# AI4CYBER

**TRUSTWORTHY ARTIFICIAL INTELLIGENCE FOR CYBERSECURITY REINFORCEMENT AND SYSTEM RESILIENCE**

# Cybersecurity Datasets – Opportunities and Challenges

Erkuden Rios, Project Coordinator, **tecnalia**

MEMBER OF BASQUE RESEARCH & TECHNOLOGY ALLIANCE

PRECINCT CONFERENCE, 16th May 2023, Brussels.

Funded by the European Union

# AI4CYBER

# AI4CYBER is a collection of AI services

AI4CYBER is an Ecosystem Framework of Next generation **AI-based services** for critical system **robustness, resilience**, and appropriate **response in the face of advanced and AI-powered cyberattacks**.

**11 Key Results** that cover **6 cybersecurity areas**

A collection of 11 innovative resilience and autonomous response services that leverage AI models and Big Data, aimed to be encapsulated in cybersecurity tools to ensure a continuum of system protection.

# Objective of the Workshop

► Push **open collaboration and open innovation** with respect to **data requirements in cybersecurity research,** by discussing the difficulties and opportunities of using and sharing cybersecurity related datasets for cybersecurity research and cybersecurity solution development.

► In the workshop we will discuss on available open datasets, sources of datasets, reusable datasets, means and principles for sharing datasets, etc.

► Suggested attendees: Cybersecurity researchers dealing with AI used in threat detection, threat simulation, code testing, incident response, etc.

# Challenges of cybersecurity datasets

▶ The **accuracy** of AI algorithms and models is directly dependant on the **quality and amount of data** used to train them.

▶ When using models and algorithms to develop AI-based cybersecurity solutions, the learning requires **big amounts of well-structed and sanitized data** which are sometimes difficult to get.

▶ Some literature works include **open datasets**, but they are difficult to reuse as a basis of further research.

▶ The **confidentiality** (and sometimes secrecy) of cybersecurity data makes it difficult to be shared.

▶ Commercial datasets are **limited**.

▶ Creating the **datasets synthetically** is often hard because it is not easy to replicate the nature of real attacks and getting realistic results is challenging.

# Opportunities of cybersecurity dataset sharing

▶ Improved cybersecurity research by learning from previous works.

  ▶ Learn how others get or build their datasets and their models.

  ▶ Learn data synthetization techniques and realistic data (quality data) creation.

  ▶ Learn anonymisation and aggregation techniques.

▶ Test your models on top of other datasets to benchmark and improve your model.

▶ Verify the quality and correctness of other researchers' solutions.

▶ Well-structured and high quality pools of FAIR data.

# FAIR Data

▶ Findability: open (meta)data shall be easily searchable and locatable -> include identifiers, keywords, version numbers and metadata.

▶ Accessibility: openness of the (meta)data -> clarify and enable means of access, identify the repository where the data will be deposited, access protocols to and any potential restriction.

▶ Interoperability: refers to the ability of the (meta)data to integrate or work together with minimal effort with other (meta)data and tools -> clarify vocabularies, standards, formats and methodologies used.

▶ Reusability: the ability of the (meta)data to be utilized in other contexts -> include attributes, clear and accessible licenses, follow standards and identify provenance.

# From ChatGPT - metrics commonly used to evaluate the quality of datasets in AI

- **Accuracy**: Measures the proportion of correctly labeled instances in a dataset. It indicates how well the dataset represents the ground truth.

- **Completeness**: Measures the extent to which the dataset covers the entire range of possible values and scenarios. It evaluates whether the dataset contains sufficient samples to adequately represent the target population.

- **Consistency**: Assesses the absence of contradictions or conflicts within the dataset. It measures the coherence and uniformity of the data.

- **Reliability:** Measures the degree of trustworthiness and dependability of the dataset. It evaluates factors such as data collection methods, sources, and potential biases.

