# Model Lake: a New Alternative for Machine Learning Models Management and Governance

Moncef Garouani [1], Franck Ravat [1], Nathalie Valles-Parlangeau [2]

[1] IRIT, UMR5505 CNRS, Université Toulouse Capitole, Toulouse, France
[2] LIUPPA, Université de Pau et des Pays de l'Adour, Anglet, France

moncef.garouani@irit.fr

# PLAN ⟩

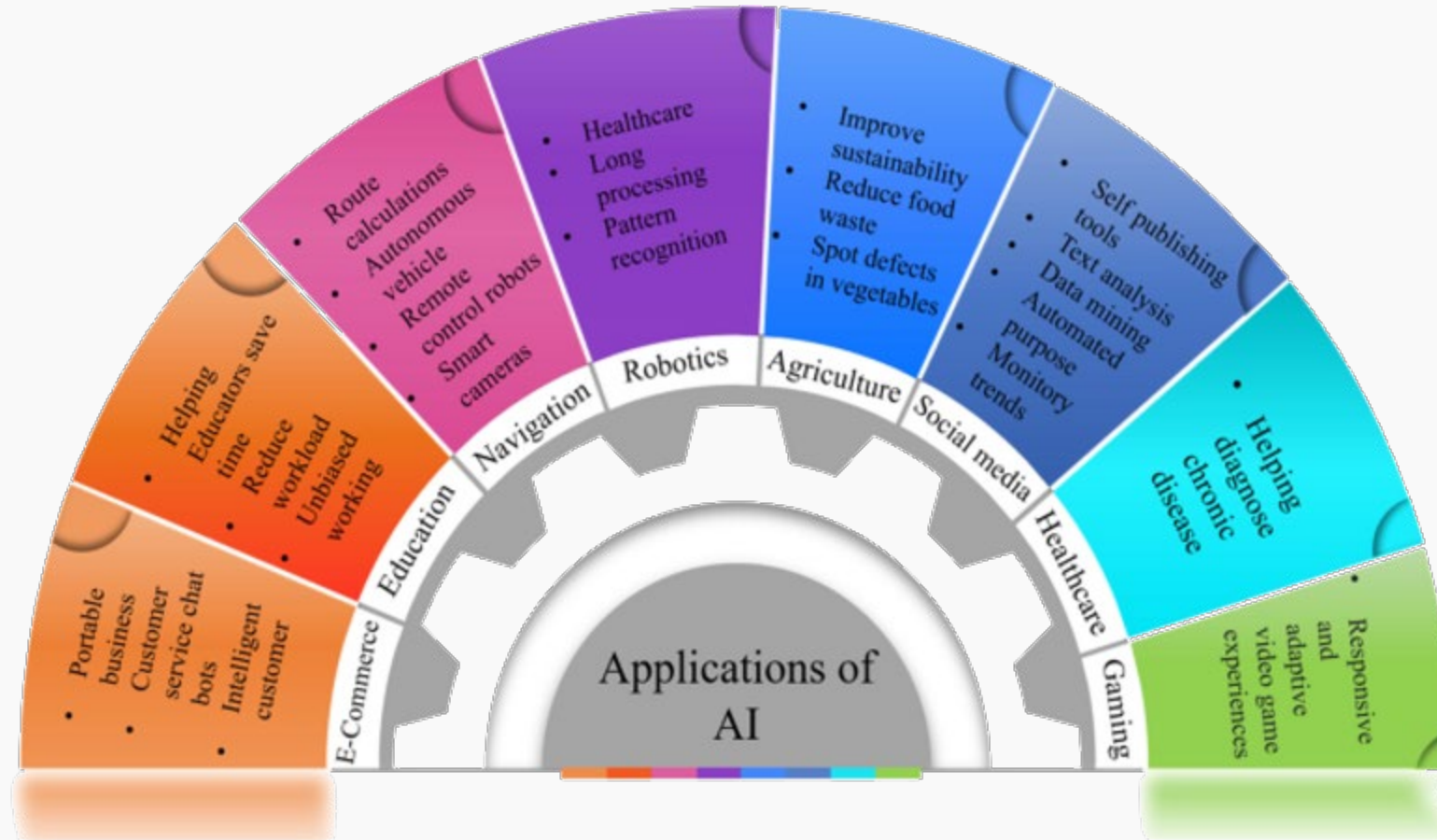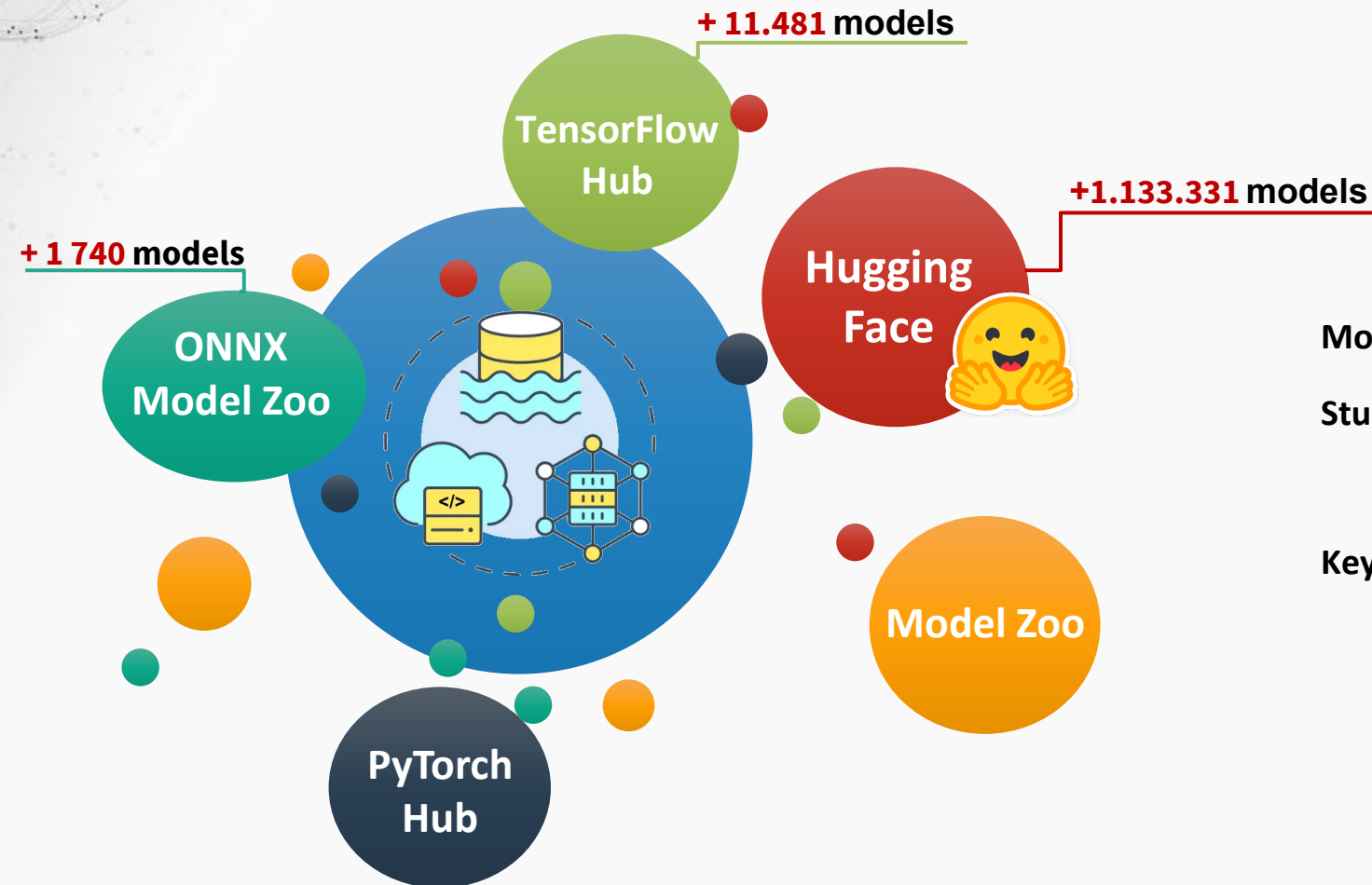# Context: Success of Machine Learning & Deep Learning (1/2)

