO 27001, SOC 2) **should** be considered to evaluate and demonstrate security of infrastructure

## Transparency for Data Sets and Documentation for Downstream Providers

### Reproducibility of datasets

- [TR-D-1] [AI-SPECIFIC] In case own or modified datasets were used for training, reusability of said datasets **should** be addressed, at least by documenting the lineage, pre-processing steps, measures and metadata (see i.e. FAIR principles [5])
- [TR-D-2] [AI-SPECIFIC] Measures taken to automate changes to datasets **should** take into consideration stochasticity and hardware-dependence of the implementations

### Documentation with regards to data sets

- [TR-D-3] [GENERAL] Documentation on datasets **should** be kept up-to date and include at least data sources, data validation procedures, provenance and results, fairness validation procedures and results and available meta-data. A standardized template for documentation is advisable (e.g. Data Sheets [6]).

### Metadata on data sets

- [TR-D-4] [GENERAL] Standardized templates for metadata generation on used datasets **should** be developed and used (consider e.g. [7])

## 4 Guidance for trustworthiness during the training process and for the AI model

In training, model calibration and validation new challenges arise. Choice of architectures and optimization algorithms as well as uncertainty quantification and explainability methods impact the trustworthiness of output, while documentation and reproducibility can support other researchers and downstream providers.

### Reliability of Training/Model

#### Design phase considerations

- [RE-M-1] [AI-SPECIFIC] Trade-offs between different models and training algorithms **should** be documented as well as a rationale for the final choice [14], e.g. via Decision Protocols

#### Validation of model output

- [RE-M-2] [GENERAL] The model output **must** be evaluated against real-world data and other forecasts
- [RE-M-3] [GENERAL] A variety of metrics and trade-offs between metrics [8] **must** be considered
- [RE-M-4] [GENERAL] Validation **must** include comparisons with traditional NWP systems (e.g. IFS) and reanalysis datasets (e.g. ERA5 data not included in training)
- [RE-M-5] [AI-SPECIFIC] Tests for physical plausibility of model output **should** be conducted, see e.g. [3]

#### Robustness of model outputs and uncertainty quantification

- [RE-M-6] [GENERAL] Uncertainty quantification methods (e.g. (deep) ensembles, Bayesian networks, Variational-Autoencoder) **should** be integrated and validated [9]
- [RE-M-7] [GENERAL] Choice of uncertainty quantification method **should** be documented with reasoning
- [RE-M-8] [GENERAL] During the development process of models – if possible - it **could** be advantageous to disentangle aleatoric and epistemic uncertainty, for example to support model optimization by accurately capturing and understanding epistemic uncertainty

### Fairness of Training/Model

#### Validation of fairness

- [FA-M-1] [GENERAL] Performance evaluation **should** include testing on subgroups and regional testing to detect potential performance differences
- [FA-M-2] [AI-SPECIFIC] If performance differences are recorded, these **should** be analysed and documented with respect to training data and differing factors
- [FA-M-3] [AI-SPECIFIC] If possible, mitigations to improve performance in lower-performing subgroups **could** be aimed at, e.g. by means of improving data quality or amount of data from these subgroups or the addition of synthetic data, if there are no other data sources available

### Security of Training/Model

#### Adversarial robustness

- [SE-M-1] [AI-SPECIFIC] Approaches to improve adversarial robustness (e.g. adversarial regularization) **could** be considered depending on the trade-off with performance loss

## Transparency of Training/Model and Documentation for Downstream Providers

### Security of infrastructure

- [SE-M-2] [GENERAL] In order to test security of infrastructure and by extension model and model output, penetration tests and/or red teaming **could** be considered as approaches to find weaknesses
- [SE-M-3] [AI-SPECIFIC] It **must** be ensured that regular backups of training data and configurations are made to support recovery

### Design phase considerations

- [TR-M-1] [AI-SPECIFIC] System architecture, model choice, choice of training algorithm, hyperparameter tuning process and safeguarding measures by design **should** be documented with rationale, e.g. via Model Cards [10] or C4 Models [13]
- [TR-M-2] [AI-SPECIFIC] The intended use cases **should** clearly be defined and documented, e.g. via Model Cards [10]

### Reproducibility of models

- [TR-M-3] [AI-SPECIFIC] Model reproducibility **should** be considered where possible and worthwhile, e.g. by containerization of models
- [TR-M-4] [AI-SPECIFIC] Reproducible implementations in common ML-frameworks **could** be considered depending on performance loss through said implementation
- [TR-M-5] [AI-SPECIFIC] Versioning and documentation for both training data and models **must** be ensured, as they are paramount for traceability, reproducibility, and quality assurance.

