



# Testing the Untestable

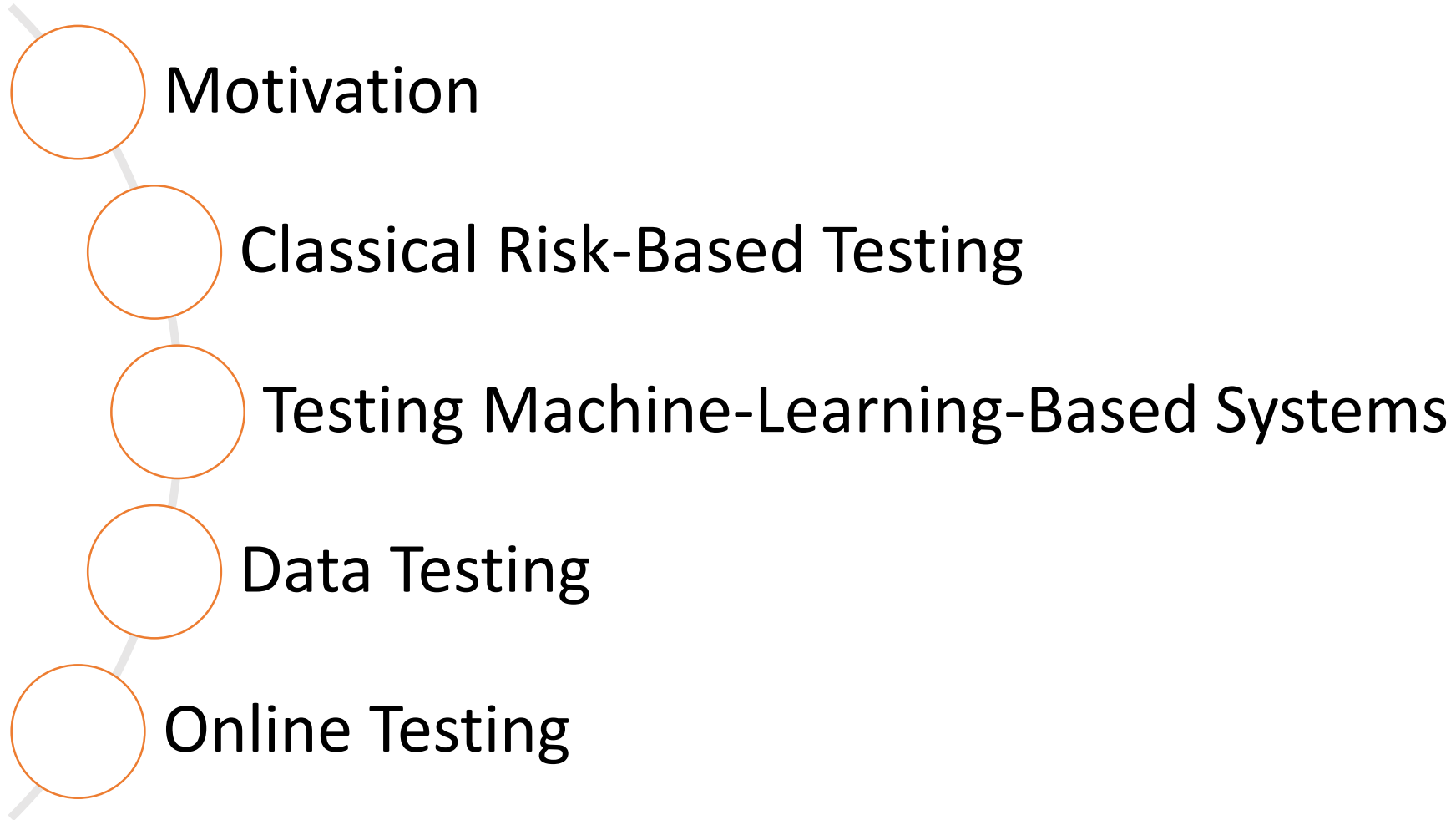
## Risikobasierte Qualitätssicherung für Machine-Learning Systeme

Prof. Dr. Michael Felderer

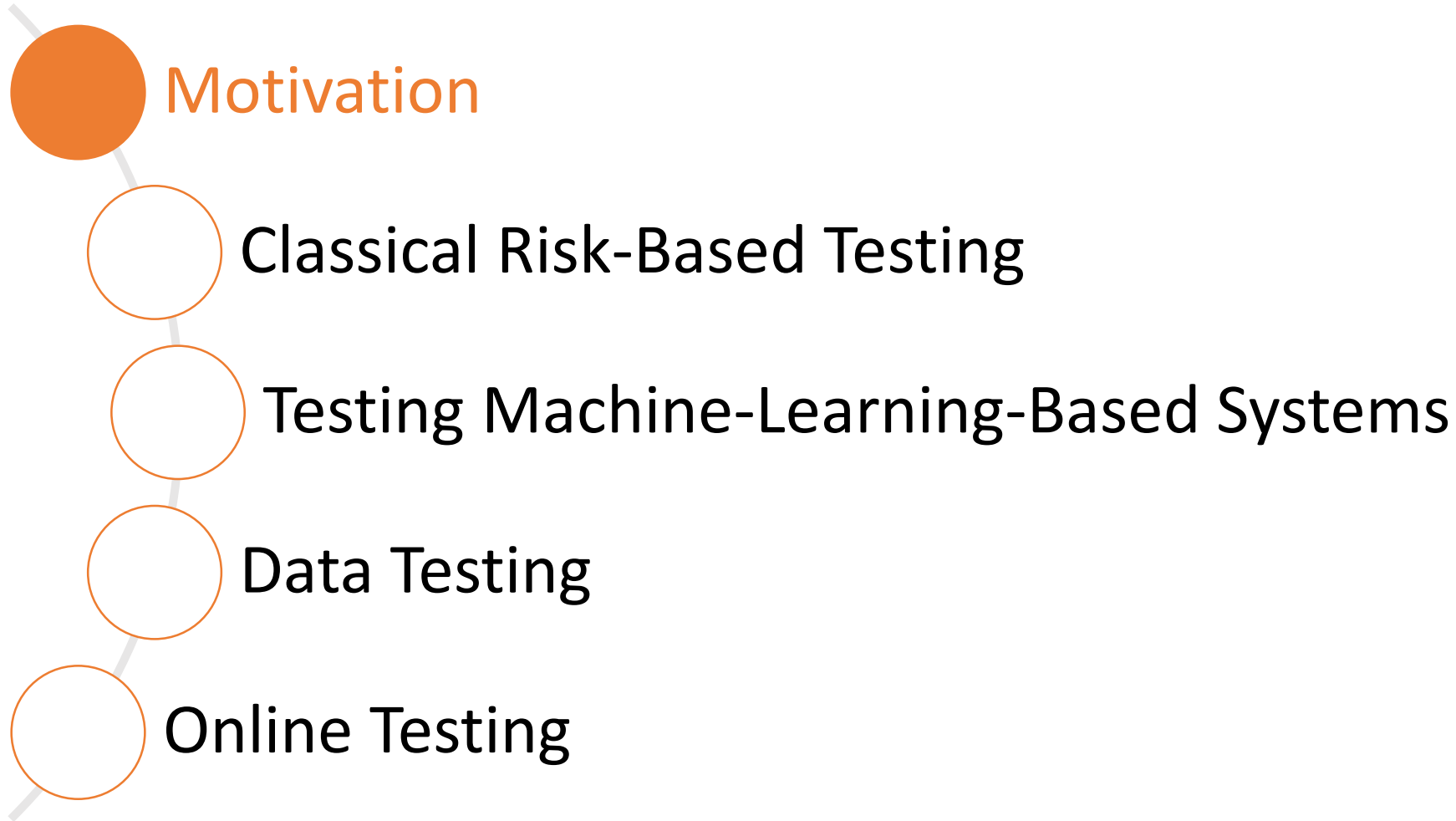
Institut für Informatik

Universität Innsbruck

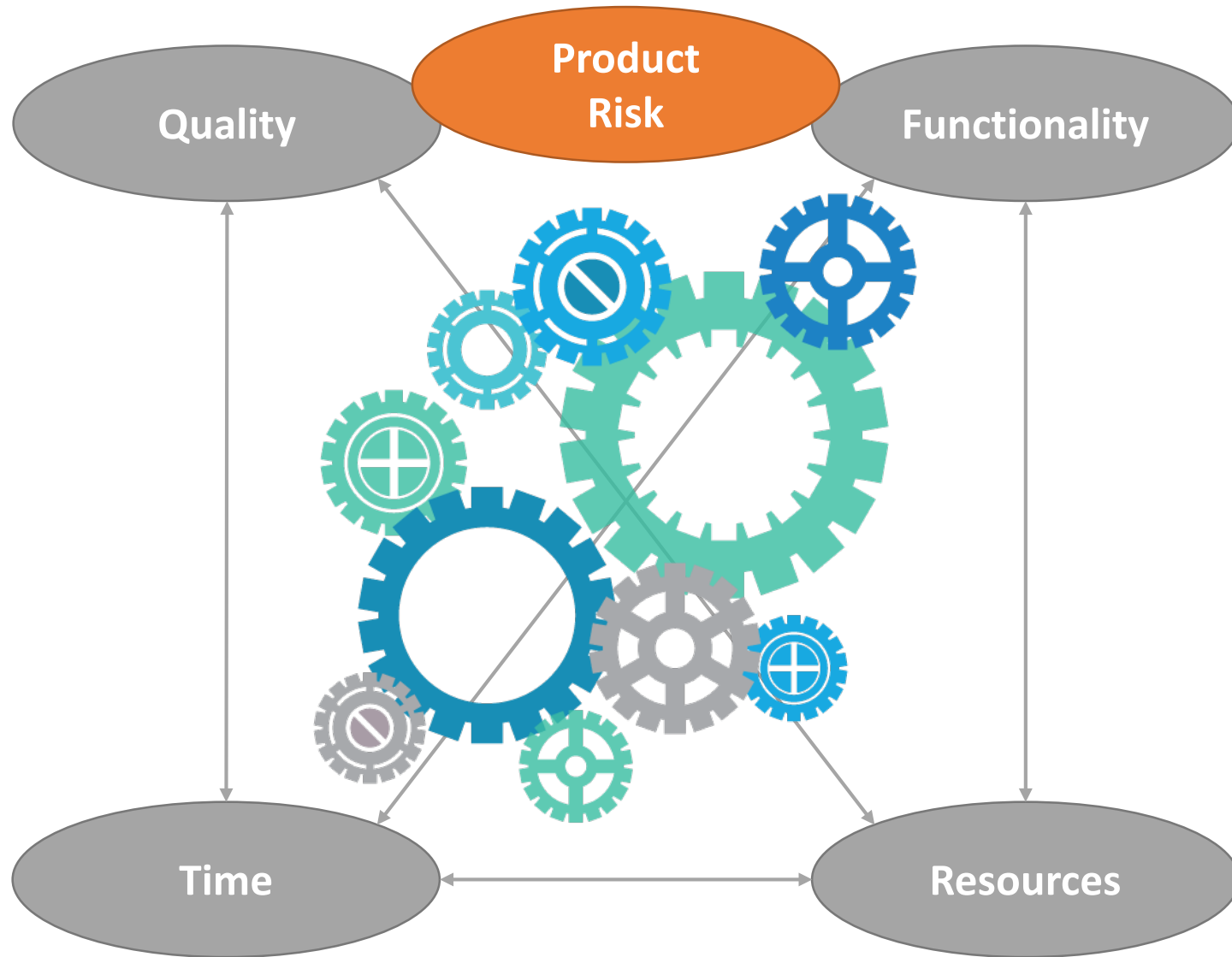
# Agenda



# Agenda




# Quality and Risk in Software Development



# Quality and Risk of ML-based Systems



 **REUTERS** Q ☰

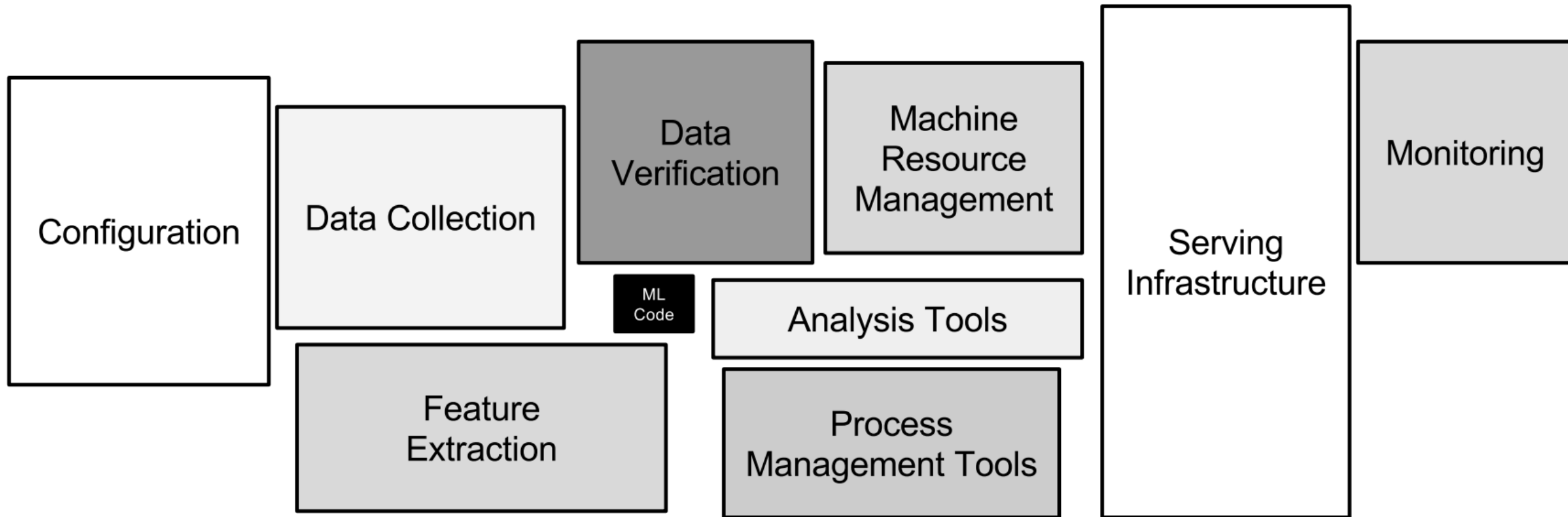
TECHNOLOGY NEWS  
OCTOBER 10, 2018 / 5:12 AM / 2 YEARS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin 🐦 f

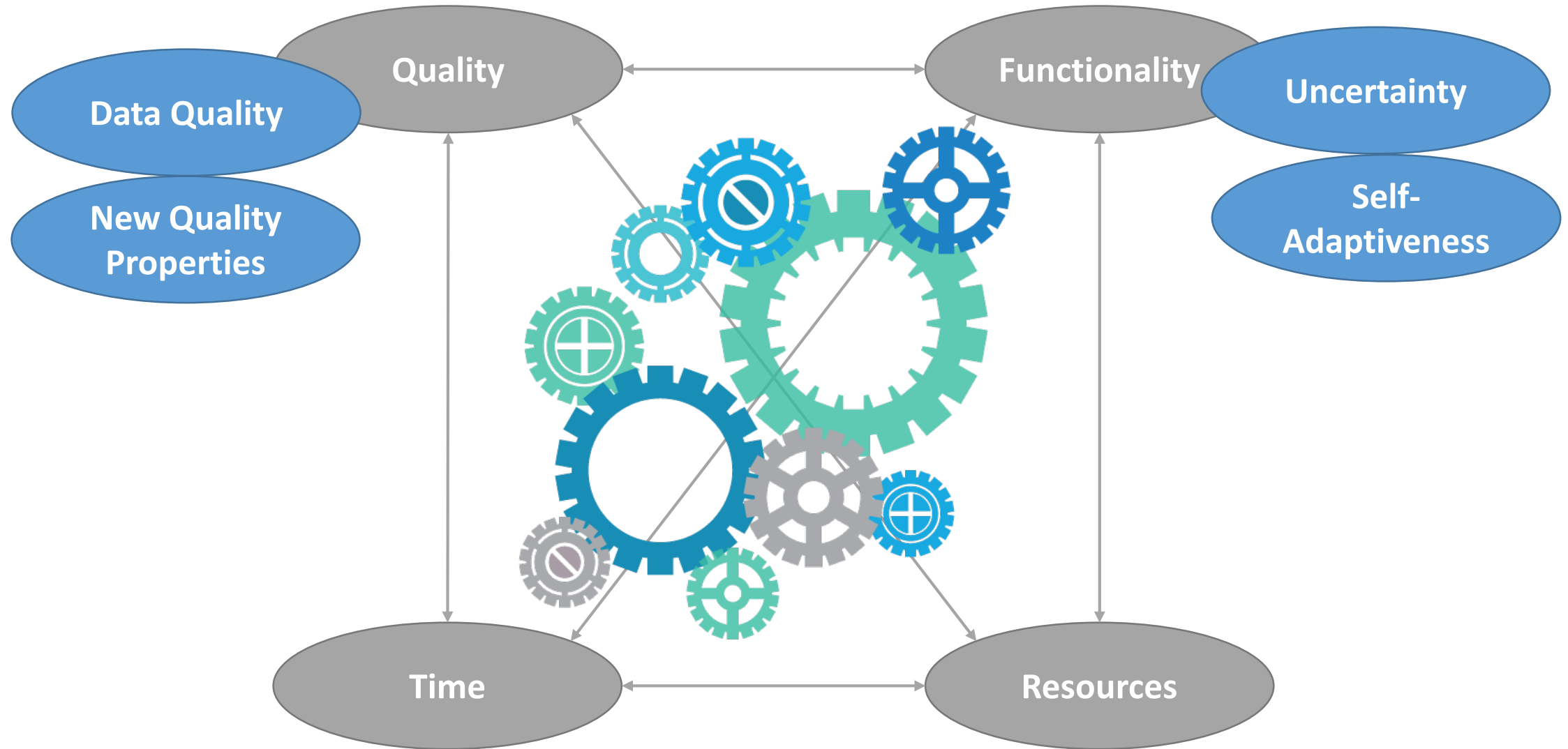
SAN FRANCISCO (Reuters) - Amazon.com Inc's ([AMZN.O](#)) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.