# Context: Management problem of ML/DL Models (2/2)



- *How can we manage and understand many available ML models?*

- *How can we know what a model does and how it was trained?*

- *How can we ensure ethical use and trace model lineage?*

- *How can we ensure compliance with regulations?*

- *How can we improve models to avoid mistakes?*

# Model Repository, Registry, Zoo, …



**+ 11.481** models

**+1.133.331** models

**+ 1 740** models

TensorFlow Hub

Hugging Face

ONNX Model Zoo

Model Zoo

PyTorch Hub

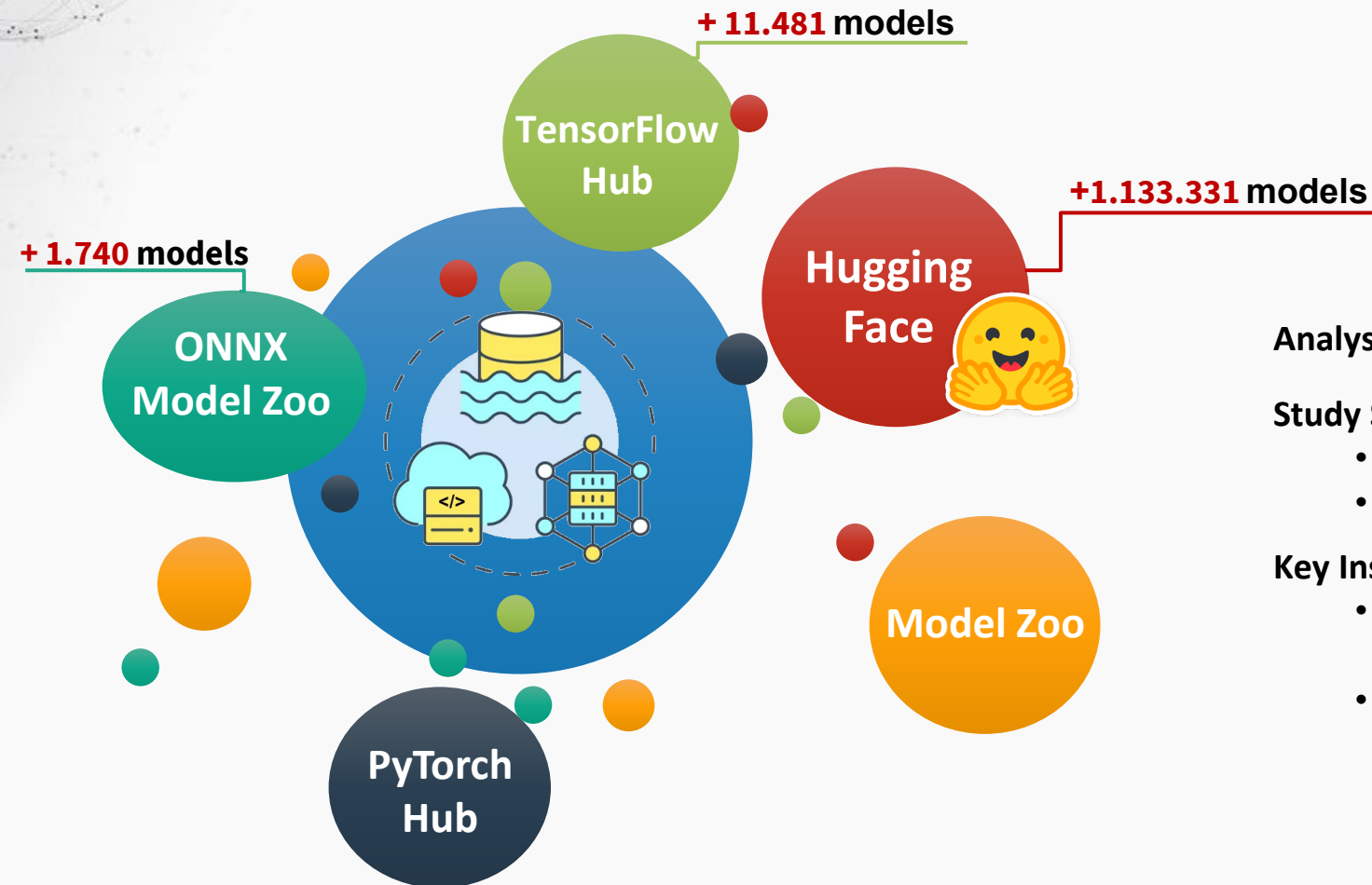**Model Reuse in the Hugging Face Registry (W. Jiang et al.):**

