

# Trust in AI, trust in public processes



S. Daniil

S. Vethman

M. Molhoek

## OUR OBJECTIVE

*To enable trustworthy use of AI by public authorities and support effective oversight.*

The use of AI tooling can facilitate the design and application of efficient and socially aware policies by governmental bodies, but involves the risk of disregarding *trustworthiness* as a priority, thus risking losing the trust of society towards the government. Guidelines for Trustworthy AI have been proposed, but require practical understanding from developers of AI systems in order to be adequately implemented.

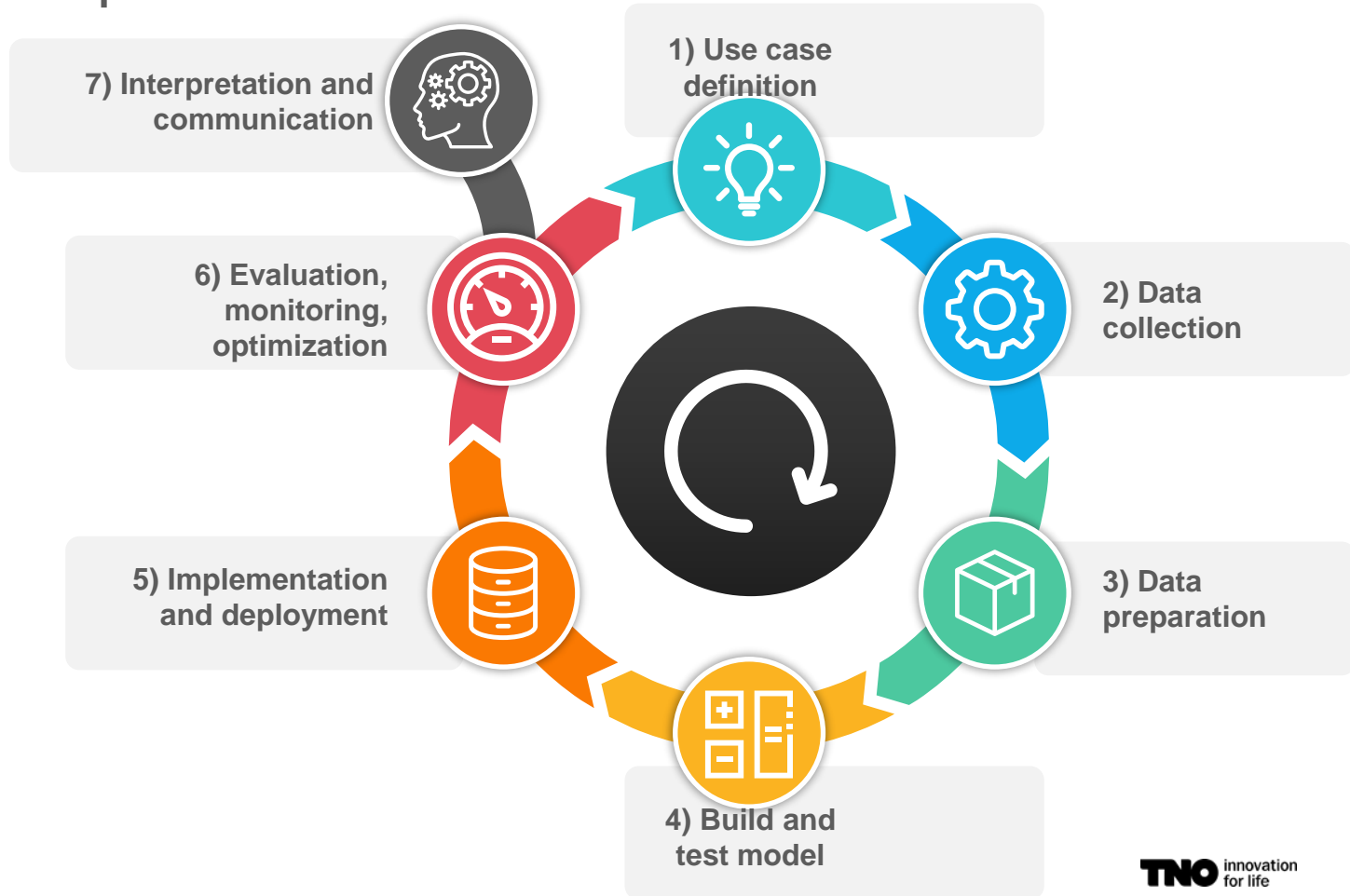
# Menu

|  |  |
|--|--|
| <b>AI life cycle and trustworthy AI principles</b> | <ul style="list-style-type: none"><li>• AI lifecycle</li><li>• Overview of important Trustworthy AI principles</li><li>• Selection of principle Reproducibility</li></ul>  |
| <b>Reproducibility</b>                             | <ul style="list-style-type: none"><li>• Importance of Reproducibility</li><li>• Views on Reproducibility</li><li>• Required Randomness</li></ul>   |
| <b>Handbook</b>                                    | <ul style="list-style-type: none"><li>• Reproducibility framework mapped to the AI lifecycle</li><li>• Best practices for implementation<ul style="list-style-type: none"><li>• Step 3: Data pre-processing</li><li>• Step 4: Build and test model</li></ul></li></ul> |
| <b>Other causes</b>                                | <ul style="list-style-type: none"><li>• Inherent randomness in hardware</li></ul>  |
| <b>Next steps</b>                                  | <ul style="list-style-type: none"><li>• Vision of the future of AI Oversight Lab handbook</li></ul>  |

# AI life cycle components

The AI development cycle passes through 7 stages. From project scoping (1-2), design and build phase (3-4) up to continuous improvement in production (5-6).