### Explainability of models

- [TR-M-6] [AI-SPECIFIC] Methods of explainable AI, such as SHAP, LIME, CAM, or LRP, **could** be considered depending on the concrete use case, in which case the choice and reasoning **should** be documented

### Documentation of model training

- [TR-M-7] [AI-SPECIFIC] All training steps, including algorithm choice, training data and preprocessing, hyperparameter tuning, and optimization and safeguarding strategies **must** be documented, e.g. via Model Cards [10]
- [TR-M-8] [AI-SPECIFIC] Training logs and configuration files **should** be versioned and stored to support reproducibility and auditing

### Generation of meta-data

- [TR-M-9] [AI-SPECIFIC] Metadata **should** include information like dataset versions, preprocessing steps, model versions, etc.

### Documentation of validation procedures

- [TR-M-10] [AI-SPECIFIC] Documentation of model validation and metrics used **must** be documented in a standardized manner, e.g. via Model Cards [10]
- [TR-M-11] [AI-SPECIFIC] Trade-offs between different validation metrics **should** be documented, with results of varied metrics considered as well as comparisons to benchmark forecasting models [14]- e.g. via Decision Protocols

## 5 Guidance for trustworthiness during deployment and operation

Operative models are subject to potential changes, both in data distributions underlying and new external requirements. Validation procedures and documentation measures should be chosen accordingly to assure ethical and trustworthiness after deployment.

### Reliability in Operation

#### Validation of inference data

- [RE-O-1] [GENERAL] The distribution of inference data **should** be monitored, e.g. using statistical tests to determine distribution differences (Kolmogorov-Smirnov-Tests, Population stability index, Wasserstein metrics, detecting ML-models, ... [11], [12])

#### Validation of model output

- [RE-O-2] [AI-SPECIFIC] In the case of continuously learning models, tests for model drift **should** be conducted to detect potential changes in models' behaviour by monitoring performance metrics over time
- [RE-O-3] [AI-SPECIFIC] Upon findings of degradation of model performance, the (re-)training data **must** be revisited
- [RE-O-4] [GENERAL] Relevant accuracy, robustness, and uncertainty metrics **should** be continuously monitored, evaluated, and documented

#### Robustness of model outputs and uncertainty quantification

- [RE-O-5] [GENERAL] Uncertainty quantification performance **should** be monitored continuously or in fixed time intervals
- [RE-O-6] [AI-SPECIFIC] If automated retraining is part of the AI system, appropriate triggers **should** be defined and documented, either continuously, in fixed time intervals or following drift detection or performance degradation thresholds

### Fairness in Operation

#### Validation of fairness

- [FA-O-1] [AI-SPECIFIC] Fairness metrics on inference and output data **could** be monitored constantly or in regular time intervals to detect disparities in forecast accuracy e.g. across geographic regions and climate regimes
- [FA-O-2] [AI-SPECIFIC] Upon findings of newly arising performance inequalities, the (re-)training dataset **should** be reviewed and, if necessary, be adapted to address the identified disparities
- [FA-O-3] [AI-SPECIFIC] Known biases in inference data and/or output and their potential impact on vulnerable regions **should** be documented

### Security in Operation

#### Adversarial robustness

- [SE-O-1] [AI-SPECIFIC] Inference data **should** be subject to constant monitoring processes aimed at detecting adversarial attacks, including checks for physical plausibility
- [SE-O-2] [AI-SPECIFIC] Presence and handling of outliers **should** be documented and – if necessary – reviewed and reported to security responsables
- [SE-O-3] [AI-SPECIFIC] Red-teaming exercises **could** be conducted to simulate adversarial scenarios and identify vulnerabilities