# ML Systems are IT Systems not only Algorithms



Sculley et al.: Hidden Technical Debt in Machine Learning Systems, NIPS 2012

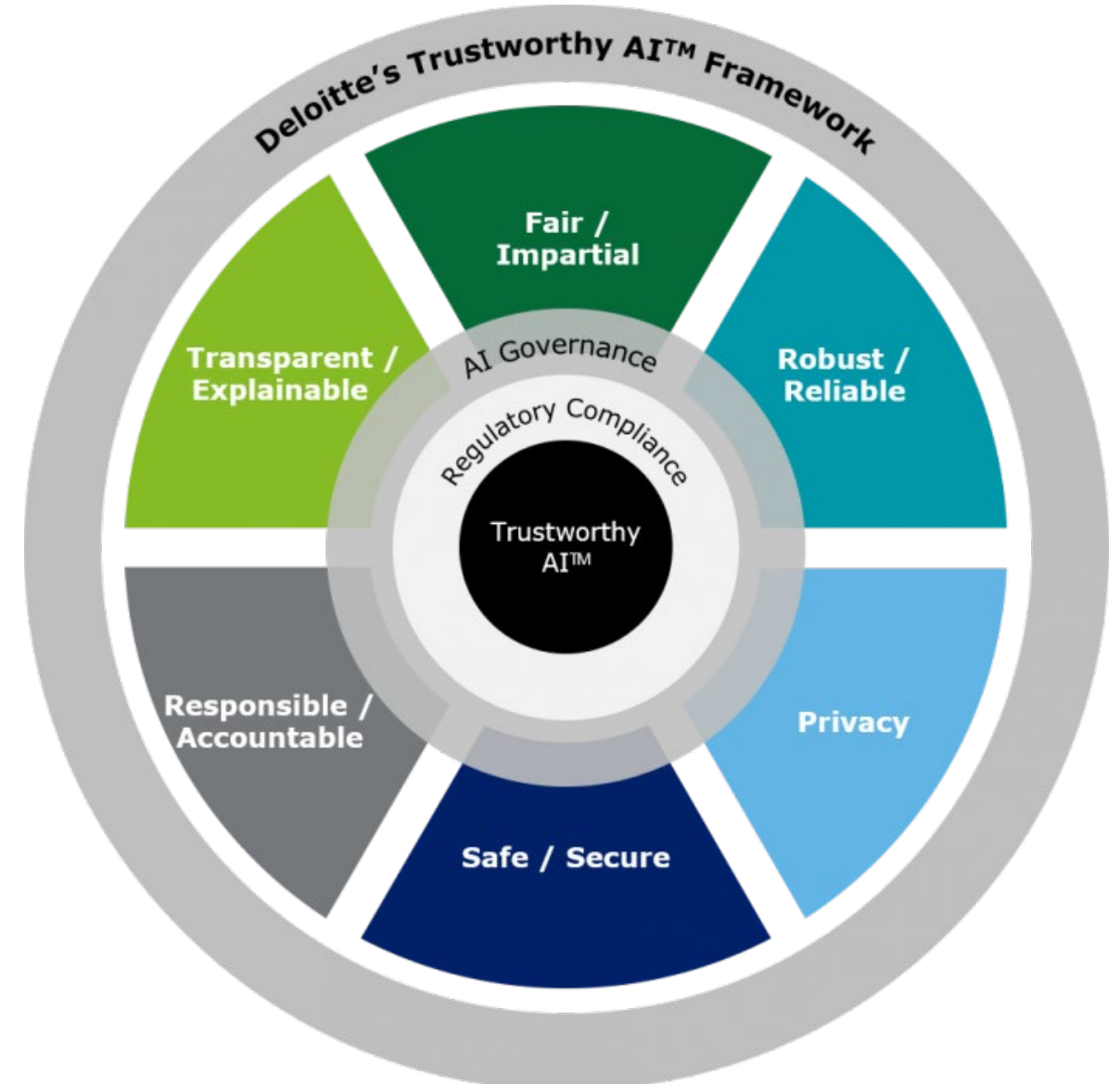
# What Makes (AI)ML-Based Systems Different?



# New Quality Properties: Trustworthiness of AI



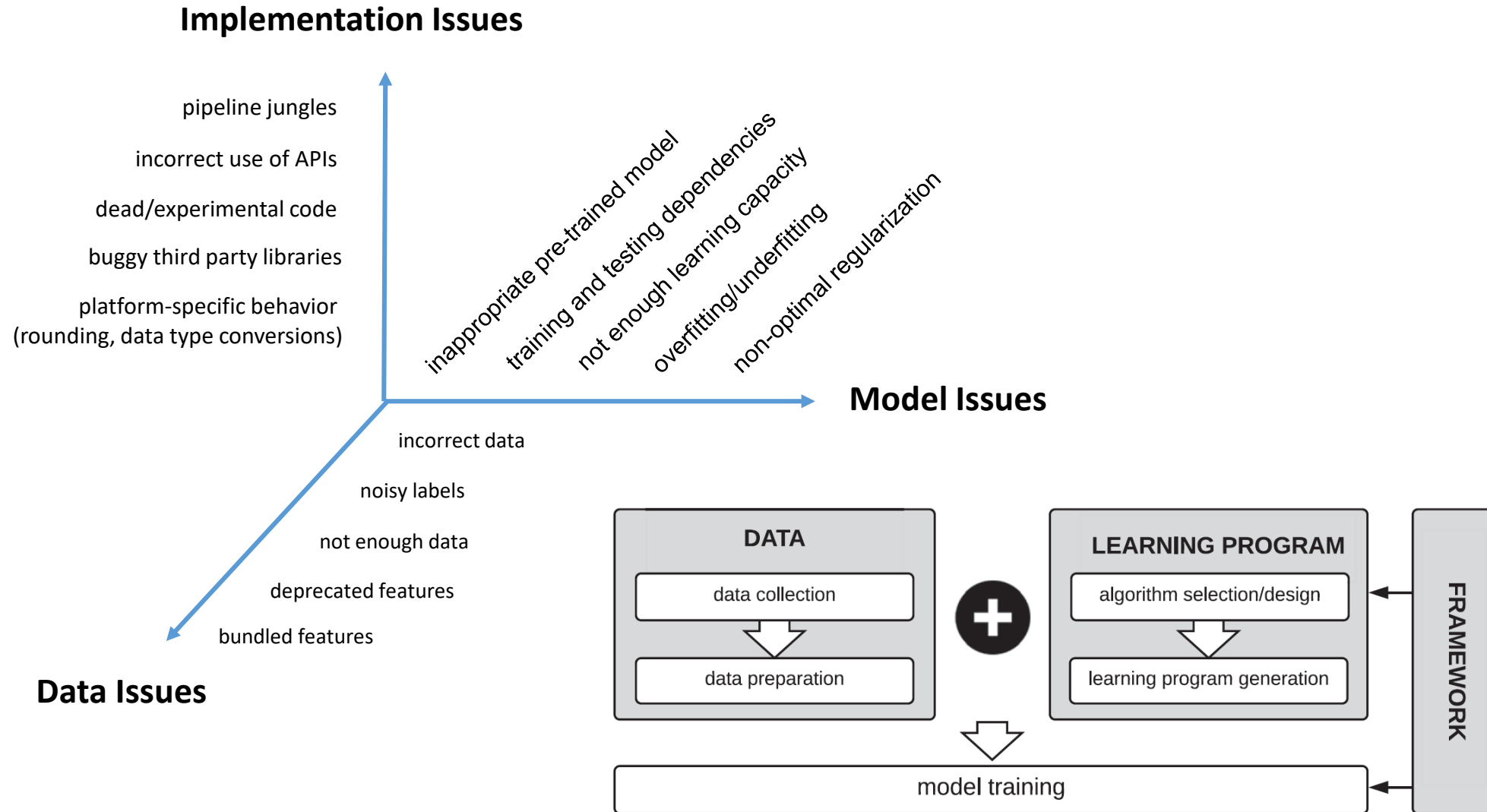
EU Ethics Guidelines for Trustworthy AI



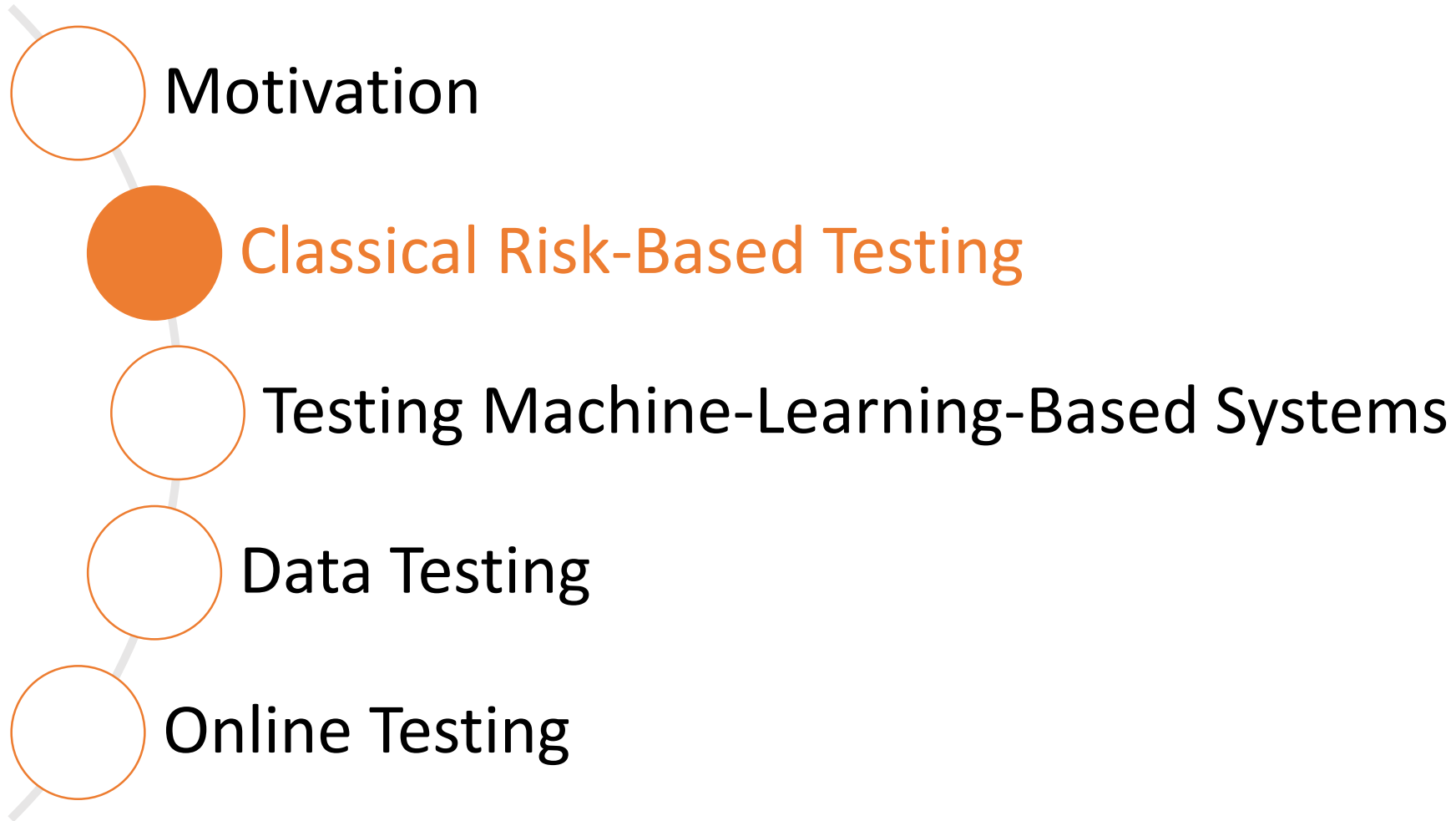
Deloitte's Trustworthy AI Framework



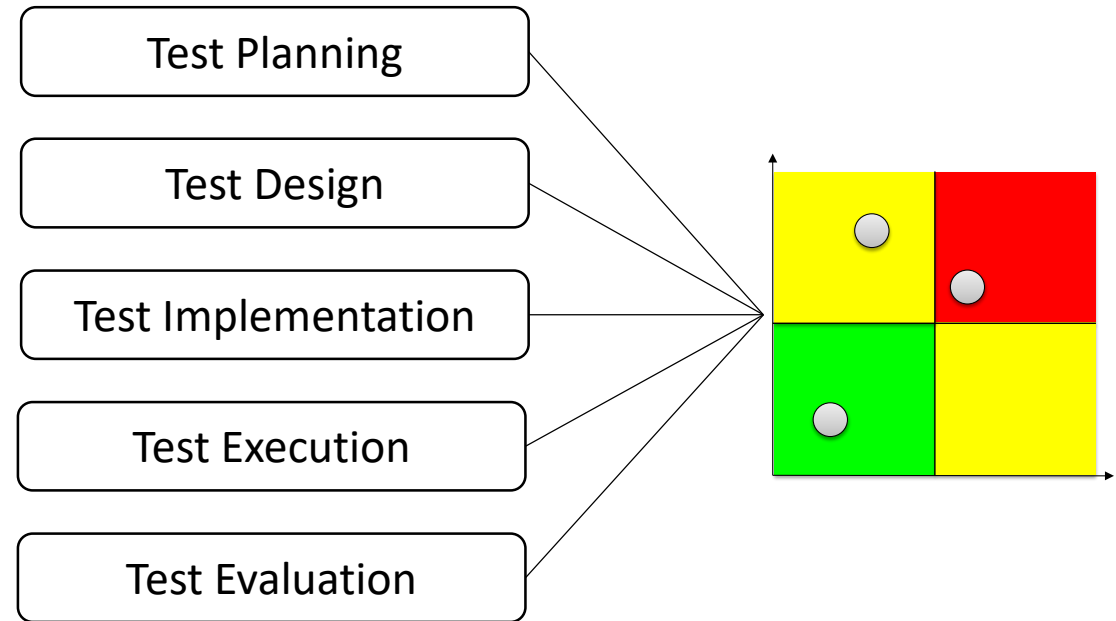
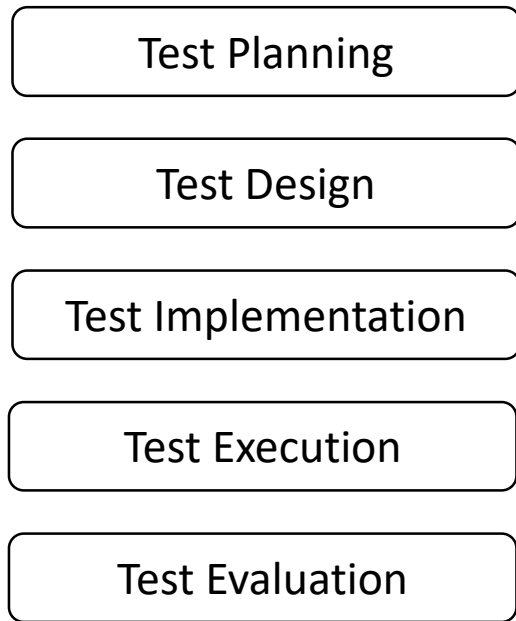
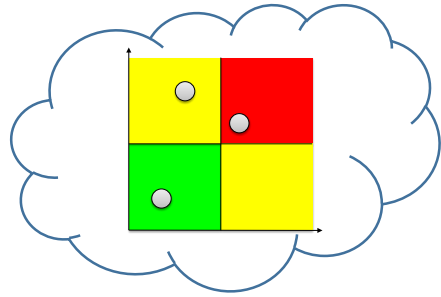
# Machine-Learning Systems as Software Systems



# Agenda



# Risk-Based Testing (Risk-Based Quality Assurance)

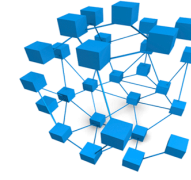
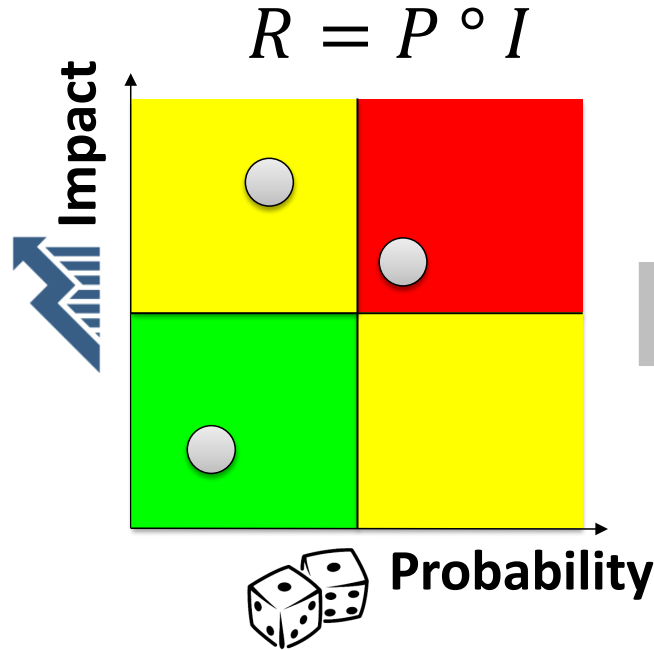


# Risk Concept in Software Quality Engineering

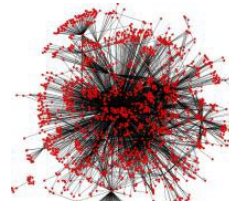
Business-Oriented Criteria



What is the impact if a defect occurs?