The 7<sup>th</sup> stage is specifically added as interpretation and communication by operators is society-dependent and can change over time.

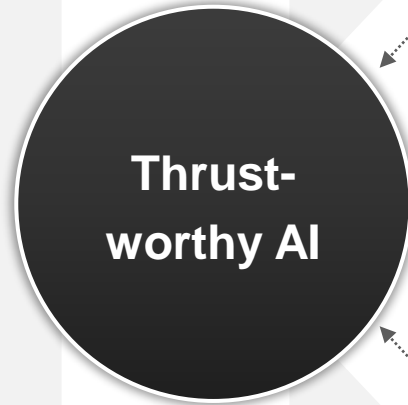


# Trustworthy AI principles:

The European ethics guidelines on trustworthy AI state that it should be:

- Lawful – respecting all applicable laws and regulations
- Ethical - respecting ethical principles and values
- Robust – both from a technical perspective while taking into account its social environment

These guidelines translate into 7 key requirements that AI systems should meet, the principles of trustworthy AI.



Human agency & Oversight

Technical Robustness & Safety

Privacy & Data Governance

Transparency

Non-discrimination & fairness

Societal & environmental well-being

Accountability

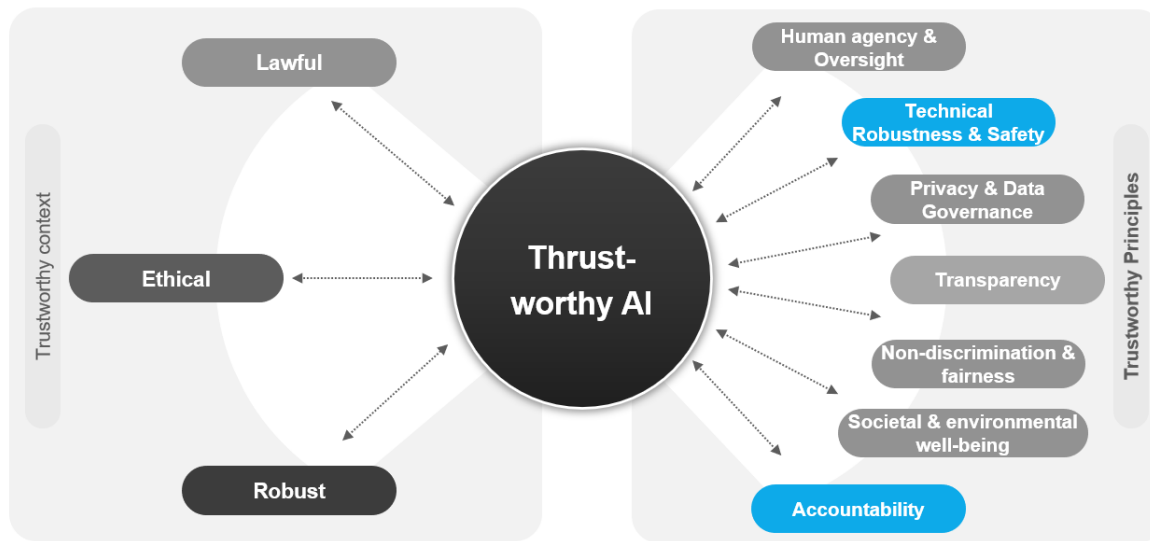
Trustworthy principles

# Begin with ... Reproducibility

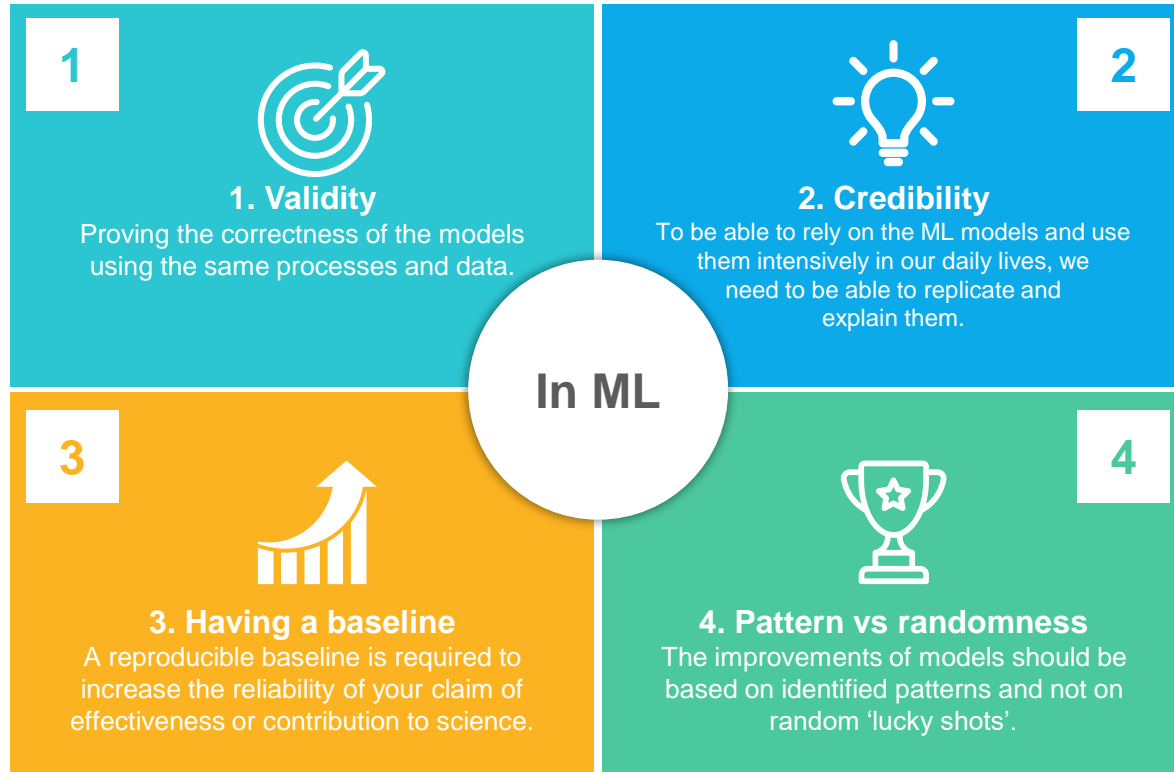
*Reproducibility* serves as the foundation for the Trustworthy principles *Robustness* and *Accountability*. Its relevance is often underestimated and has not been thoroughly investigated in the context of AI, despite the crisis that it seems to face.

Reproducibility of AI systems highly influences their performance and applicability in social context, since:

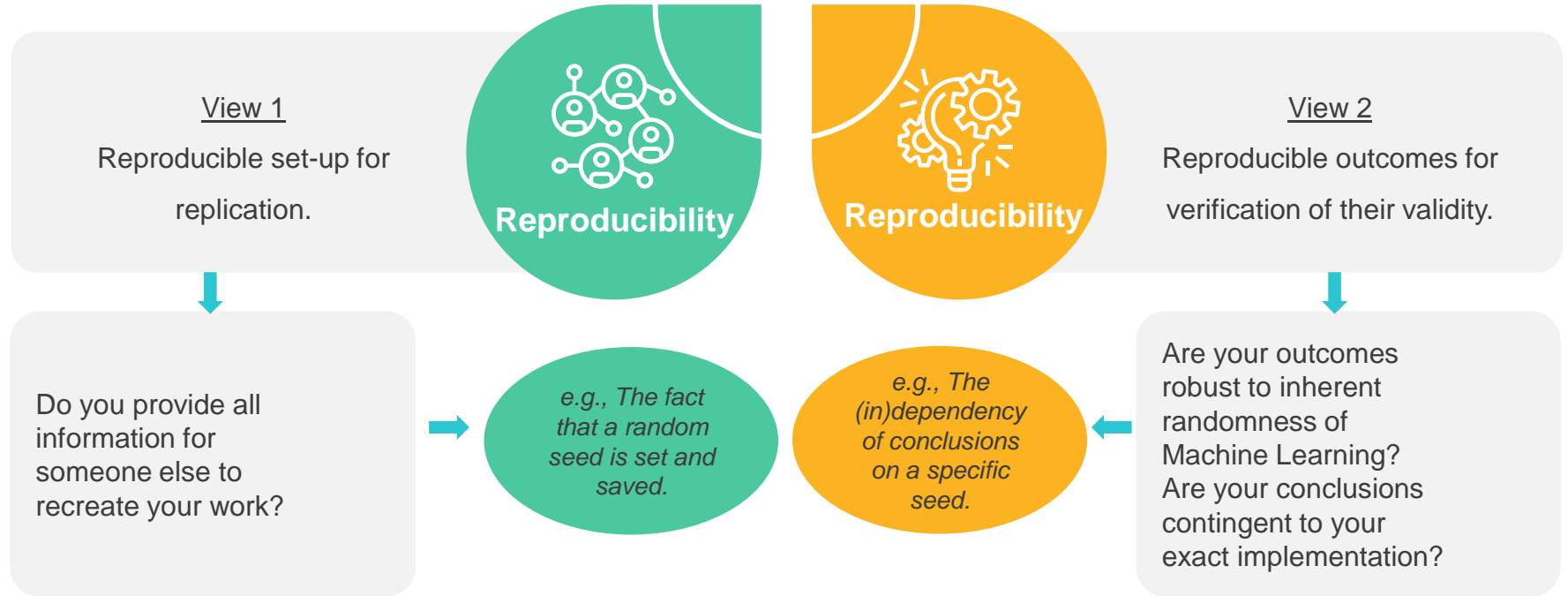
1. It is a prerequisite for technical robustness, and
2. It is crucial in rendering systems verifiable and thus fit for criticism and improvement.