## Security of infrastructure

- [SE-O-4] [AI-SPECIFIC] Access to data used to (re-)train or run the model while in operation **must** be restricted<sup>2</sup>
- [SE-O-5] [GENERAL] Transmission of inference data **must** be encrypted and secured
- [SE-O-6] [GENERAL] Logging or traceability measures for editing of training or inference data **should** be implemented
- [SE-O-7] [GENERAL] In order to test the security of the infrastructure, penetration tests and/or red-teaming **could** be considered as approaches to find weaknesses
- [SE-O-8] [GENERAL] Security certifications (i.e. ISO 27001, SOC 2) **should** be considered to evaluate and demonstrate security of infrastructure

## Fallbacks and availability

- [SE-O-9] [GENERAL] In the case of (some of the) inference data not being readily available, protocols to maintain forecast availability **should** be available and followed

## Transparency in Operation and Documentation for Downstream Providers

### Documentation of validation procedures in operation

- [TR-O-1] [AI-SPECIFIC] Processes of validation in operation, monitoring of inference data distributions and model output **must** be documented, ideally in a standardized and automated manner
- [TR-O-2] [AI-SPECIFIC] For the purpose of downstream compliance and risk assessment, retraining triggers **should** be specified and documented
- [TR-O-3] [AI-SPECIFIC] Metadata regarding inference data and model **should** be produced on a continuous basis regarding new inference data and forecasts
- [TR-O-4] [AI-SPECIFIC] Testing and operating conditions, (e.g. hardware), should be documented and communicated, including eventual hardware dependencies in code that might impact reproducibility on different hardware

### Versioning changes

- [TR-O-4] [GENERAL] Upon releasing new versions, release notes **must** be drafted, including recent changes and known implications of these changes, e.g. differences in training data or model behaviour
- [TR-O-5] [GENERAL] If significant differences are noticed, the notes **should** contain comparisons of evaluations of previous and current version as well as changes in other documented metadata
- [TR-O-6] [GENERAL] Notification to downstream suppliers **should** contain impact assessments for changes in API access and required hardware to run models
- [TR-O-7] [GENERAL] Dependency registries (i.e. structured lists of downstream providers with model versions used) **should** be held and kept up to date

---

<sup>2</sup> This requirement aligns with the ISO/IEC 27001 standard, which includes controls for access management to sensitive information and systems as part of an Information Security Management System (ISMS). Implementing ISO/IEC 27001 can help ensure that access to operational and training data is appropriately restricted and monitored.

## Sources

- [1] M. Poretschkin *et al.*, "Guideline for Trustworthy Artificial Intelligence -- AI Assessment Catalog," 2023, *arXiv*. doi: 10.48550/ARXIV.2307.03681.
- [2] Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (Artificial Intelligence Act) [2024] OJ L2024/1689
- [3] G. J. Hakim and S. Masanam, "Dynamical Tests of a Deep Learning Weather Prediction Model," *Artificial Intelligence for the Earth Systems*, vol. 3, no. 3, Jul. 2024, doi: 10.1175/aies-d-23-0090.1.
- [4] Y. Zhou, F. Tu, K. Sha, J. Ding, and H. Chen, "A Survey on Data Quality Dimensions and Tools for Machine Learning," Jun. 28, 2024, *arXiv*: arXiv:2406.19614. doi: 10.48550/arXiv.2406.19614.
- [5] M. D. Wilkinson *et al.*, "The FAIR Guiding Principles for scientific data management and stewardship," *Sci Data*, vol. 3, no. 1, p. 160018, Mar. 2016, doi: 10.1038/sdata.2016.18.
- [6] T. Gebru *et al.*, "Datasheets for Datasets," Dec. 01, 2021, *arXiv*: arXiv:1803.09010. doi: 10.48550/arXiv.1803.09010.
- [7] B.-I. Manola *et al.*, "Best Practice Guidelines for Climate Data and Metadata Formatting, Quality Control and Submission," 2020, doi: 10.24381/KCTK-8J22.
- [8] Zied Ben Bouallègue and ECMWF AIFS Team, "Accuracy versus activity," *ECMWF Newsletter*, 2024, doi: 10.21957/8B50609A0F.
- [9] J. Gawlikowski *et al.*, "A survey of uncertainty in deep neural networks," *Artif Intell Rev*, vol. 56, no. S1, pp. 1513–1589, Oct. 2023, doi: 10.1007/s10462-023-10562-9.
- [10] M. Mitchell *et al.*, "Model Cards for Model Reporting," in *Proceedings of the Conference on Fairness, Accountability, and Transparency*, Atlanta GA USA: ACM, Jan. 2019, pp. 220–229. doi: 10.1145/3287560.3287596.
- [11] S. Arora, R. Rani, and N. Saxena, "A systematic review on detection and adaptation of concept drift in streaming data using machine learning techniques," *WIREs Data Min & Knowl*, vol. 14, no. 4, p. e1536, Jul. 2024, doi: 10.1002/widm.1536.
- [12] F. Hinder, V. Vaquet, and B. Hammer, "One or Two Things We know about Concept Drift -- A Survey on Monitoring Evolving Environments," Oct. 24, 2023, *arXiv*: arXiv:2310.15826. doi: 10.48550/arXiv.2310.15826.
- [13] Brown, S. (n.d.). The C4 model for visualising software architecture. Retrieved September 3, 2025, from <https://c4model.com>
- [14] Chmielinski, K. et al. (2024). The CLeAR Documentation Framework for AI Transparency: Recommendations for Practitioners & Context for Policymakers. Shorenstein Center on Media, Politics and Public Policy, Harvard Kennedy School.