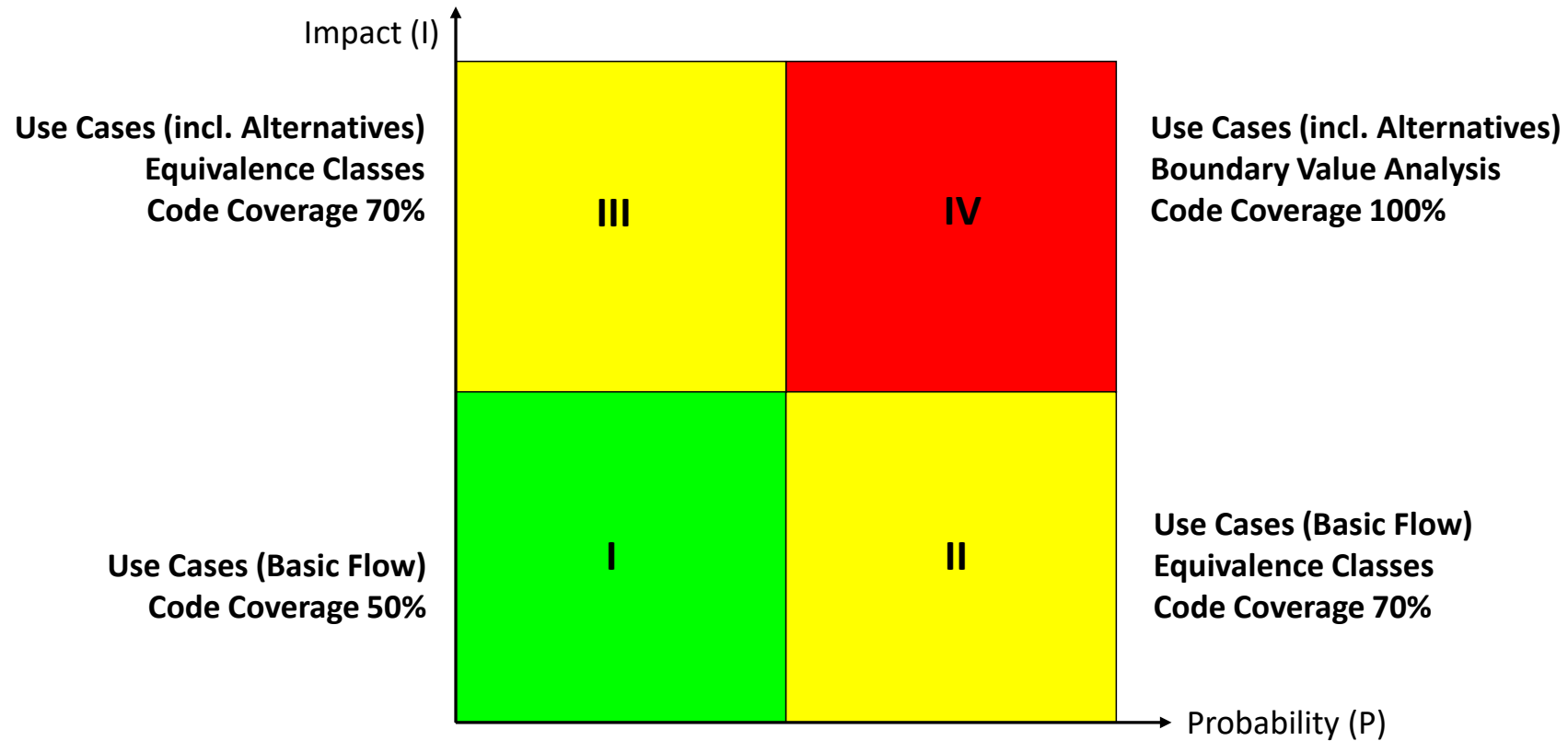
Asset (Risk Item)



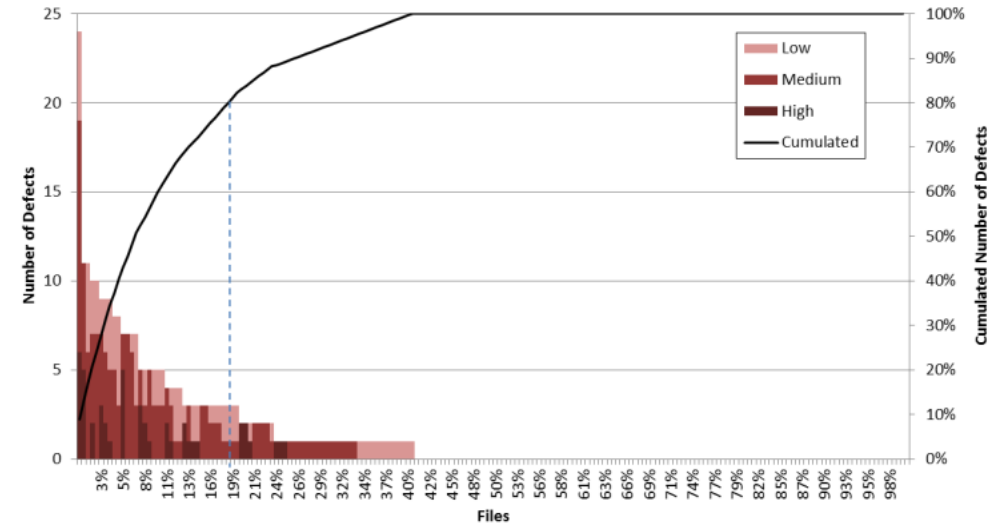
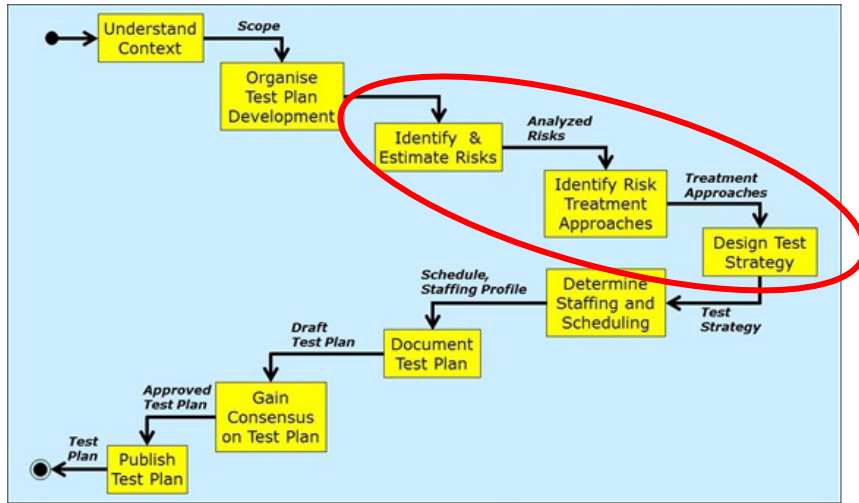
What is the probability a defect will occur?

Technology-Oriented Criteria

# Risk-Based Test Strategy



# Effectiveness and Efficiency of Risk-Based Testing



**MORE TEST LESS**

## Compliance

- Aligned with standard
- Minimizing critical risks
- ...

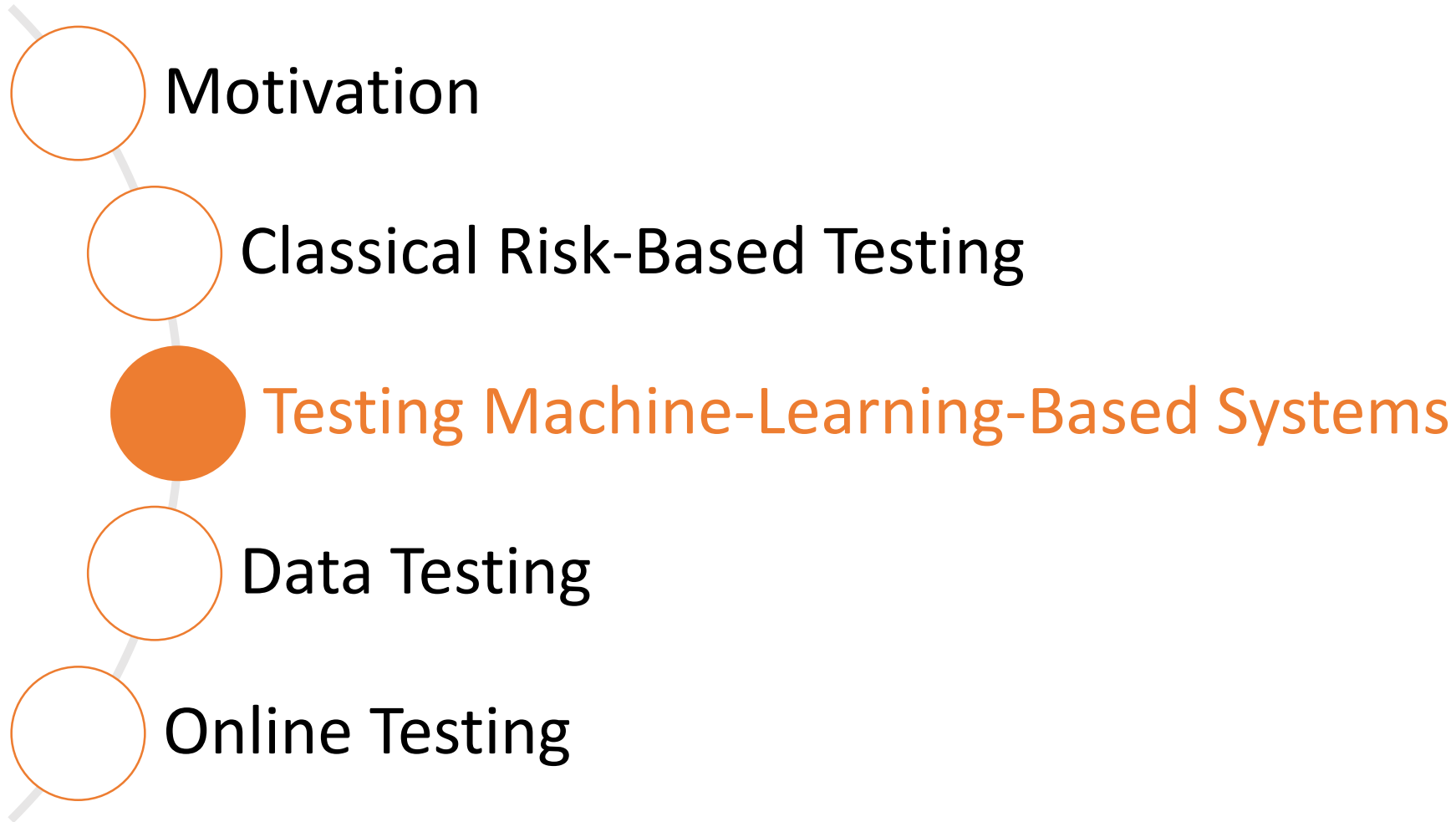
## Effectiveness

- Finding more defects
- Finding defects earlier
- Finding the critical defects
- ...

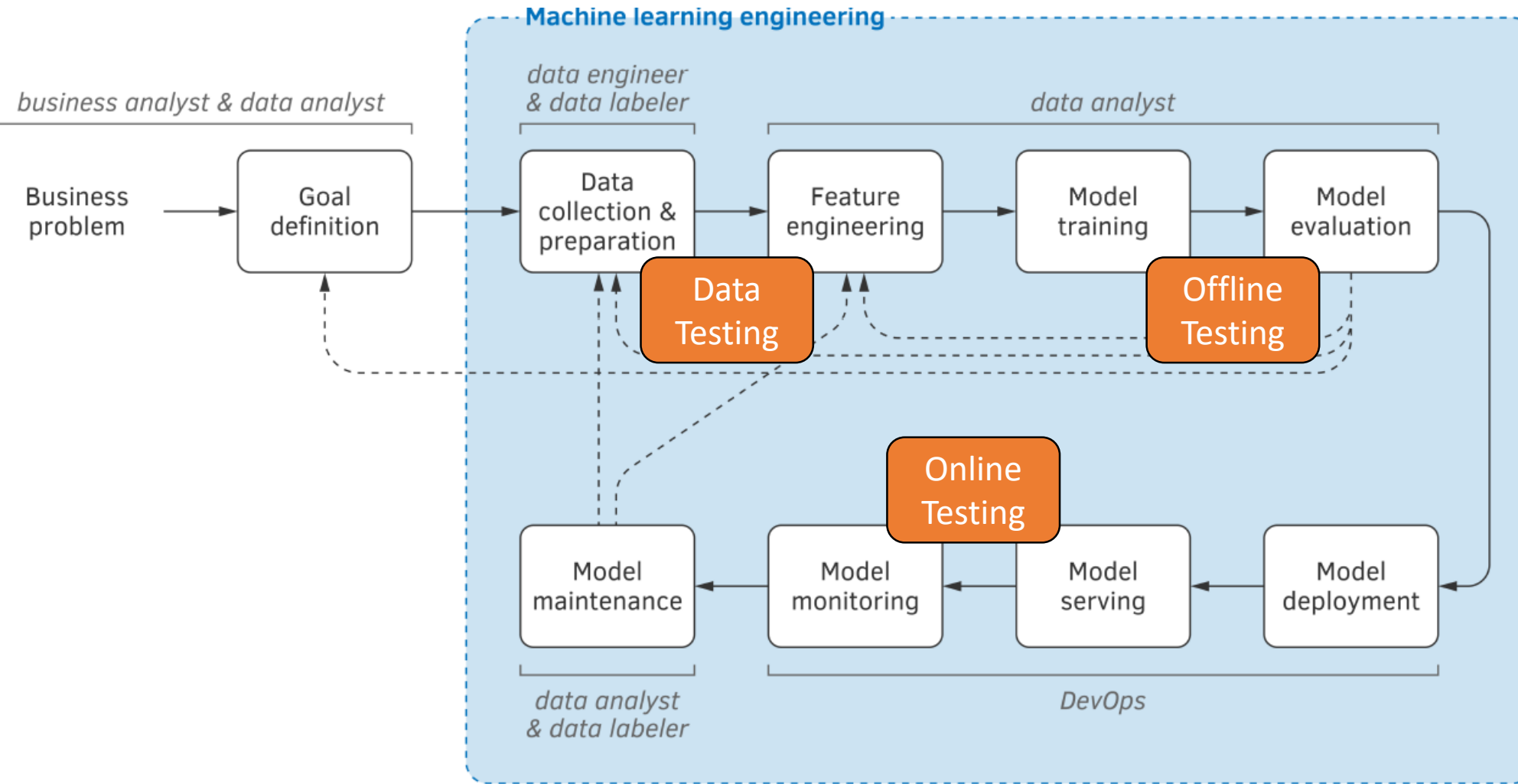
## Efficiency

- Reduce time of testing
- Reduce cost of testing
- ...