- **Relevance**: Determines the degree to which the dataset aligns with the objectives and requirements of the AI task. It evaluates the suitability and applicability of the dataset for the specific problem at hand.

- **Precision**: Measures the proportion of true positives out of the total positive instances. It is particularly relevant for datasets used in classification tasks.

- **Recall**: Measures the proportion of true positives identified out of the total actual positive instances. It is also important in classification tasks, especially when it is crucial to minimize false negatives.

- **Bias**: Evaluates the presence of systematic errors or prejudices in the dataset that may lead to unfair or discriminatory outcomes. It is essential to assess and mitigate biases in AI systems.

- **Noise**: Assesses the presence of irrelevant or misleading data within the dataset. It measures the level of interference or distortion that may impact the performance of AI models.

- **Diversity**: Measures the variety and heterogeneity of the dataset in terms of different attributes, categories, or scenarios. A diverse dataset helps ensure the model's generalization capability.

- **Balance**: Evaluates the distribution of instances across different classes or categories. It assesses whether the dataset has an equal representation of different classes or if it is skewed towards certain categories.

- **Annotation Quality**: Applies specifically to labeled datasets. It assesses the accuracy and consistency of the annotations or labels assigned to the instances in the dataset.
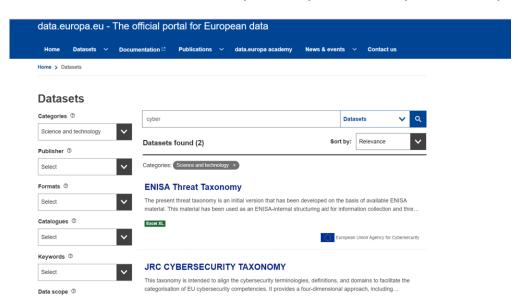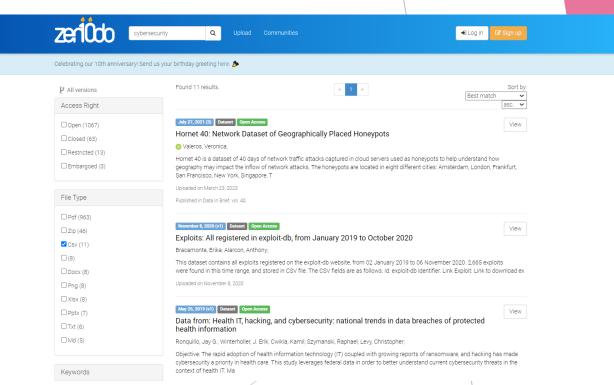
# European data pools

▶ European Data Portal: https://data.europa.eu/en

    ▶ "Science & technology", "cyber" ->

    ENISA Threat taxonomy, JRC Cybersecutity Taxonomy

▶ Zenodo: https://zenodo.org/

    ▶ "cybersecurity" -> 11 csv, 8 xlsx, some are studies and listings.





**Other data pools** (e.g. IEEE DataPort) require subscription for accessing most datasets.

# Open discussion

► What are the **major barriers** to sharing?

► What possible **solutions** do you see for increasing dataset sharing?

► How to **avoid sharing sensitive information** in the datasets? Anonymisation, pseudonimisation, etc. do always work?

► What **trustworthy open data pools** do you know so as the datasets can be shared? Zenodo (Open Aire), JRC, etc.

► Could an **EU cybersecurity researchers' community** be created with the goal of sharing datasets and information about the usage of the datasets, in which researchers contribute to with datasets and related info?

# AI4CYBER

**TRUSTWORTHY ARTIFICIAL INTELLIGENCE FOR CYBERSECURITY REINFORCEMENT AND SYSTEM RESILIENCE**

https://ai4cyber.eu

https://twitter.com/Ai4Cyber

https://www.linkedin.com/company/ai4cyber/

Erkuden Rios

Project Coordinator

erkuden.ríos@tecnalia.com

# Thank you for your attention!

Funded by the European Union