# The importance of reproducibility



# What is reproducibility?



# The reproducible scale

There are various strictness levels of reproducibility for another person to get to the same results. All variations can occur in between. What is the reproducibility if different hardware is used? What is the effect of coding in R versus in Python? What is the result of different library versions? How do random forest and gradient boost affect reproducibility? Depending on the need of the user, the strictness of reproducibility can be chosen.



## Reproducibility

With **the same** hardware,  
the **original** code,  
the **same** language,  
the **same** AI methods,  
the result is reproducible for  
another person.



## Reproducibility

With **various** hardware,  
**no original** code available,  
**different** languages,  
**comparable** AI methods,  
the result is reproducible for  
another person





# Reproducibility: Complementing views

Easy to reproduce

Hard to reproduce

View 1: Replicability

## REQUIRES VALIDATION

The use of the algorithm is **accountable**, as the settings and outcomes can be repeated such that previous execution can be corroborated.

However, the outcome **may be dependent on the exact implementation** of the code, **more validation is needed** before one considers the outcomes of the algorithm in process affecting society.

## IDEAL

The algorithm is **accountable** as information for reproducibility for all stages in the AI life cycles is appropriately stored.

The **outcomes are verifiable or are verified** with different models and settings such that they can be **interpreted with appropriate level of confidence**.

## NO GO

Outcomes of the algorithm cannot be reproduced by any (3<sup>rd</sup>) party. Randomness in the algorithm or contingencies of the algorithm are ignored in the interpretation of the outcomes. Hence, application of the algorithm is in many real life processes ill-advised as **the outcomes do not have a robust interpretation nor are they accountable**.

## REQUIRES TRUST

Outcomes of the algorithm cannot be reproduced by 3<sup>rd</sup> parties due to omission of information on parameters and packages or privacy aspects of even the characteristics of the data.

The outcomes are presented to be verified by multiple approaches and classifiers, however **due to the inability to exactly replicate the outcomes, the developer need be trusted**.

View 2: Robustness/Stability/Reliable

Contingent

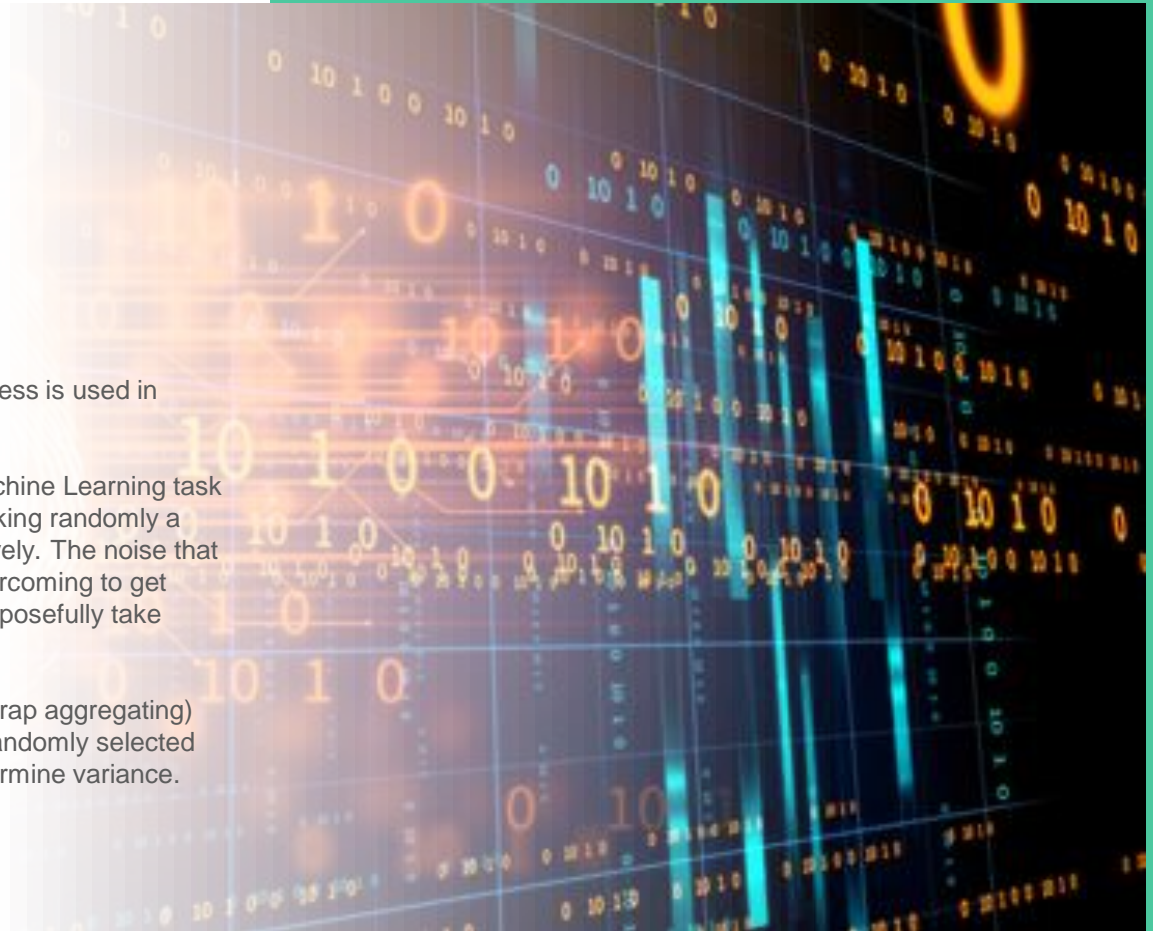
Verified/Verifiable

# Importance of randomness

Due to memory and time constraints randomness is used in machine learning tasks.