# Agenda

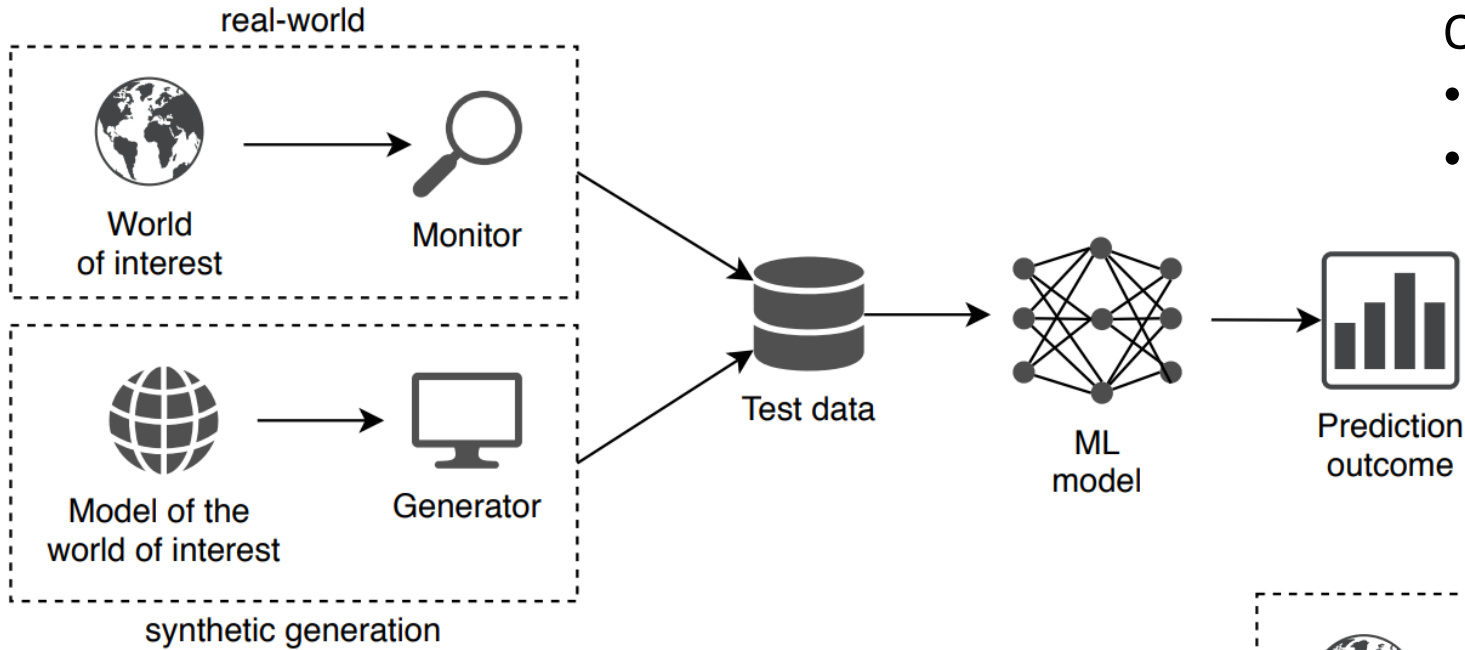


# ML Engineering and Testing





# Offline and Online Testing

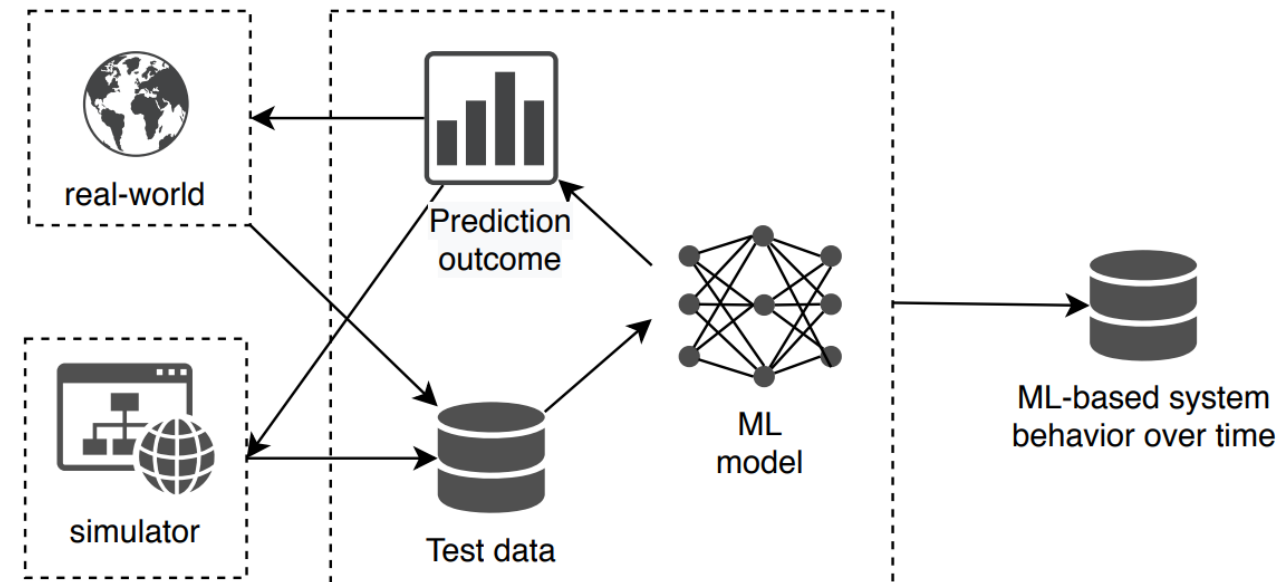


Offline Testing:

- Testing ML model as **standalone component**
- ML model tested as a unit in **open loop mode**

Online Testing:

- Testing ML model in **real or simulated environment**
- ML model tested as a unit in **closed loop mode**



# White-Box and Black-Box Testing?

## White-box

Model's  
Architecture  
& Weights



Hyper  
parameters



Neuron  
Activations



## Data-box

Training  
set



Test  
set

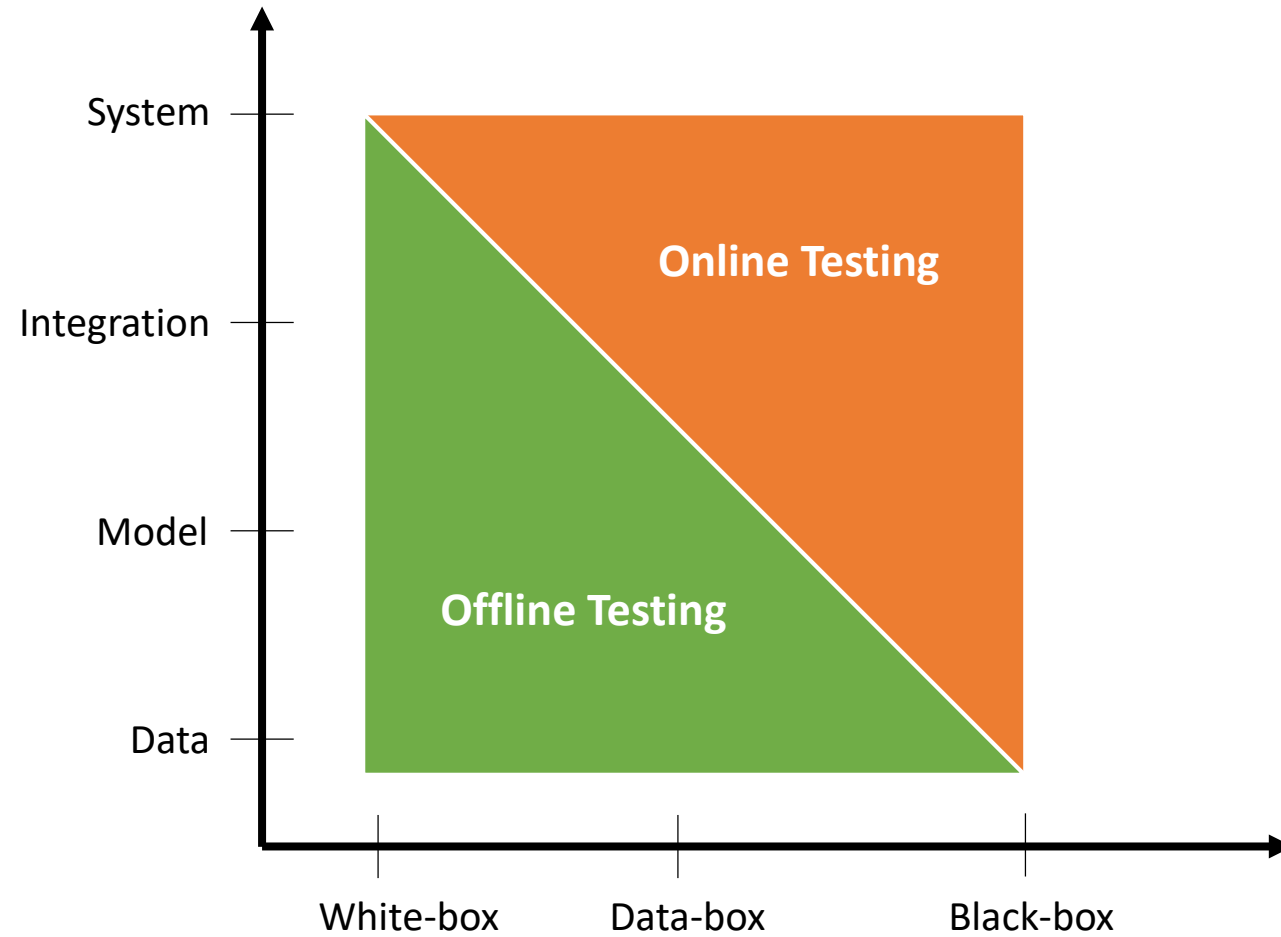


## Black-box

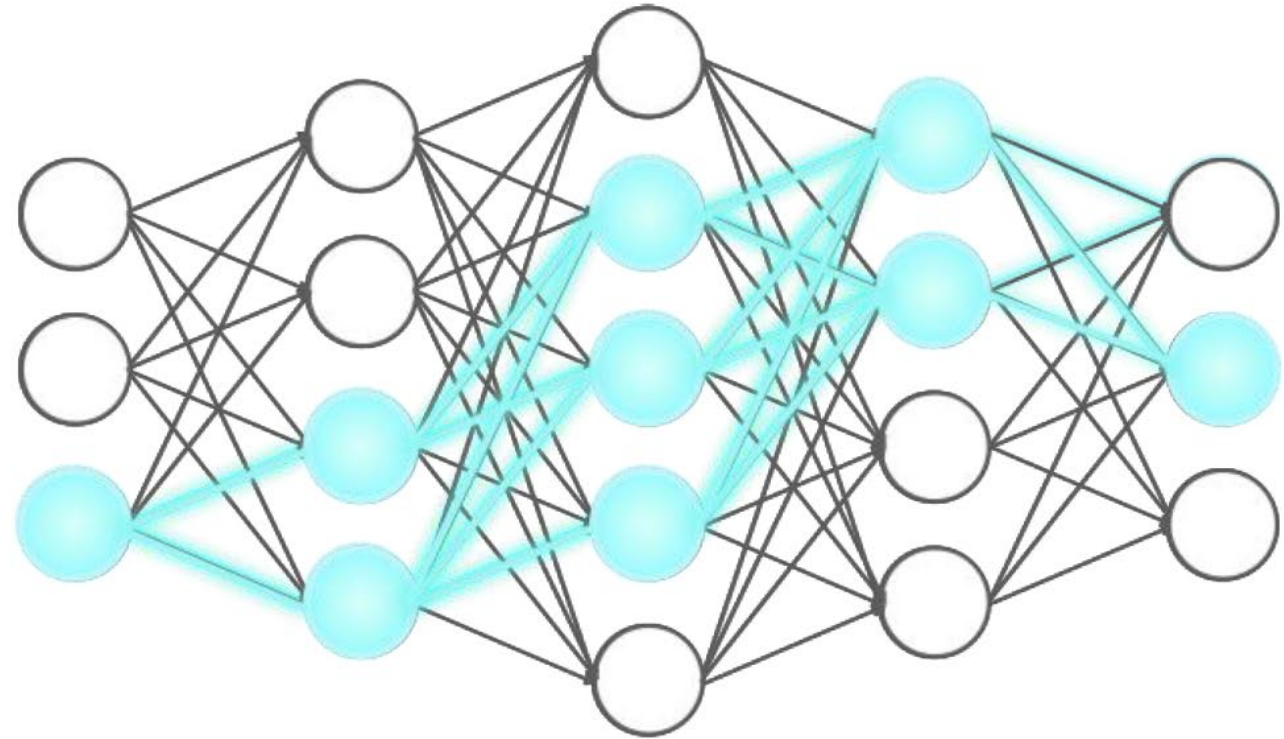
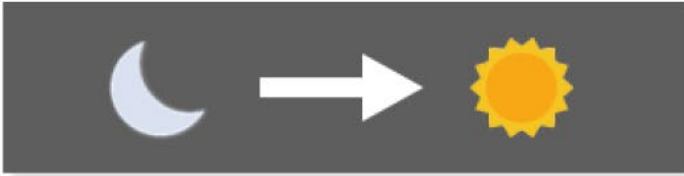
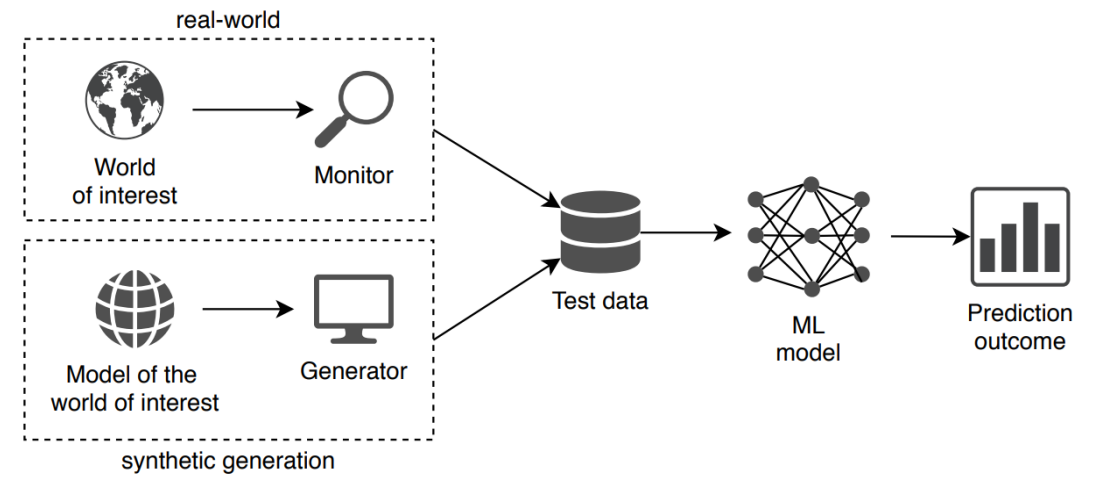
Trained Model



# Offline vs. Online Testing



# Offline Testing (White Box)



**0.66 Bicyclist**

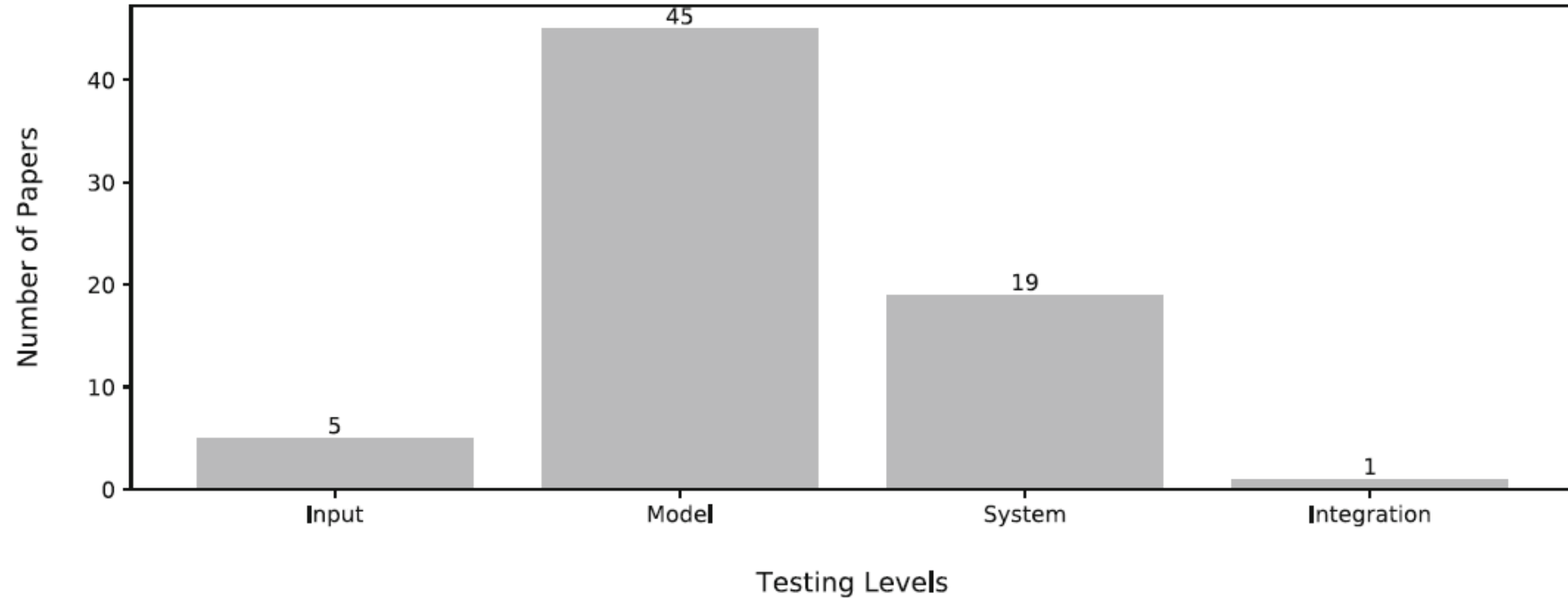
**0.93 Pedestrian**

Neuron Coverage is not strongly and positively correlated with defect detection and naturalness.