## Annex I: Ethical principles for the development of AI systems in weather forecasting

Data-driven weather forecasting introduces a series of new ethical challenges, often due to the black-box nature of models and the inherent dependency of model's behaviours and their training data. The purpose of this document is to lay out basic ethical considerations in the development of AI systems in weather forecasting in a concise manner. For structured and in-depth information consider the practical guidelines.

### **Reliability: How can we measure and ensure that the AI system is making reliable predictions?**

The impact of weather forecasting systems on downstream systems as well as the reliance of weather services on its outputs make reliability a paramount aspect of trustworthiness. Broad testing across different metrics and implementation of uncertainty quantification measures are crucial. In machine learning models, performance often depends on underlying training data, making thorough evaluation of datasets necessary.

### **Transparency: How can AI systems and results be made more transparent for stakeholder?**

When developing an AI system for weather forecasting, it is essential to document every stage — including datasets, code, modelling, evaluation, deployment, as well as monitoring and logging during operation. This not only supports downstream providers and internal developers in building their own applications, but also makes it easier to fix errors, conduct audits, and reproduce results. Where possible, additional technical measures such as explainability methods or choice of reproducible algorithms can help manage the black-box character of AI models. Transparency as a guiding principle can strengthen trust in forecasts, raise acceptance among decision makers and support adoption of this new technology.

### **Fairness: How can we make sure that the output of AI systems is not discriminating?**

Bias in training datasets or uneven coverage can impact the models themselves and lead to performance differences across regions. While this is often the case due to differing density of measurements, having the issue in mind and evaluating datasets for fairness can still be worthwhile. On the one hand, it allows to by consider technical mitigation, e.g. more data from other sources and on the other hand knowing and communicating potential weaknesses can generate trust and support decision makers.

### **Safety & Security: What can be done to keep models secure?**

The public and political relevance, particularly of extreme weather and climate predictions, make forecasting systems a potential target of adversarial attacks. This makes securing model training, monitoring forecasting systems and access management crucial. While a lot of security measures are not AI-specific, awareness for certain types of adversarial attacks is important to take preventive measures.