To calculate the gradient in optimizing the Machine Learning task often Stochastic Gradient Descent is used picking randomly a part of the data to update its parameter iteratively. The noise that is added during these updates can help in overcoming to get stuck in a local minimum. The model may purposefully take random steps to seek a better state.

Another example is the use of bagging (bootstrap aggregating) which trains multiple models on overlapping randomly selected subsets of data to increase accuracy and determine variance.



# TNO's focus in the landscape

## Existing frameworks/tools

EU ethics guidelines and frameworks WP1



TILT handbook

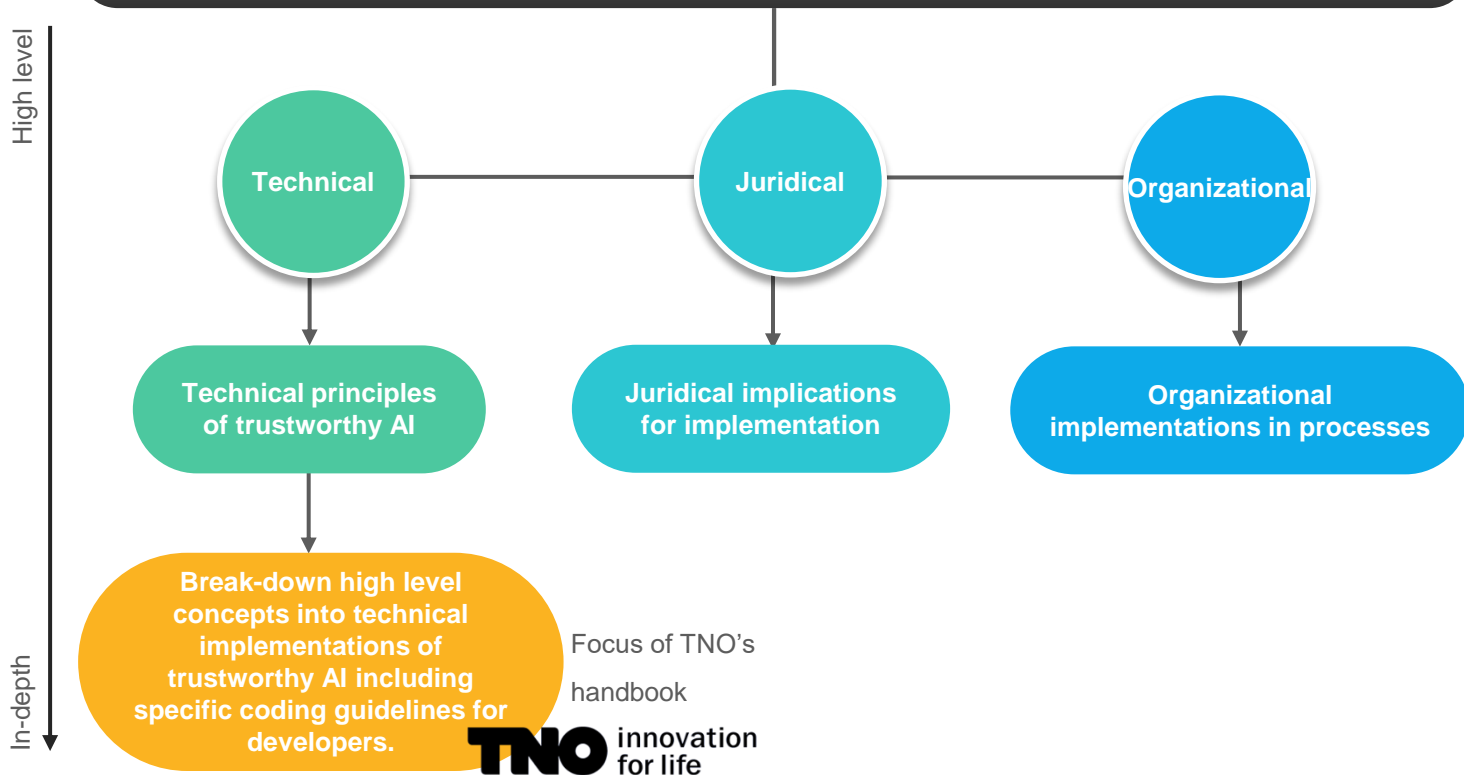


Trustworthy AI Questionnaire



V29 questionnaire

## Trustworthy AI guidelines and handbooks



# AI life cycle

This step includes:

## AI Life cycle components

**1) Use case definition:** setting the goal and the necessity of the project, define envisioned impact, set the required assumptions for this goal, define end-user and define success criteria.