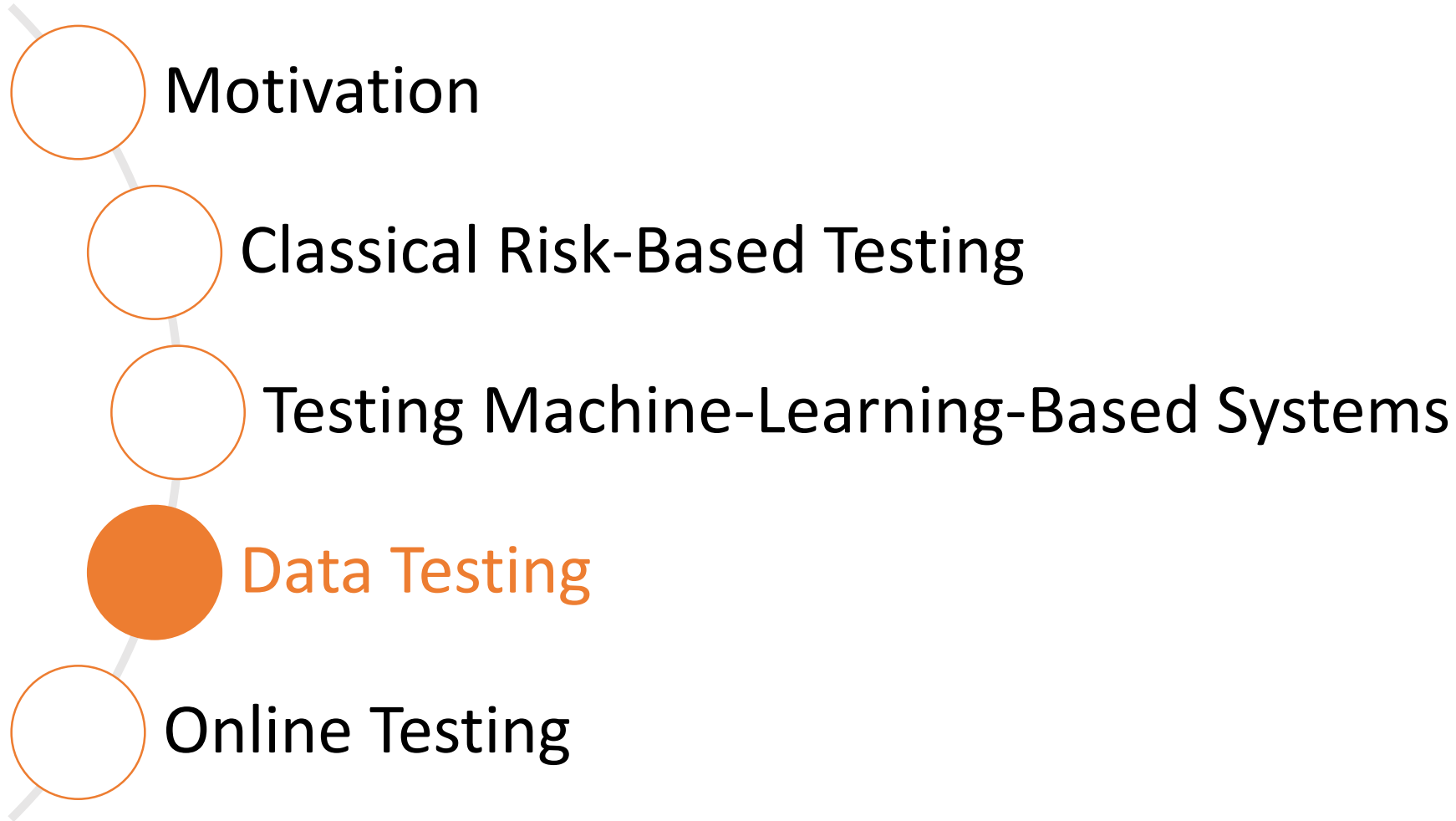
# Test Levels for Machine-Learning Systems

## Four Test Levels

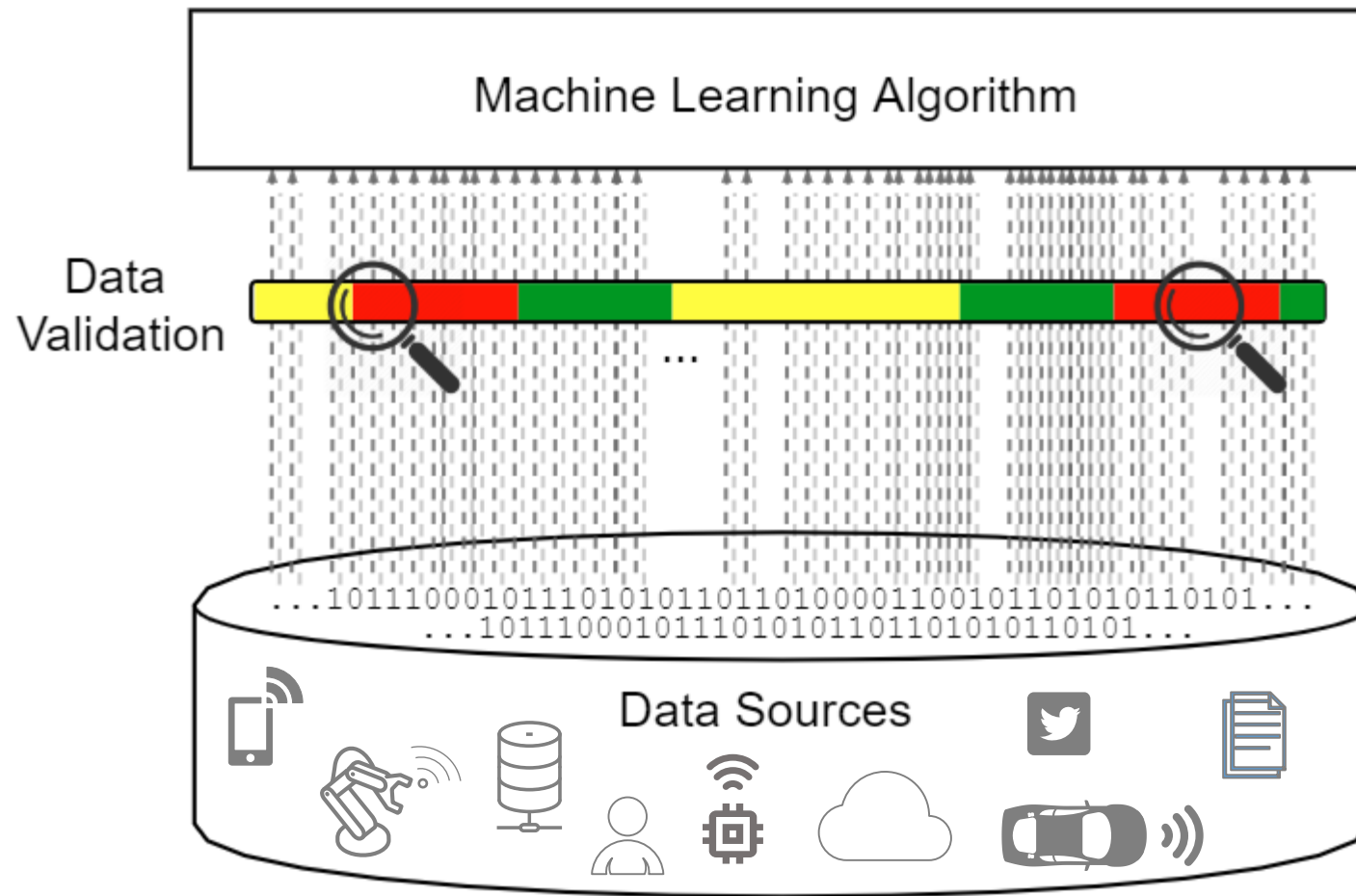
- Input Testing
- Model Testing
- Integration Testing
- System Testing



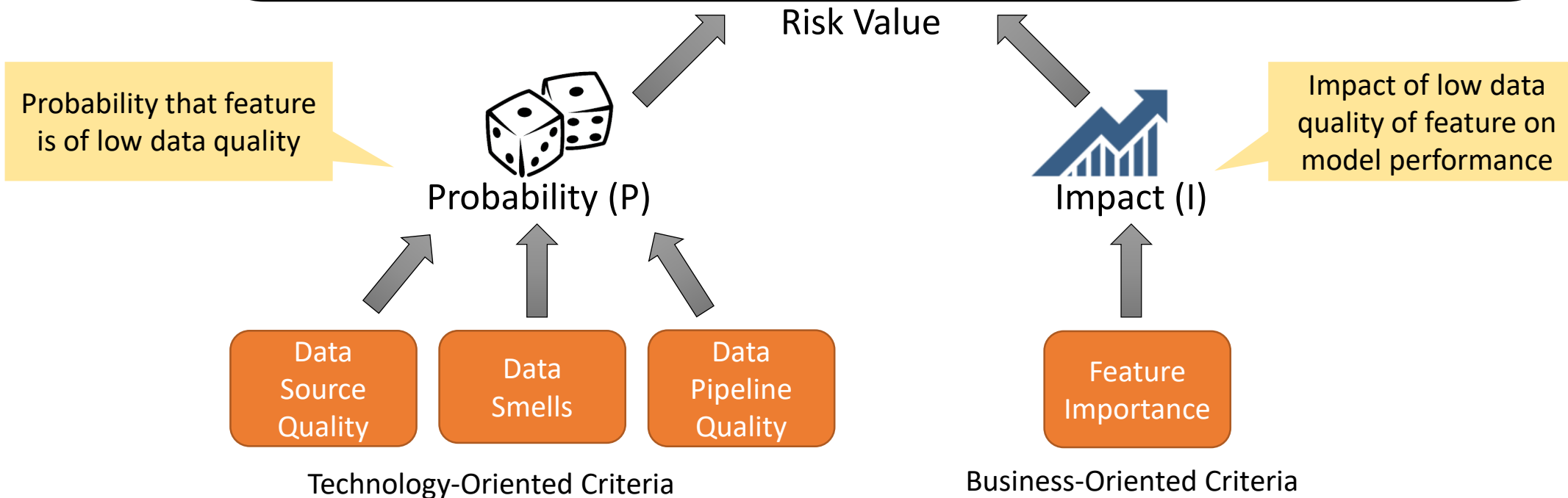
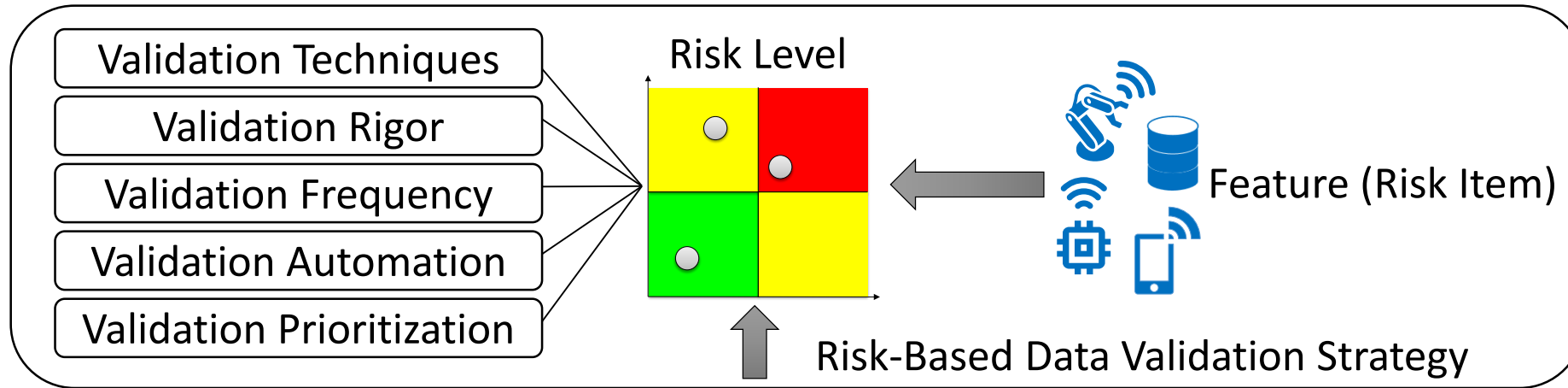
# Agenda



# Idea of Risk-Based Data Testing



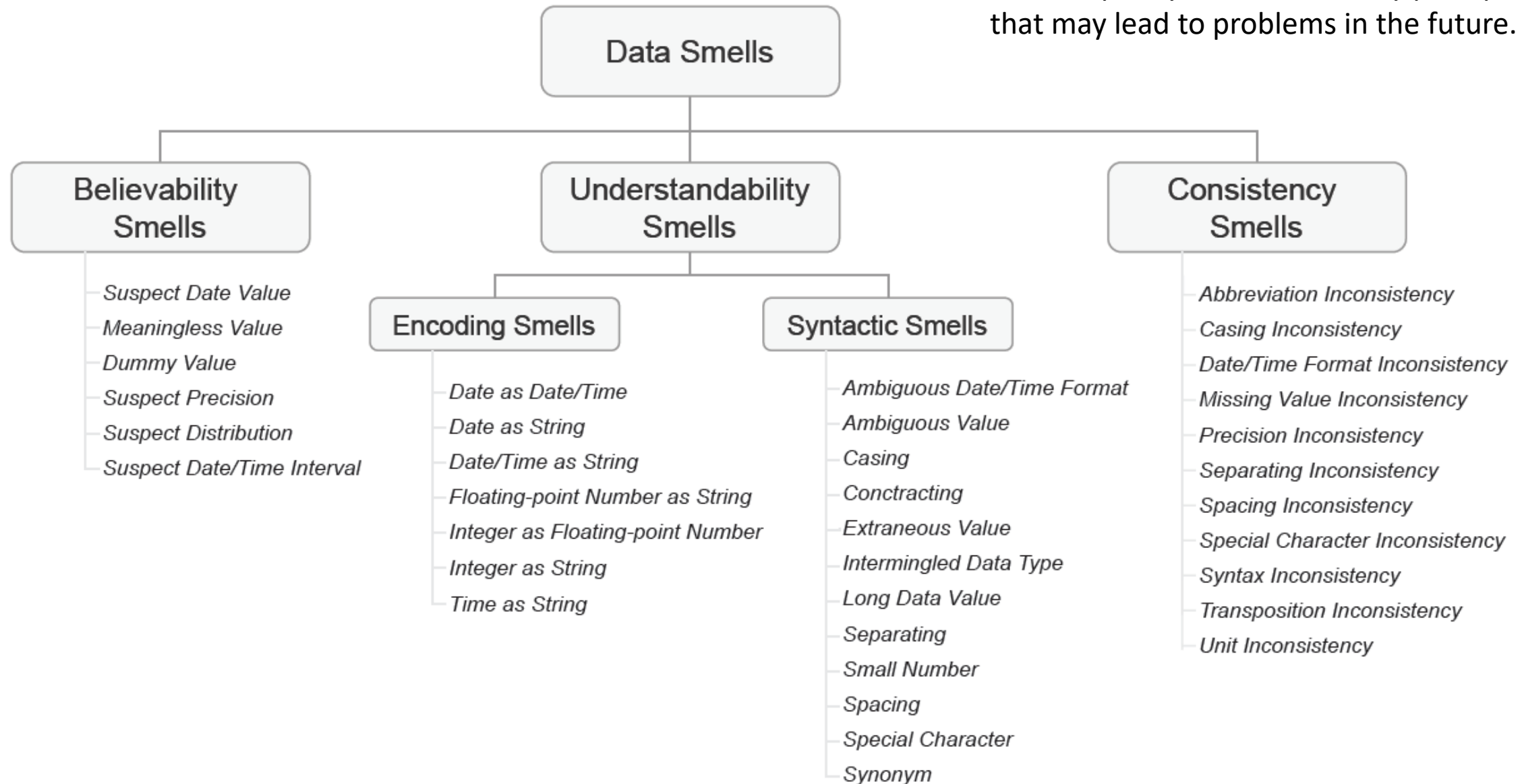
# Risk-Based Data Testing Framework





# Data Smell Types

Context-independent, data value-based indications of latent data quality issues caused by poor practices that may lead to problems in the future.



# Example Smells

Appointme...	Appointme...	# Age
5638447	2016-04-29T00:00:00Z	21
5629123	2016-04-29T00:00:00Z	19
5630213	2016-04-29T00:00:00Z	30
5620163	2016-04-29T00:00:00Z	29
5634718	2016-04-29T00:00:00Z	22
5636249	2016-04-29T00:00:00Z	28

Date as Date/Time Smell

<https://www.kaggle.com/joniarroba/noshowappointments>

Date	# Daily Confi...	# Total Confi...
11-Feb	0	3
12-Feb	0	3
13-Feb	0	3
14-Feb	0	3
15-Feb	0	3
16-Feb	0	3
17-Feb	0	3
18-Feb	0	3

Ambiguous Date/Time Format Smell

<https://www.kaggle.com/ravichaubey1506/covid19-india>

# Data Smell Detection Tools

The screenshot shows the DSD web interface. The top navigation bar includes 'Upload file', 'Customize', and 'Results' buttons. The main content area displays 'Smell Results By Column Names' with a table of results.

DATA SMELL TYPE	TOTAL ELEMENT COUNT	FAULTY ELEMENT COUNT	FAULTY ELEMENT O
Long Data Value Smell	891	687	[nan, nan, nan, nar
Casing Smell	891	24	['C23 C25 C27', 'F C
Integer as Floating Point	891	103	['C123', 'G6', 'C23 C

Rule-based detection

The screenshot shows the ML Data Smell Detection web interface. The main heading is 'Data Smell Detection with Machine Learning'. It displays 'LSTM Detection Results' for a sample dataset.

### LSTM Detection Results

**Agent:** Sample: Date Smell Classification using LSTMs  
**Dataset:** Sample Dataset: LSTM Date Classification

[Download Results](#)

### LSTM Classification

Shown below is the class distribution of your data as well as examples for each class. By default, the classes are labeled by the corresponding data smells according to the research paper. This can be disabled on the analyze page. If the class distribution does not match up your expectations, please download the corresponding dataset to further inspect the classification.