**Study Approach:**
- Interviews with 12 practitioners to identify challenges.
- Analysis of 63,182 models in the registry.

**Key Challenges:**
- Missing attributes: model lineage, training data.
- Disparities in claimed vs. actual performance.
- Privacy and ethical concerns due to opaque data lineage.

*W. Jiang et al. An Empirical Study of Pre-Trained Model Reuse in the Hugging Face Deep Learning Model Registry. 2023. arXiv: 2303.02552.*

# Model Repository, Registry, Zoo, ...



**+ 11.481** models

**TensorFlow Hub**

**+1.133.331** models

**Hugging Face**

**+ 1.740** models

**ONNX Model Zoo**

**Model Zoo**

**PyTorch Hub**

**Analysis of Hugging Face Model Cards (W. Liang et al.):**

**Study Scope:**
- Analyzed 74.970 model repositories from 20.455 user accounts.
- Found only 32.111 repositories (44.2%) include model cards.

**Key Insights:**
- Over 56% of models lack proper documentation, reducing reliability.
- **Highlights the need for data-centric approaches to improve model quality and support responsible AI development.**
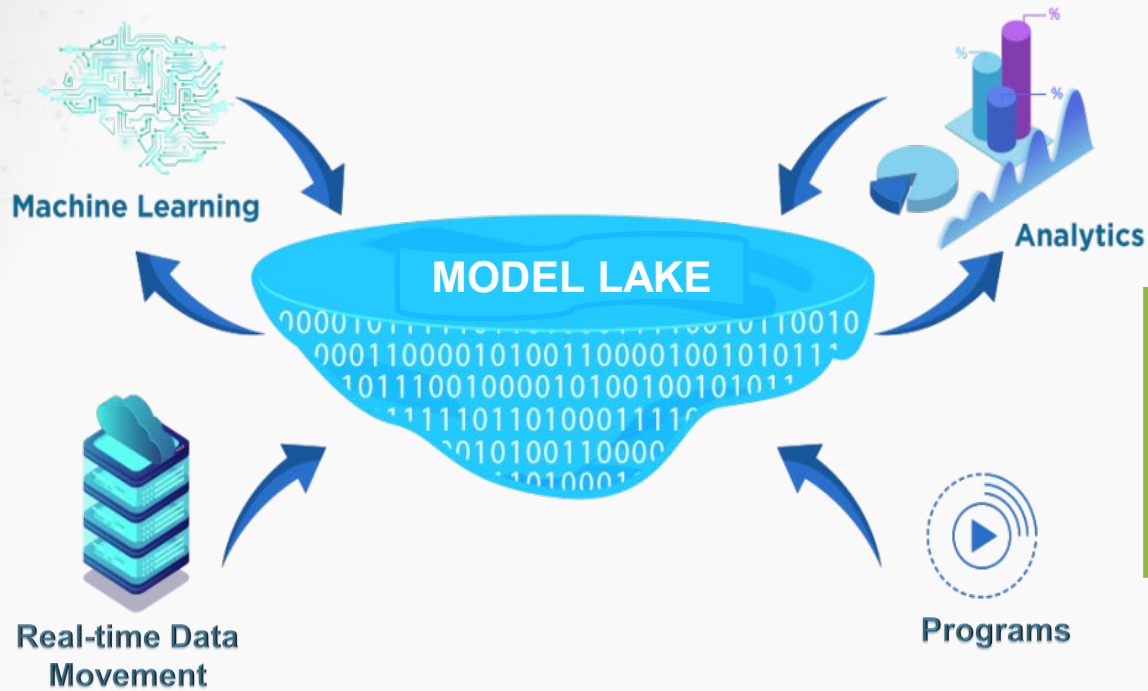
W. Liang et al. What's documented in AI? Systematic Analysis of 32K AIModel Cards. 2024. arXiv: 2402.05160.