**2) Data collection:** define data sources, protocols to safely collect and store data and handle boundary conditions for collecting the data such as juridical implications.

**3) Data pre-processing:** data cleansing, instances selection and partitioning, feature tuning, representation transformation, feature extraction, feature selection, feature construction, coupling of datasets and data labelling. Starts with raw data and ends up with a ML ready dataset.

**4) Build and Test model:** model(s) selection, train and test split, testing and evaluating the model on pre-defined success criteria.

**5) Implementation and deployment:** practical tests first in a sandbox environment, model adjustment and applying in process, document restrictions and conditions for use.

**6) Evaluation, monitoring optimization:** define implementation strategy, evaluation on defined goal/impact, actions to improve, continues monitoring and improvement.

**7) Interpretation and communication:** define user interpretation and expected actions, evaluate universality of interpretation, monitor and act upon changes in actors/environment acting on outcome, adjust process accordingly.

# Framework Reproducibility

|       |                  | Life cycle components   |  |  |  |   |   |  |
|-------|------------------|---|--|--|--|---|---|--|
|       |                  | Use case  | Data collection  | Data preparation   | Build and test model   | Implementation and deployment   | Evaluation, monitoring optimization   | Interpretation and communication   |
| Level | Definition       | <ul style="list-style-type: none"> <li>• Problem statement</li> <li>• Define goal and scope</li> <li>• Define intended use</li> </ul> | <ul style="list-style-type: none"> <li>• Define data sources</li> <li>• Procedure and boundaries of data collection</li> </ul> | <ul style="list-style-type: none"> <li>• From raw data to ML input</li> <li>• Data cleaning and preprocessing</li> </ul> | <ul style="list-style-type: none"> <li>• Model selection</li> <li>• Test procedures including train test split.</li> </ul>   | <ul style="list-style-type: none"> <li>• Sandbox first</li> <li>• Practical application</li> <li>• Restrictions and conditions for use</li> </ul> | <ul style="list-style-type: none"> <li>• Evaluation with defined goal</li> <li>• Actions to improve</li> <li>• Continuous monitoring</li> </ul> | <ul style="list-style-type: none"> <li>• User interpretation and expected actions.</li> </ul>              |
|       | Problems sources | <ul style="list-style-type: none"> <li>• What level and type of reproducibility is necessary?</li> </ul>                              | <ul style="list-style-type: none"> <li>• Data size</li> <li>• Computing power requirements</li> </ul>                          | <ul style="list-style-type: none"> <li>• Dynamic data</li> <li>• Feature engineering</li> <li>• Missing data</li> </ul>  | <ul style="list-style-type: none"> <li>• Random sampling</li> <li>• Parameter choices</li> <li>• Inherent randomness</li> <li>• Libraries</li> </ul>                         | <ul style="list-style-type: none"> <li>• Documentation by user</li> </ul>   | <ul style="list-style-type: none"> <li>• Metrics</li> </ul>   | <ul style="list-style-type: none"> <li>• Visualization</li> <li>• Non-universal standardization</li> </ul> |
|       | Solutions        | <ul style="list-style-type: none"> <li>• Use our handbook to select your needs for your use case.</li> </ul>                          | <ul style="list-style-type: none"> <li>• Learning Curve</li> <li>• Benchmark set</li> </ul>                                    | <ul style="list-style-type: none"> <li>• Timestamps</li> <li>• Recording of preliminary data sets.</li> </ul>            | <ul style="list-style-type: none"> <li>• Documentation of seed.</li> <li>• Version control.</li> <li>• Automated testing</li> <li>• Environment and docker files.</li> </ul> | <ul style="list-style-type: none"> <li>• Automated documentation.</li> </ul>  | <ul style="list-style-type: none"> <li>• Checklist</li> <li>• Dynamic dashboard.</li> </ul>   | <ul style="list-style-type: none"> <li>• Benchmark</li> <li>• Examples</li> </ul>                          |