Class DateTime as String Smell	Class Date as DateTime Smell
This class contains <b>48</b> total samples, or <b>3.45%</b> of the total data. See some examples below for what data has been classified in this class.	This class contains <b>48</b> total samples, or <b>3.45%</b> of the total data. See some examples below for what data has been classified in this class, and w

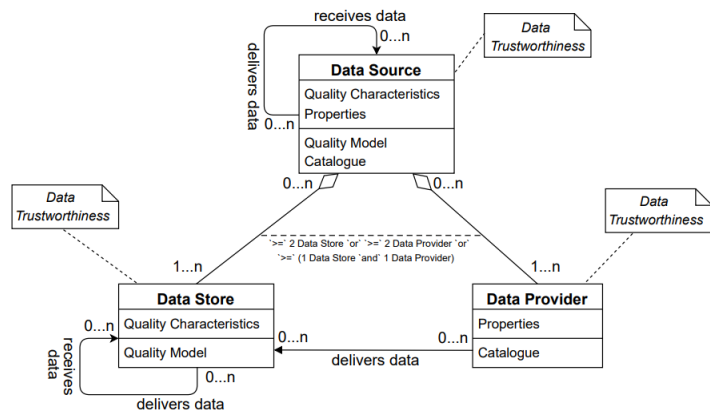
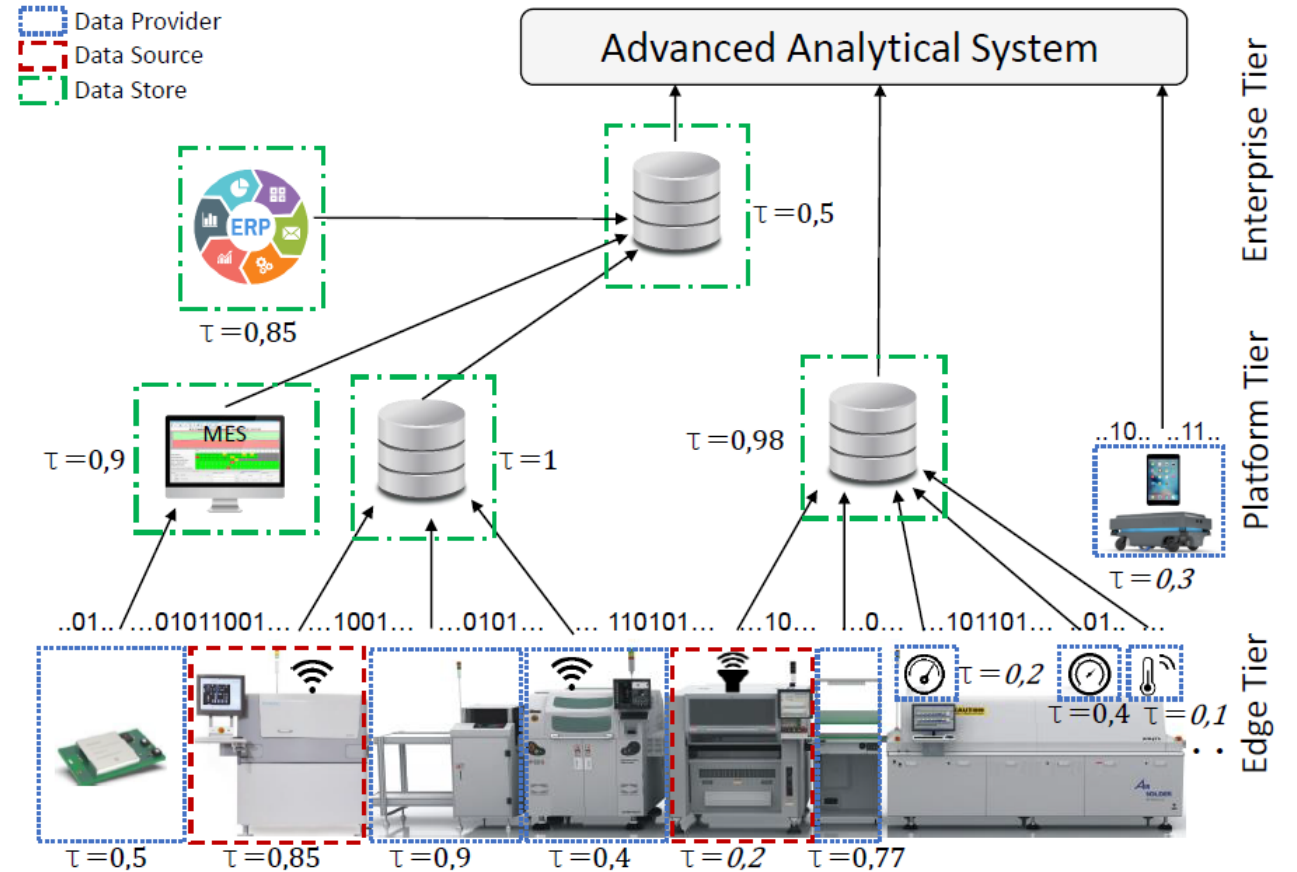
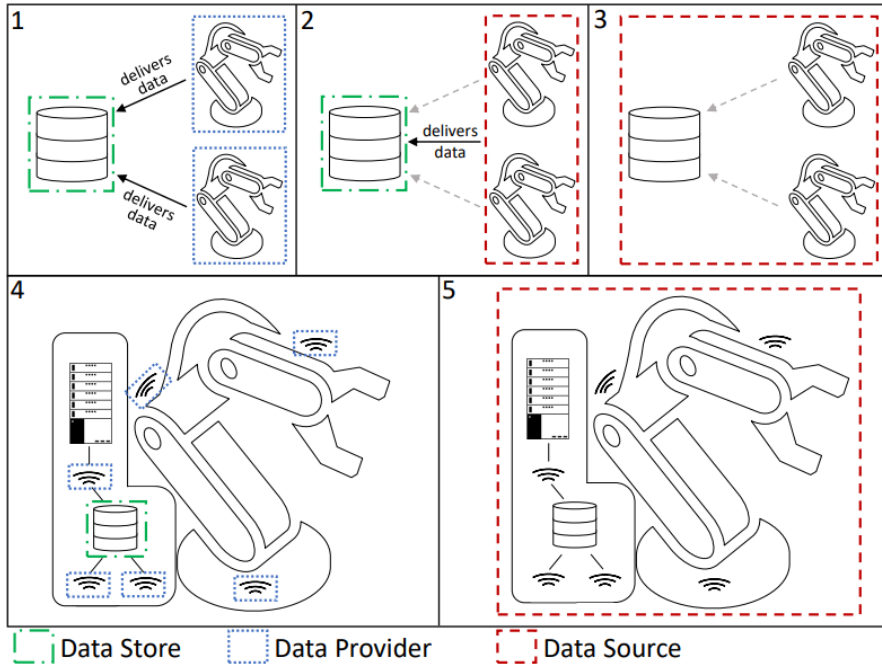
- "03-MAR-1994 09:57PM+06:00" (100.0%)
- "25-APR-2001 00:14+10:00" (100.0%)
- "11-JUN-1977 02:03:44.0051AM +06:00" (100.0%)
- "18-JUN-2015 09:09:25.0453AM +04:00" (100.0%)
- "20-APR-1999 02:51:26AM+03:00" (100.0%)
- 10-APR-1992 00:00+00:00 (99.93%)
- 15-DEC-1993 00:00+00:00 (99.92%)
- 01-OCT-1999 12:00:00.0000AM +00:00 (100.0%)
- 30-MAR-1995 12:00:00.0000AM +00:00 (100.0%)
- 06-JAN-1976 00:00:00+00:00 (99.93%)

Machine learning-based detection

# Data Source Quality Model

Quality Characteristics		Description	Properties and Subquality Factors
Representational Data Store Quality	<b>Representational Adequacy</b>	Degree to which a data store presents data in a <i>concise and organized way</i> .	Schema Minimality Schema Normalization Schema Pertinence
	<b>Representational Consistency</b>	Degree to which a data store presents data <i>always in the same format and compatible with previous data</i> .	Data Format Variety Data Type Variety Schema Change Proneness
	<b>Understandability</b>	Degree to which <i>users can understand the data</i> provided by a data store.	Data Format Complexity Documentation Degree Metadata Quality Schema Readability
Dynamical Data Store Quality	<b>Accessibility</b>	Degree to which data are <i>easily and quickly retrievable</i> .	Access Maturity Operability Retrievability
	<b>Availability</b>	Degree to <i>which data are available</i> from a data store.	Durability Fault Tolerance Recoverability Scalability Uptime
	<b>Security</b>	Degree to which <i>access to data for unauthorized persons is restricted</i> by a data store.	Authentication/Encryption Authorization Policy
	<b>Timeliness</b>	Degree to which a data store <i>provides up-to-date data in a timely manner</i> .	Refresh Rate Response Time
Statical Data Store Quality	<b>Completeness</b>	Degree to which a data store is able to represent <i>every meaningful state of the real world</i> .	Schema Completeness Schema Correctness
	<b>Contactability</b>	Degree to which a data store provides <i>contact information for further inquiries</i> .	Support Degree Complexity
	<b>Trustworthiness</b>	Degree to which a <i>data store can be trusted</i> .	Data Governance Maturity Verifiability

# Data Source, Provider, and Store



# Data Pipeline Quality

Issues of the implemented data pipeline that can affect the quality of the processed data

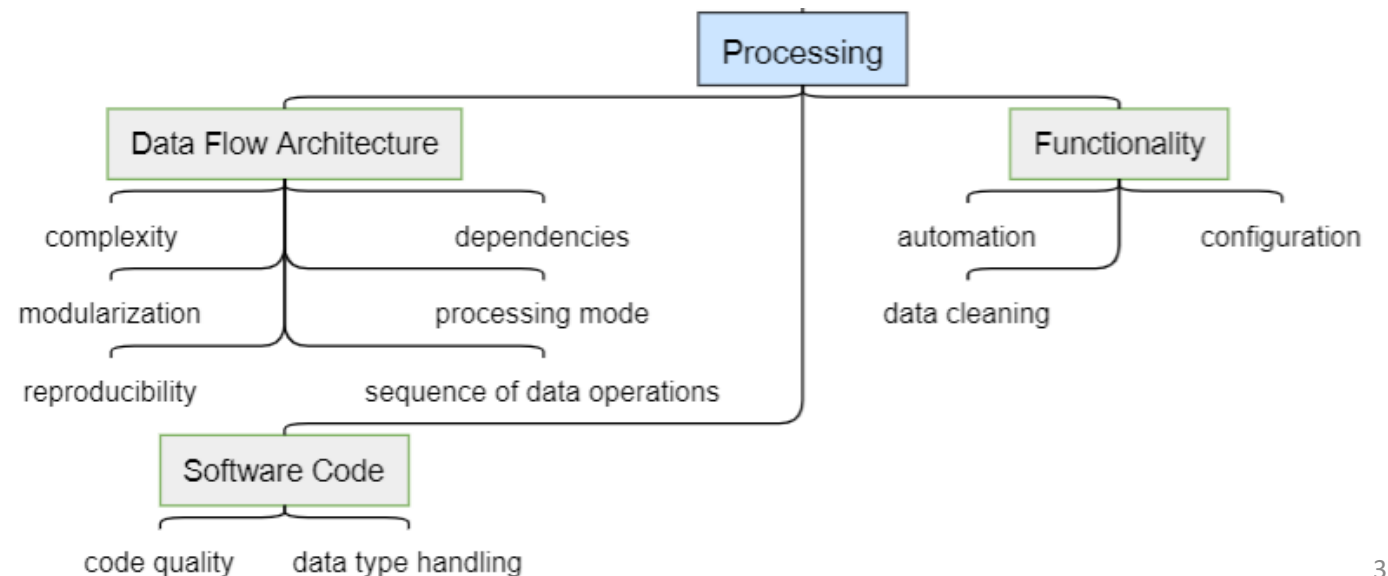
Code smells

- libraries Pandas vs. Numpy handling missing values differently by default (e.g. `pandas.DataFrame.max` vs. `numpy.amax`)

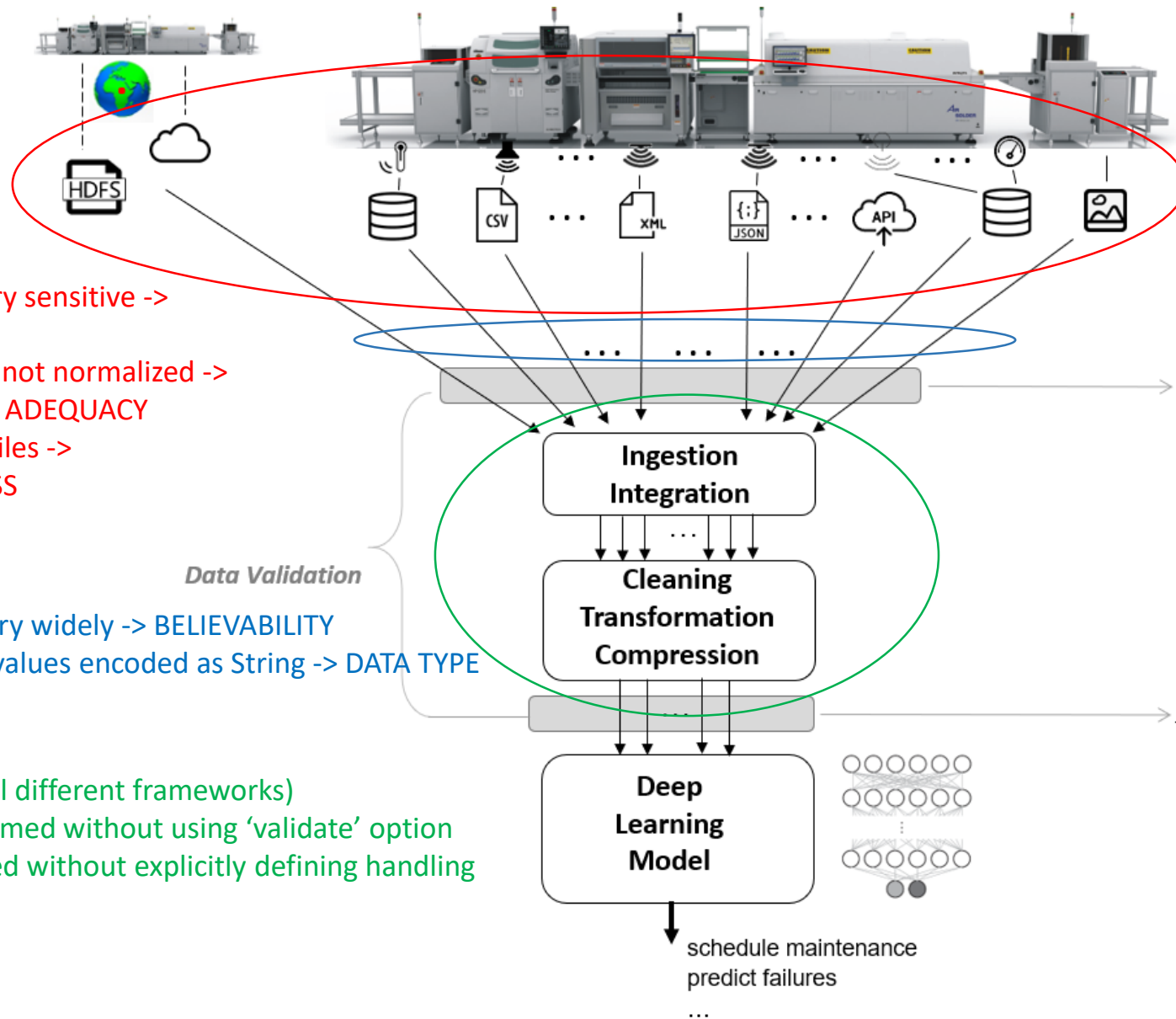
```
DataFrame.max(self, axis=None, skipna=None, level=None, numeric_only=None, **kwargs)
```

**skipna** : *bool*, default *True*

Exclude NA/null values when computing the result.



# Predictive Maintenance Example



## Data Sources:

- capacitance sensors
- current sensors
- acoustic emission sensors
- ...
- databases
- flat files
- ...

## Data Signals:

- acoustic data signals
- temperature data signals
- vibration data signals
- ultrasonic signals
- voltage signals
- ...

## Features:

- rotation speed
- vibration of generator shaft
- bearing wear
- voltage imbalance
- ...

## Data Source Quality

- capacitance sensor very sensitive -> AVAILABILITY
- database of generator not normalized -> REPRESENTATIONAL ADEQUACY
- Personal data in XML files -> TRUSTWORTHINESS
- ...

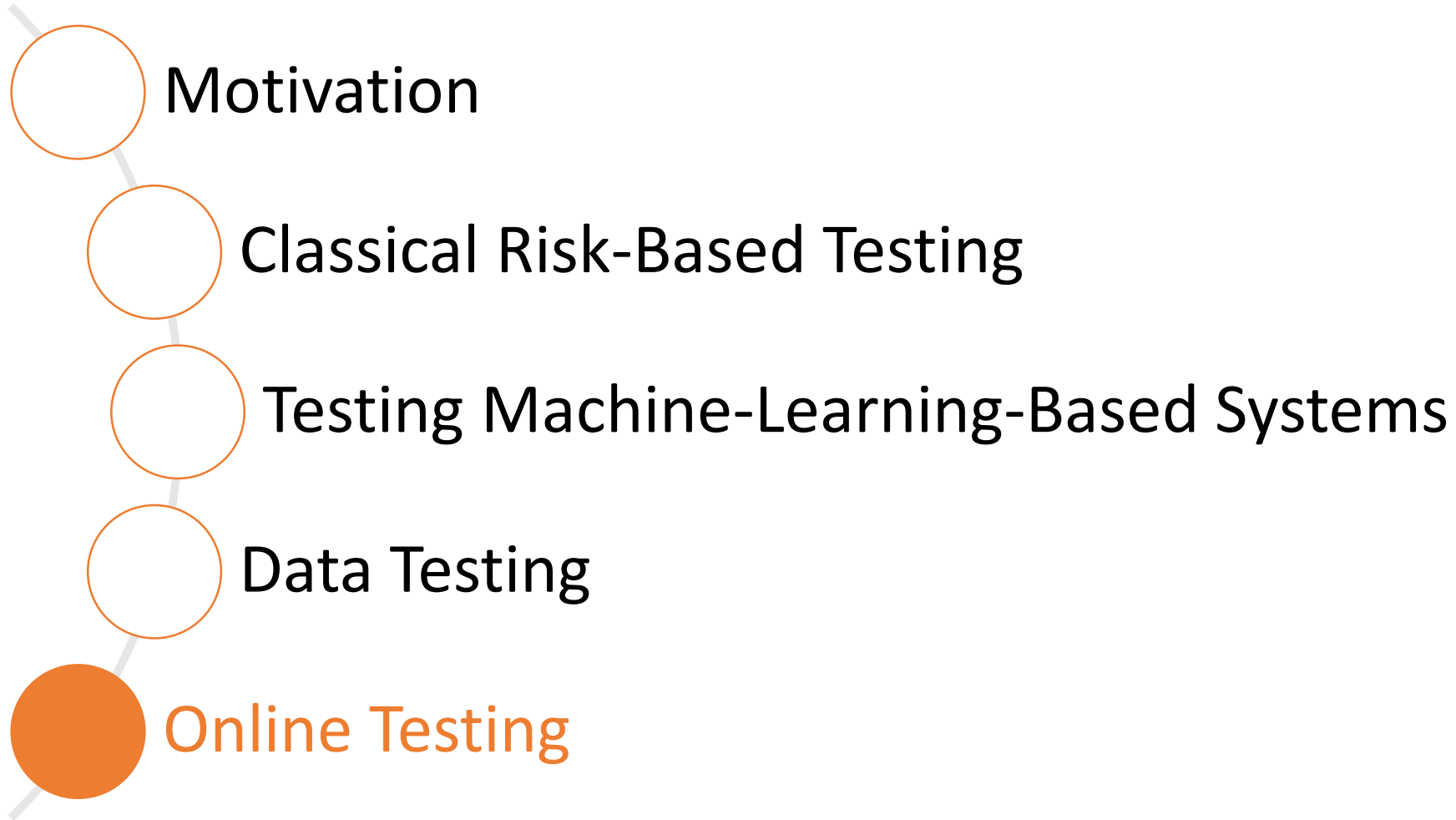
## Data Smells

- temperature values vary widely -> BELIEVABILITY
- timestamp of current values encoded as String -> DATA TYPE
- ...