# Annex II: Visualization of the structure of the guidelines

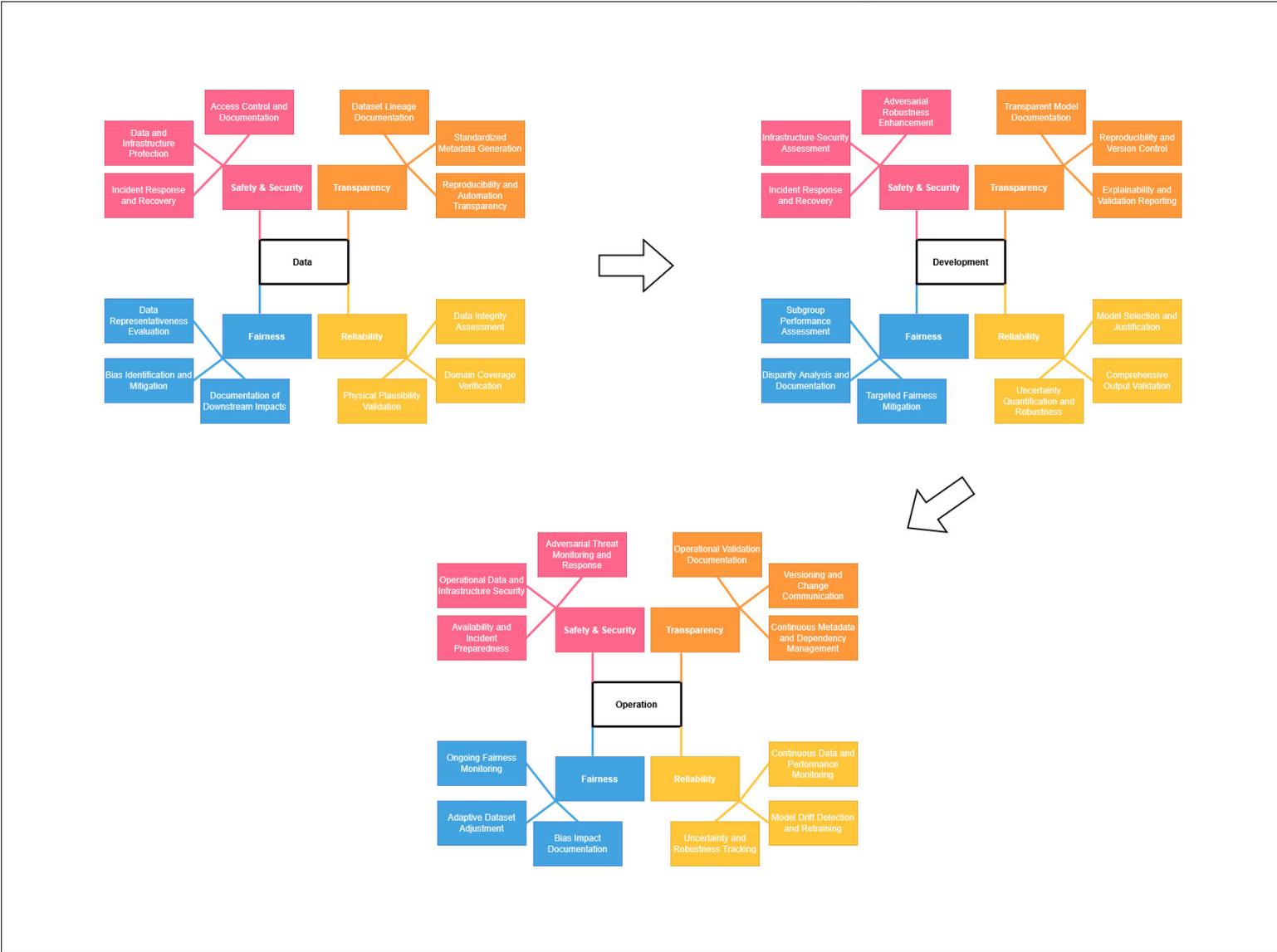


Fig. 1: Visualization of Areas of Recommendations by lifecycle stages – the four trustworthiness dimensions that were identified are linked to the most relevant examples for each dimension and lifecycle stage.

*This document has been produced in the context of the Destination Earth Initiative and relates to tasks entrusted by the European Union to the European Centre for Medium-Range Weather Forecasts implementing part of this Initiative.*

*This document is funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the European Commission can be held responsible for them. The European Centre for Medium-Range Weather Forecasts is not liable in respect of this document and gives no warranty for the information provided.*

UK: (Headquarters) ECMWF, Shinfield  
Park, Shinfield Road, Reading,  
RG2 9AX, UK

Italy: ECMWF, Tecnopolo di Bologna,  
Via Stalingrado 84/3, 40128 Bologna,  
Italia

Germany: ECMWF, Robert-Schuman-  
Platz 3, 53175 Bonn, Deutschland