# Build and test model:

|                  | Problems / Sources   | Solutions   |
|------------------|--|---|
| Train-test split | <ul style="list-style-type: none"><li>• <b>A random allocation of each data point or subject into two separate sets</b> is executed to split the data into a train or test set. This randomness may influence the results of the algorithm.<ul style="list-style-type: none"><li>• In python, this can e.g. be performed by a function of the package scikit learn called: <code>train_test_split</code>.</li></ul></li></ul>  | <ul style="list-style-type: none"><li>• <b>Setting and saving a seed</b> to the train test split ensures that creation of the train and test is repeatable.<ul style="list-style-type: none"><li>• In python, the package scikit learn often allows the random state to be given to the function.</li></ul></li><li>• <b>Cross validation</b> is a more extensive version on a train-test split with multiple splits per execution such that the influence of the randomness becomes smaller. This increases the reproducibility w.r.t. View 2 of the results as it decreases the reliance of a specific seed.</li><li>• <b>Saving (and sharing) the index of the train and test splits</b> (folds) increase the reproducibility of results w.r.t. a train-test split further. In this way, you can exactly recreate the splits irrespective of another random number generator and if applicable others may test the findings on the exact same data splits.</li></ul> |
| Parameter choice | <ul style="list-style-type: none"><li>• <b>Unknown parameter choices</b>, i.e. missing records of the settings used in an algorithmic application, are a risk for reproducibility.</li><li>• <b>Default parameters in packages</b> are an even larger risk for the reproducibility in terms of View 2. These are often unacknowledged and not recorded, while default parameters differ between (versions of) packages and programming languages.<ul style="list-style-type: none"><li>• In python, the default number of splits in function <code>StratifiedShuffleSplit</code> of package scikit learn has changed over versions from 3 to 5. E.g. the solver for <code>RandomForest</code> in python has a default depth of 100 branches while the package for R has a default of 200 branches.</li></ul></li></ul> | <ul style="list-style-type: none"><li>• <b>Configuration</b> files or scripts aid in the recording of the parameter settings applied in testing the model. This is in contrast to parameters being set somewhere in the collection of scripts that the eventual algorithm will consist of.<ul style="list-style-type: none"><li>• It is common practice to use <code>.yaml</code> files for saving parameter choices and data paths for referencing all relevant directories.</li></ul></li><li>• <b>Version control</b> is in turn able to save the configuration files and as such the chosen parameters of each test. For this solution there is a distinction for the developer (elaborated here) and the user (elaborated in the deployment phase).</li></ul>  |

# Build and test model

## Problems / Sources

## Solutions

### Inherent algorithmic randomness (non-deterministic)

- There are **many reasons for algorithms to be non-deterministic**, see slide 19 for elaboration on inherent randomness in algorithms. Some specific origins of randomness in machine learning algorithms are listed here:
  - Many models optimize **weights** such that they capture the targeted relation in the data. To acquire a starting point of the optimization, weights are often randomly initialized.
  - Multiple features of **Neural Networks** use random sampling, one of them is the random selection in the application of dropout layers.
  - For **stochastic optimization** to find (local) optima, random sampling is used. An example for which this holds is stochastic gradient descent.
  - Randomness in **random forest** models. Each tree in the random forest selects random features on which the model is trained.

### Library

- The use of **libraries** simplifies the use of AI and machine learning but often does not guarantee an 100% reproducibility. Libraries and the underlying functions change each version, including their dependencies to other libraries.

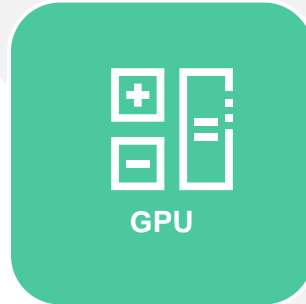
- **Test automation:** To gain the benefits from algorithmic randomness and limit the drawbacks, automatic testing of the algorithm for e.g. a range of seeds may be performed.
  - Do the results for different seeds fall into an acceptable range?
  - Are the conclusions the same for all most seeds?
- **Continuous integration:** A step further in test automation is the use of continuous integration of testing during the building and testing of the model. Before any version of the model is approved or even any merge from branch to main is requested, automated testing can be integrated. As such, version updates and merge request are accompanied with a testing report.
- **Alternative solver, same results?** In many applications, there are many options to solve the objective function, i.e. to optimize the classifier. To increase the reproducibility of outcomes w.r.t. View 2, it is advised to see the outcomes for multiple solvers with similar functionalities s.t. the dependency on one component is limited.
- **Statistical significance tests** between results and benchmark take into account random noise in data and if applied correctly contribute to robustness of a model's results.
- **Environment**, in python, it is common practice to create a designated virtual environment per project or per version of the model. These virtual environment constitute a.o. the python version and the versions of all installed packages. The creation and preservation of python environments enhance the reproducibility by ensuring execution of the code on similar software settings. The operation system is however not included in the virtual environment.
- **Docker file**, the operation system (Linux, Windows 2000) on which is directly related to whether certain algorithms can be executed. To increase reproducibility, a docker file, "a virtual operating system often containing a virtual environment" can be created and shared.

# Inherent randomness in Hardware

There is inherent randomness and therefore irreproducibility in hardware selection as well. Using different hardware types may affect the results.



Complex computations require Giga or even Tetra floating point operations. To decrease the timeframe and increase efficiency the processes are run in parallel on multiple central processing units and/or graphics processing units.