## Data Pipeline Quality

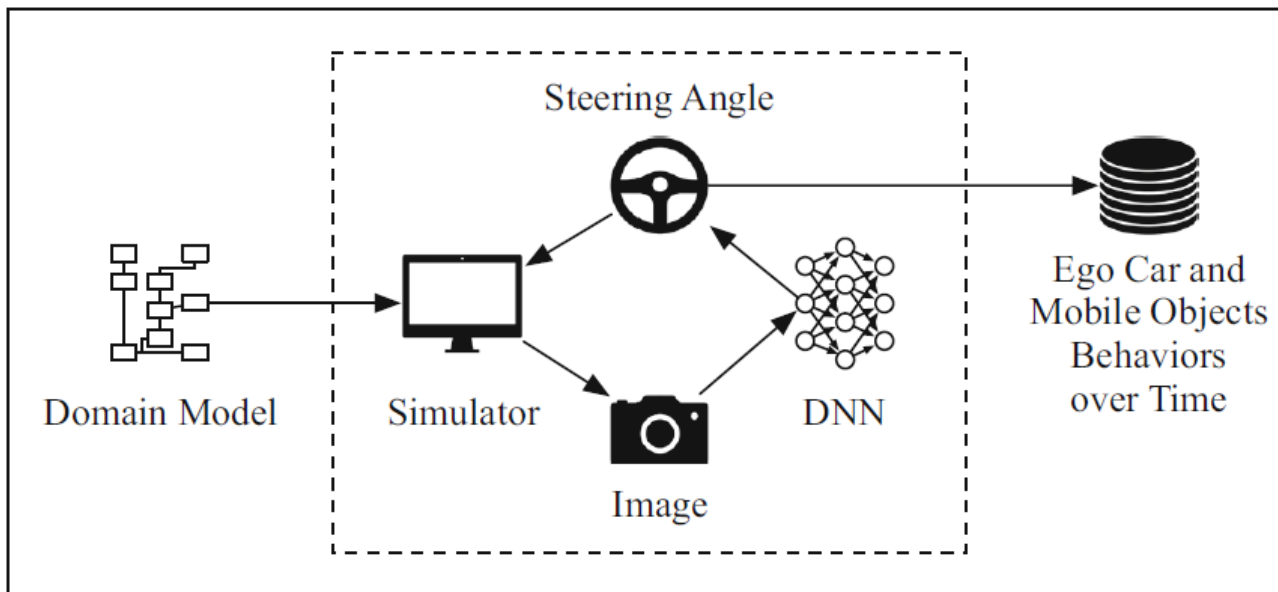
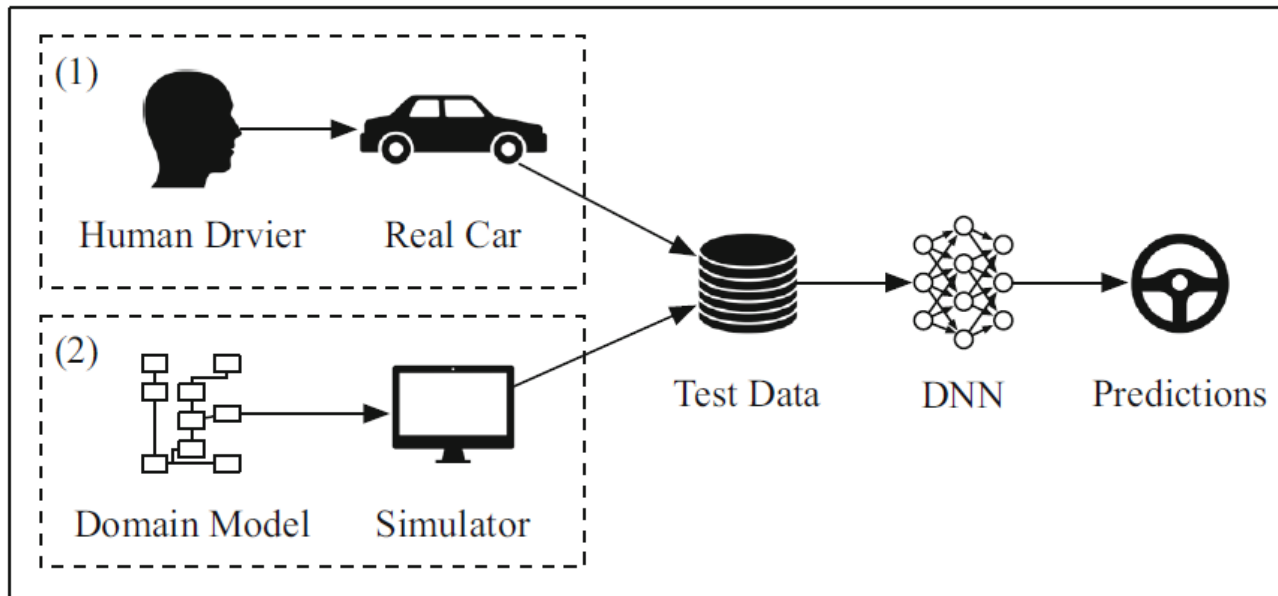
- pipeline jungle (several different frameworks)
- vibration data transformed without using 'validate' option
- voltage data normalized without explicitly defining handling of missing values
- ...

# Agenda

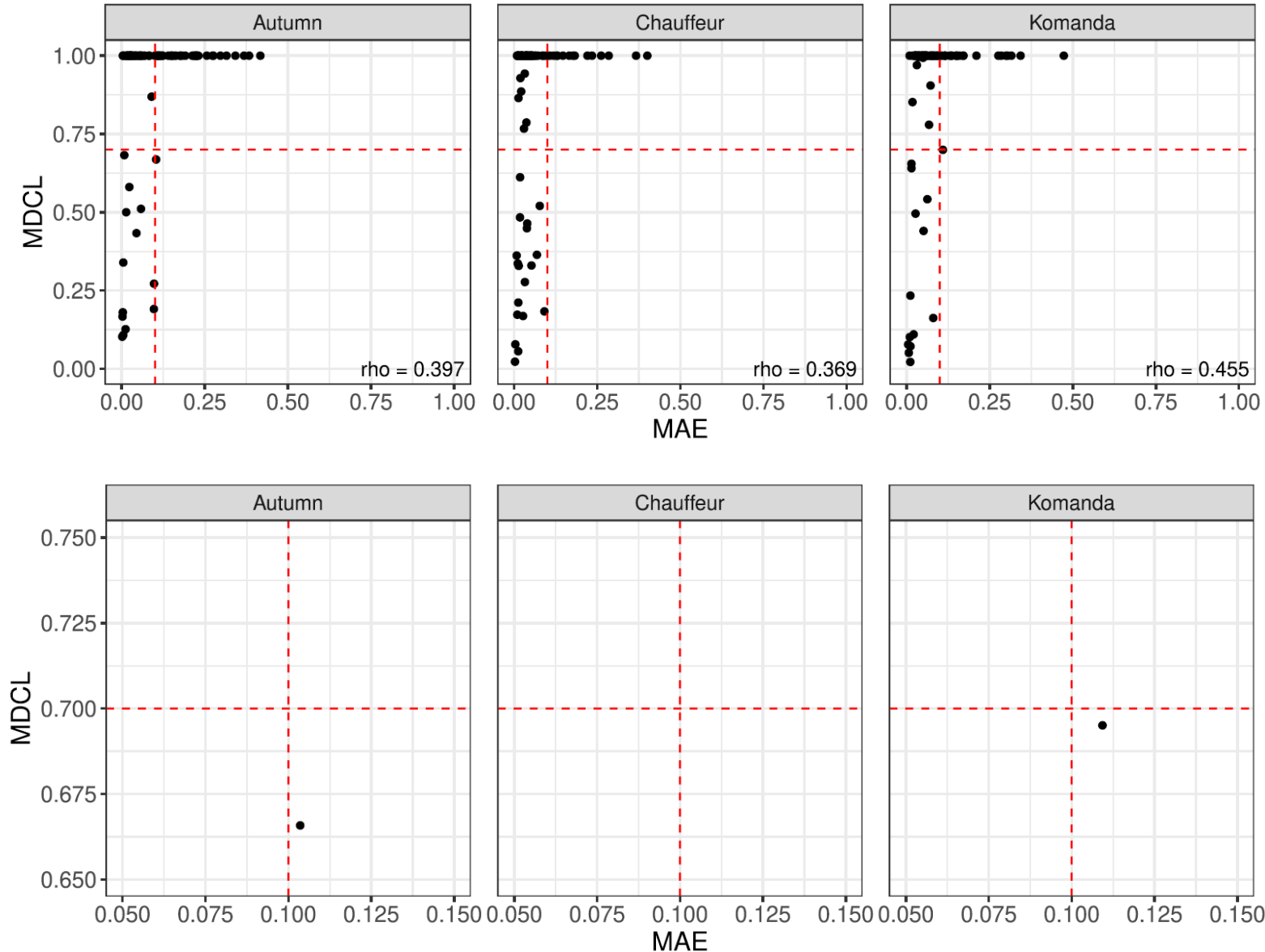




# Offline and Online Testing



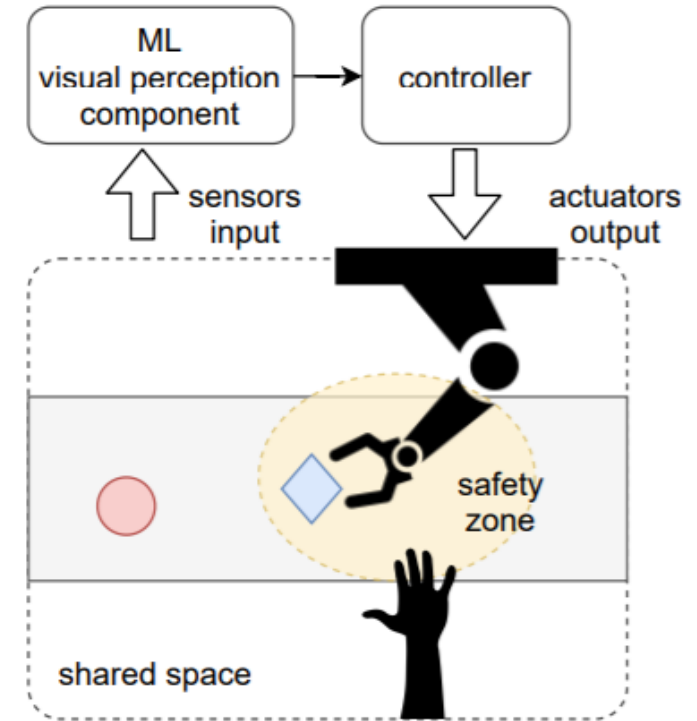
# Comparison Offline and Online Testing Results



Many safety violations identified by online testing could not be identified by offline testing

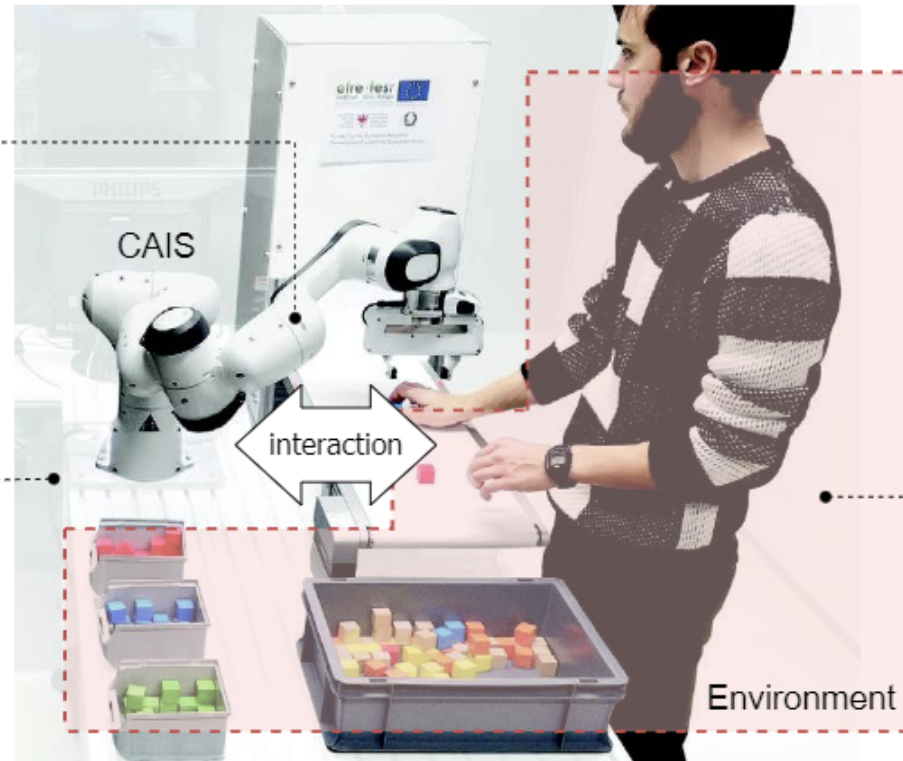
Offline testing cannot properly reveal safety violations

# Collaborative AI Setting and Safety



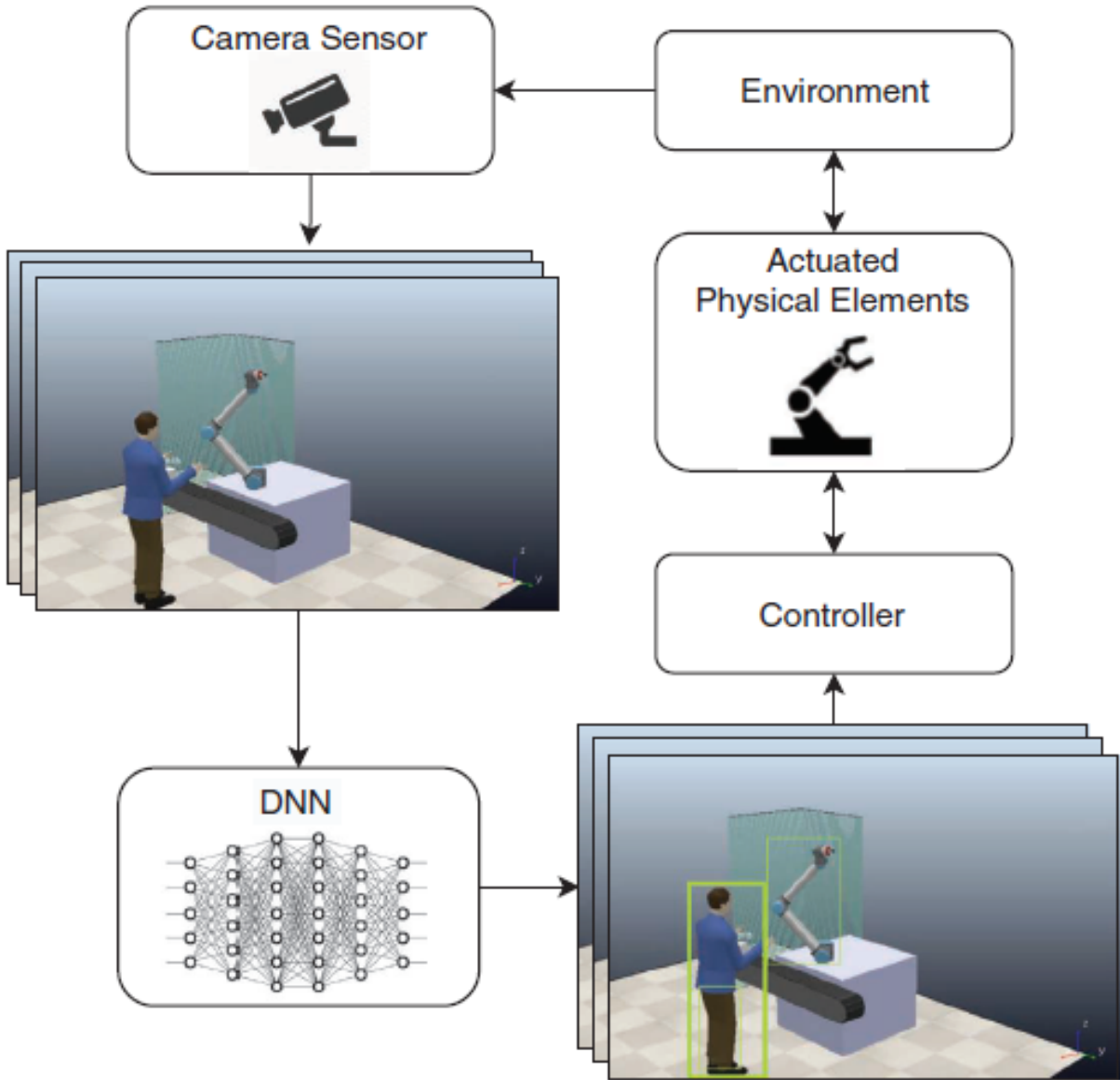
**Online learning:**  
object classification,  
human classification,  
motion direction,  
motion speed