# Model Lake

❑ **Data Lake efficiency**

# Model Lake

❑ Model Lake Definition



**Machine Learning**

**MODEL LAKE**

**Analytics**

**Real-time Data Movement**

**Programs**

**Definition**. Model Lake stands as an integrated ecosystem encompassing respectively the input, process, output and governance aspects of both mined data and developed models. It acts as a centralized hub and management system accommodating diverse data and model types, meeting the requirements of various stakeholders including data engineers, data scientists, data analysts, and business intelligence professionals.

**Key Functionalities:**

- Raw data ingestion, processing, storage, and governance.
- Model training, fine-tuning, review, and monitoring.
- Data, model, and code provenance, management, and governance.

# Model Lake
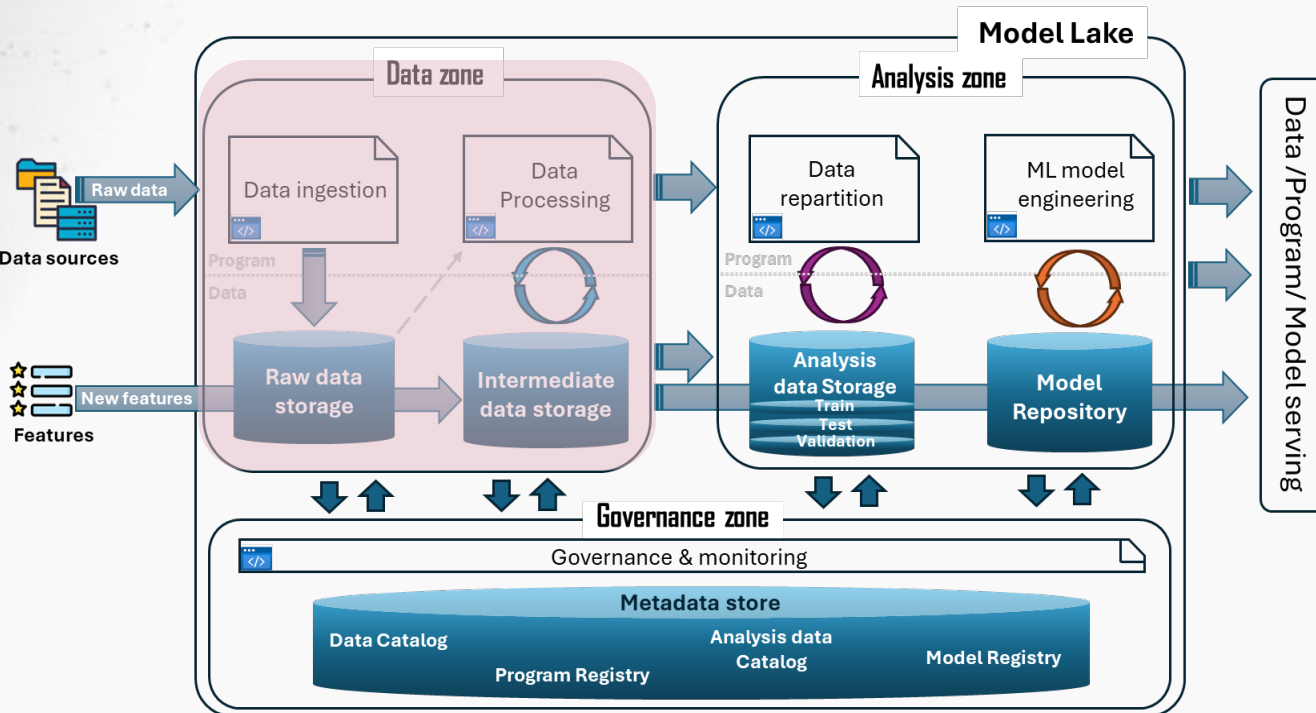
## ❏ Model Lake Architecture



*Fig. 1. The proposed Model Lake architecture.*

**Data Zone :**

❏ **Purpose**: Manages data ingestion, processing, and intermediate storage.

  •**Ingestion**: Connects to data sources for extraction and change tracking.
  •**Processing**: Standardizes raw data with operations like integration, cleaning, transformation, and reduction.
  •**Intermediate Storage**: Saves processed data and metadata for lineage tracking.

❏ **Key Role**: Prepares data for analysis.

# Model Lake

❑ Model Lake Architecture



*Fig. 1. The proposed Model Lake architecture.*

**Analysis Zone :**

❑ **Purpose**: Central hub for data exploration and ML model development.

❑ **Key Features:**
- Advanced data exploration (meta-features, attributes, transformations).
- ML model development and evaluation.
- Model storage for production.

❑ **Additional Capabilities:**
- Continuous monitoring and feedback for performance and reliability.
- Lineage tracking, model comparison, auditing, and compliance.

# Model Lake

❑ Model Lake Architecture



Fig. 1. The proposed Model Lake architecture.

**Governance and Management Zone:**

•**Purpose**: ensure Data, Program, Model security, lifecycle management, access, and metadata management.

•**Preventing Model Swamp**:
  •Maintains accessibility and usability of the Model Lake.
•**Metadata Store**:
  •Records metadata for all ML workflow tasks and iterations.
  •Tracks job details (e.g., training date, artifact sources).
•**Model Lineage**:
  •Combines data, model, and code lineage.
  •Tracks metadata like feature data sources, parameters, and performance metrics.
  •Ensures full traceability for each registered model.

# Model Lake Metadata Management

❑ **Metadata Model on Data zone**

We adopt the 5W1H (What, Who, Where, When, why, how) method to facilitate a systematic understanding of data ingestion and processing. This method prompts the following inquiries:

– **What**: Identifying external data sources and the nature of ingestion activities (ingested datasets, their quality, security level, and interrelations).
– **Who**: Determining ownership of the source data, as well as the individuals responsible for data ingestion and processing.
– **Where**: Locating the storage sites for ingested and processed datasets and associated data ingestion/processing code.
– **When**: Establishing timelines for the ingestion and processing of datasets.
– **Why**: Understanding the purpose behind the data processing activities.
– **How**: Understanding the ingestion and processing operations.

# Model Lake Metadata Management

❑ Metadata Model on Data zone

# Model Lake Metadata Management

❑ Metadata Model on Data Analysis

# Model Lake Metadata Management

# Conclusion & Perspectives

## Conclusion

- Rapid ML model growth presents both opportunities and challenges,
- Lack of standardized management risks limiting their full potential.
- **Model Lakes**: A promising solution for centralized, scalable model management.
- Success requires collaboration across the ecosystem and commitment to responsible AI.

## Perspectives

- Expand model lake system to include additional analysis types and ML pipeline artifacts.
- Develop a **recommender system** to enhance data and model search and discovery.

THANK YOU!

# Model Lake: a New Alternative for Machine Learning Models Management and Governance

Moncef Garouani [1], Franck Ravat [1], Nathalie Valles-Parlangeau [2]

[1] IRIT, UMR5505 CNRS, Université Toulouse Capitole, Toulouse, France
[2] LIUPPA, Université de Pau et des Pays de l'Adour, Anglet, France

moncef.garouani@irit.fr