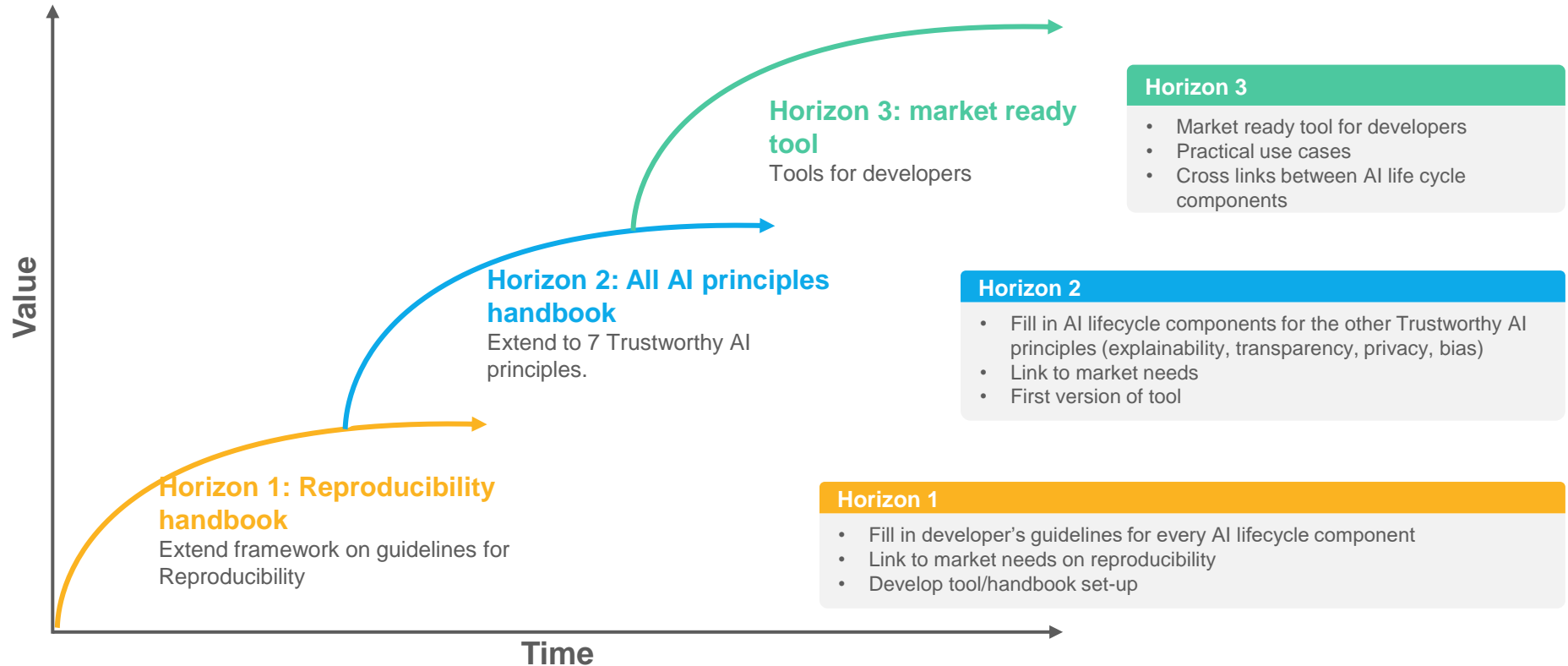
Uses Intra-Ops and Inter-ops parallelism on CPU, which can result in different results. By using floating points in solving the result differs.



Uses stream multiprocessing units which performs asynchronous computation which results in different runs.



# Next steps towards trustworthy AI



# Trust in AI, trust in public processes



S. Daniil

S. Vethman

M. Molhoek

## OUR OBJECTIVE

*To enable trustworthy use of AI by public authorities and support effective oversight.*

The use of AI tooling can facilitate the design and application of efficient and socially aware policies by governmental bodies, but involves the risk of disregarding *trustworthiness* as a priority, thus risking losing the trust of society towards the government. Guidelines for Trustworthy AI have been proposed, but require practical understanding from developers of AI systems in order to be adequately implemented.

# References

- [1] P. Hemant, "Reproducible machine learning," Apr 2020.
- [2] T. Allard, "10 top tips for reproducible machine learning," Apr 2020.
- [3] W. D. Heaven, "Ai is wrestling with a replication crisis," Nov 2020
- [4] "Reproducibility challenge @ neurips 2019.
- [5] P. Sugimura and F. Hartl, "Building a reproducible machine learning pipeline," 2018.
- [6] P. Warden, "The machine learning reproducibility crisis," Mar 2018.
- [7] jennifervilla, "Reproducibilityinml:whyitmattersandhowtoachieveit," May 2018.
- [8] J. Fjeld, N. Achten, H. Hilligoss, A. Nagy, and M. Srikumar, "Principled artificial intelligence: Mapping consensus in ethical and rights-based approaches top principles for ai," SSRN Electronic Journal, 2020.
- [9] B. Shneiderman, "Bridging the gap between ethics and practice," ACM Transactions on Interactive Intelligent Systems, vol. 10, no. 4, p. 1–31, 2020.
- [10] J. Brownlee, "Embrace randomness in machine learning," Aug 2019
- [11] S. Mall, "Realizing reproducible machine learning - with tensorflow," Dec 2019
- [12] C. Shao, "Properly setting the random seed in ml experiments. not as simple as you might imagine," April 2019.
- [13] C. Ayuya, "Reproducibility to improve machine learning," January 2021.
- [14] EU, "Ethics guidelines for trustworthy ai," April 2019.
- [15] N. Koleva, "Ai and data science lifecycle: Key steps and considerations," May2020.
- [16] Bart van der Sloot, Esther Keymolen, Merel Noorman, Yvette Wagenveld(Tilburg University), College voor de rechten van de Mens, Hilde Weerts (TU E) Bram Visser (VU Brussel) "Handreiking non-discriminatie by design," 2021
- [17] Interviews with NFI, VORtech, UvA, Leiden University