**Risks:**  
protective distance violation,  
injuring behavior,  
robot unable to classify objects,  
robot prevails over human needs

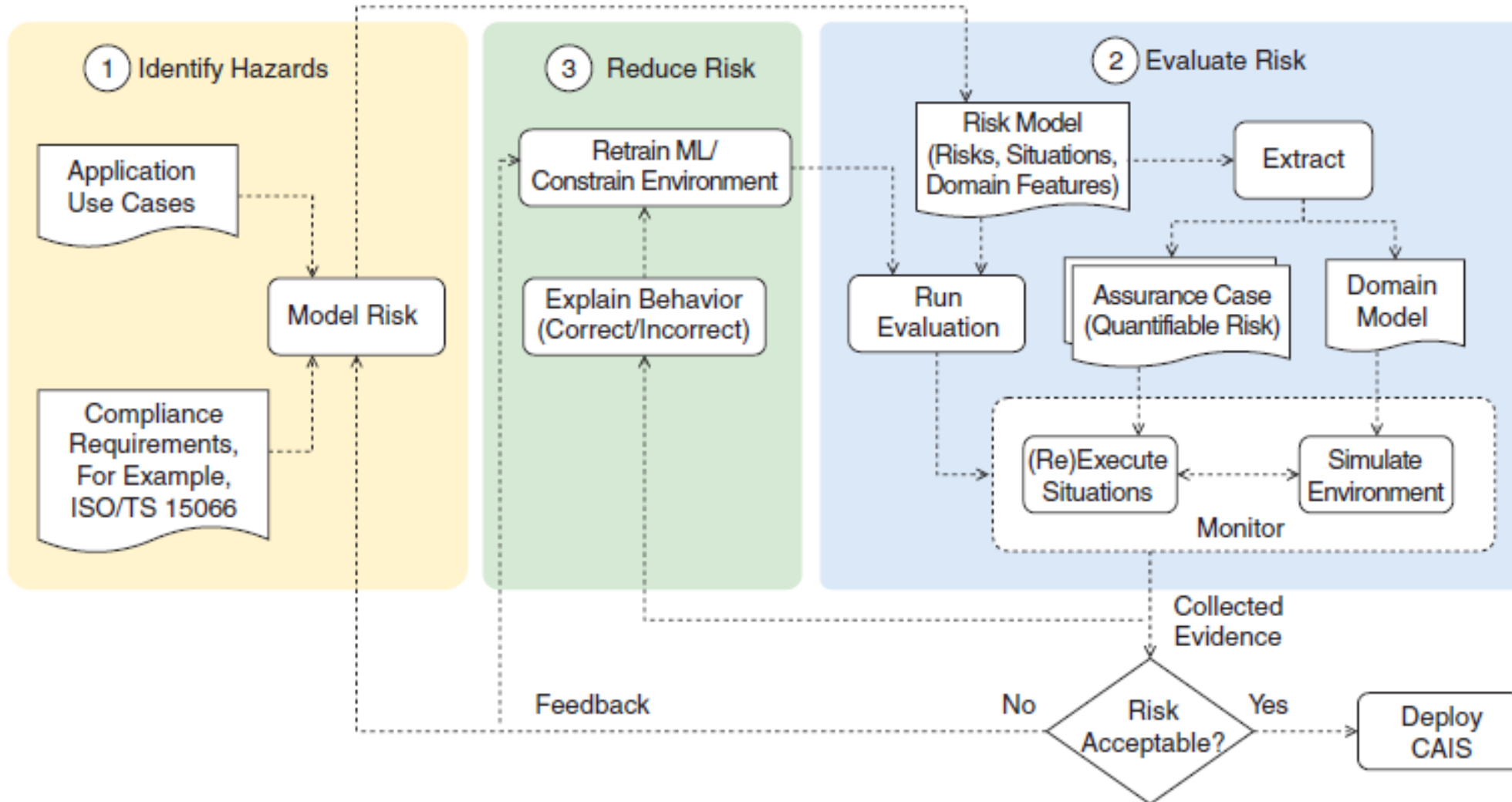
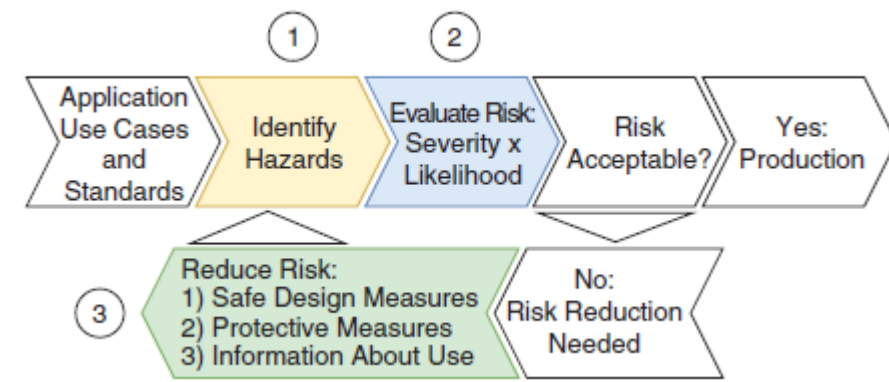


**Uncertainties:**  
human position,  
human motion speed,  
human-background contrast,  
luminance,  
shape/color of objects

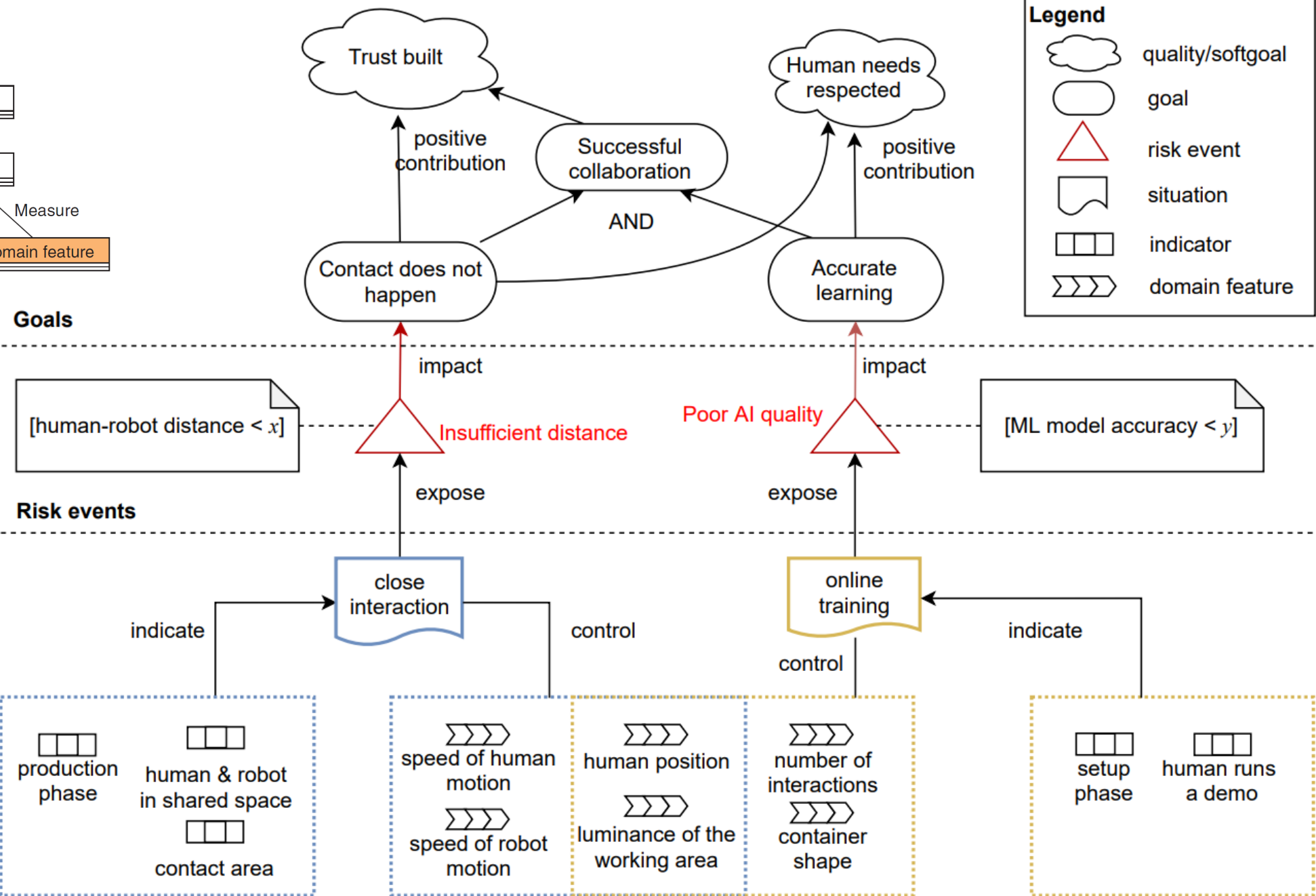
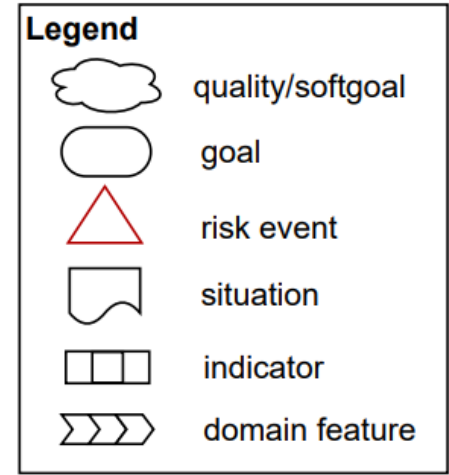
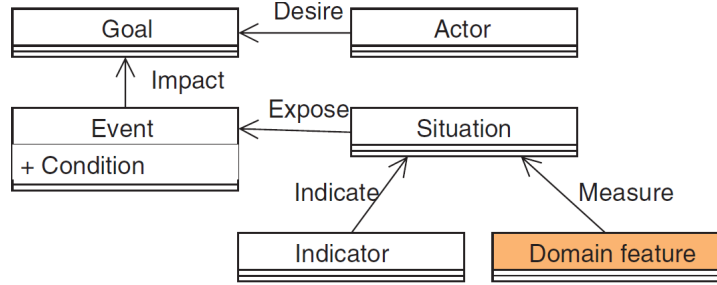
# Setting for Online Testing



# Risk-Driven Assurance Process



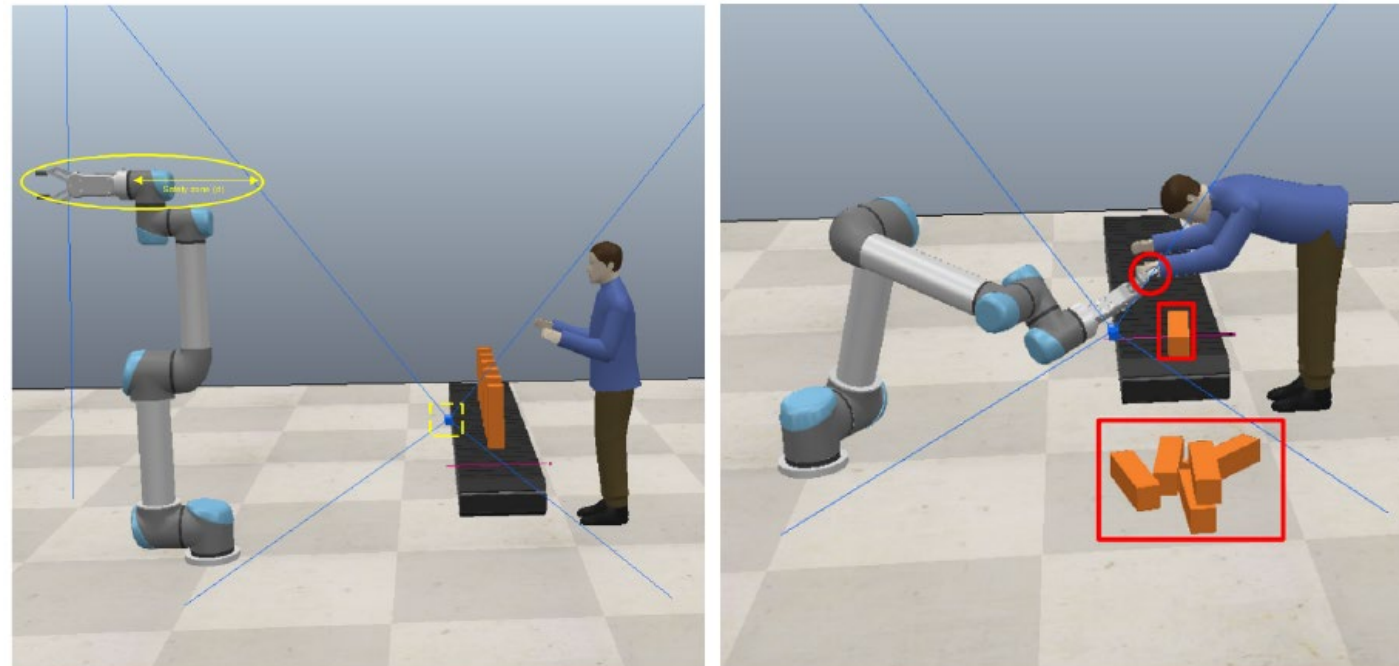
# Risk Model



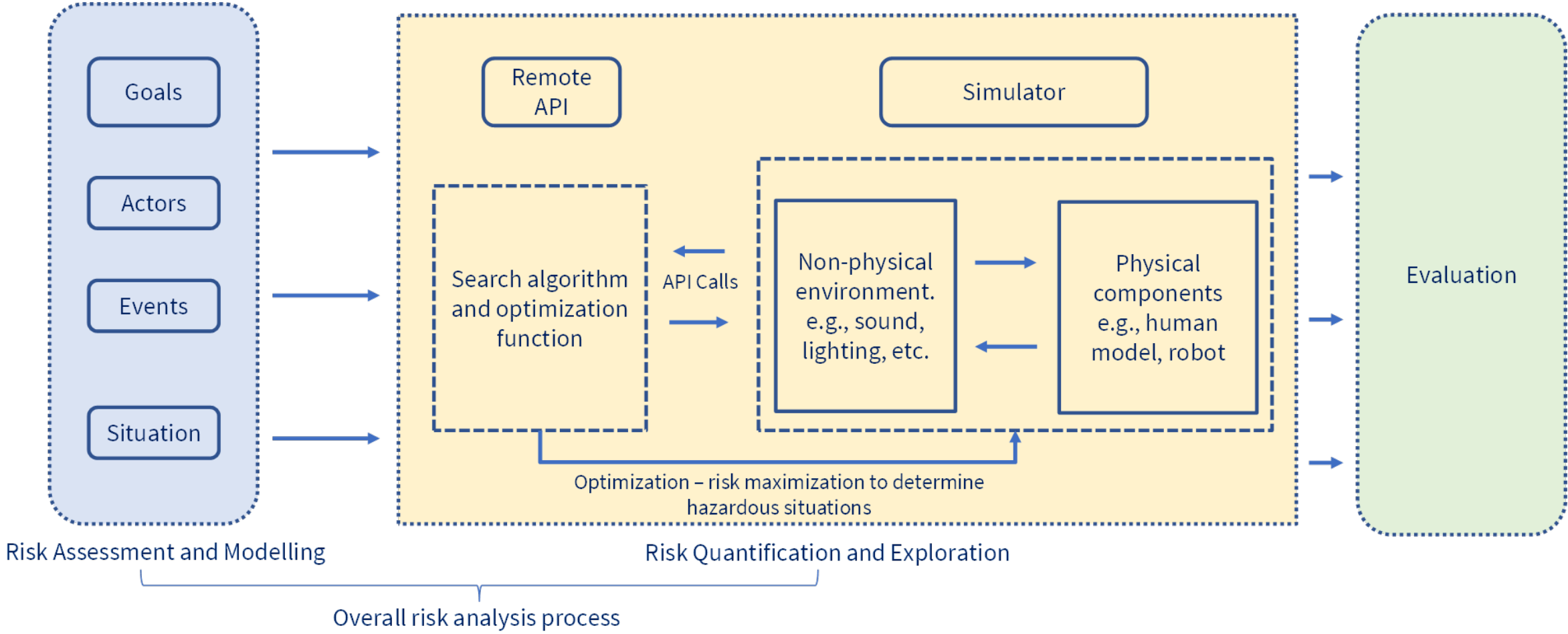
# Search Space Exploration and Risk Evaluation

- [Start]** Initialize simulator variables for path e.g., Given initial system states ( $s_0, n, u$ ) and “m” number of runs such that:
  - $s_0$  = robot arm coordinates ( $R [x, y, z]$ ), human position ( $H [x, y, z]$ ),
  - $n = 0 : 0, 1, 2 \dots n \leq m$
  - $U = []$
- [Input space]** “p” population of key simulator features (suitable solutions for the problem) which is the expected output:
  - $F' \leq F$  (e.g., environment lighting ( $L [h, s, l]$ ))
  - Path trace during its execution ( $v = (v_1, \dots, v_p)$ ), implicitly determined within simulator
- [Fitness]** Determine fitness function  $f(x)$  of each arrangement in the population that satisfies a value  $\phi$  along the path  $v$ , e.g.,
  - $d(R, H) \leq \phi$
- [Update]** The simulation rewards those path functions where an unsafe state is reached for instance with a “1” or otherwise “0”, returns:
  - feature  $F' \subseteq U$
  - $$U_r = \begin{cases} 1, & dvi(R, H) \leq \phi \\ 0, & otherwise \end{cases}$$
- [Increment]** test case  $n = n + 1$
- [Repeat/Loop]** : Step 2 until  $n = m$
- Return** outcome  $U$

Objective Functions:  
 Minimum distance between human and robot arm  
 Relative speed of human and robot arm

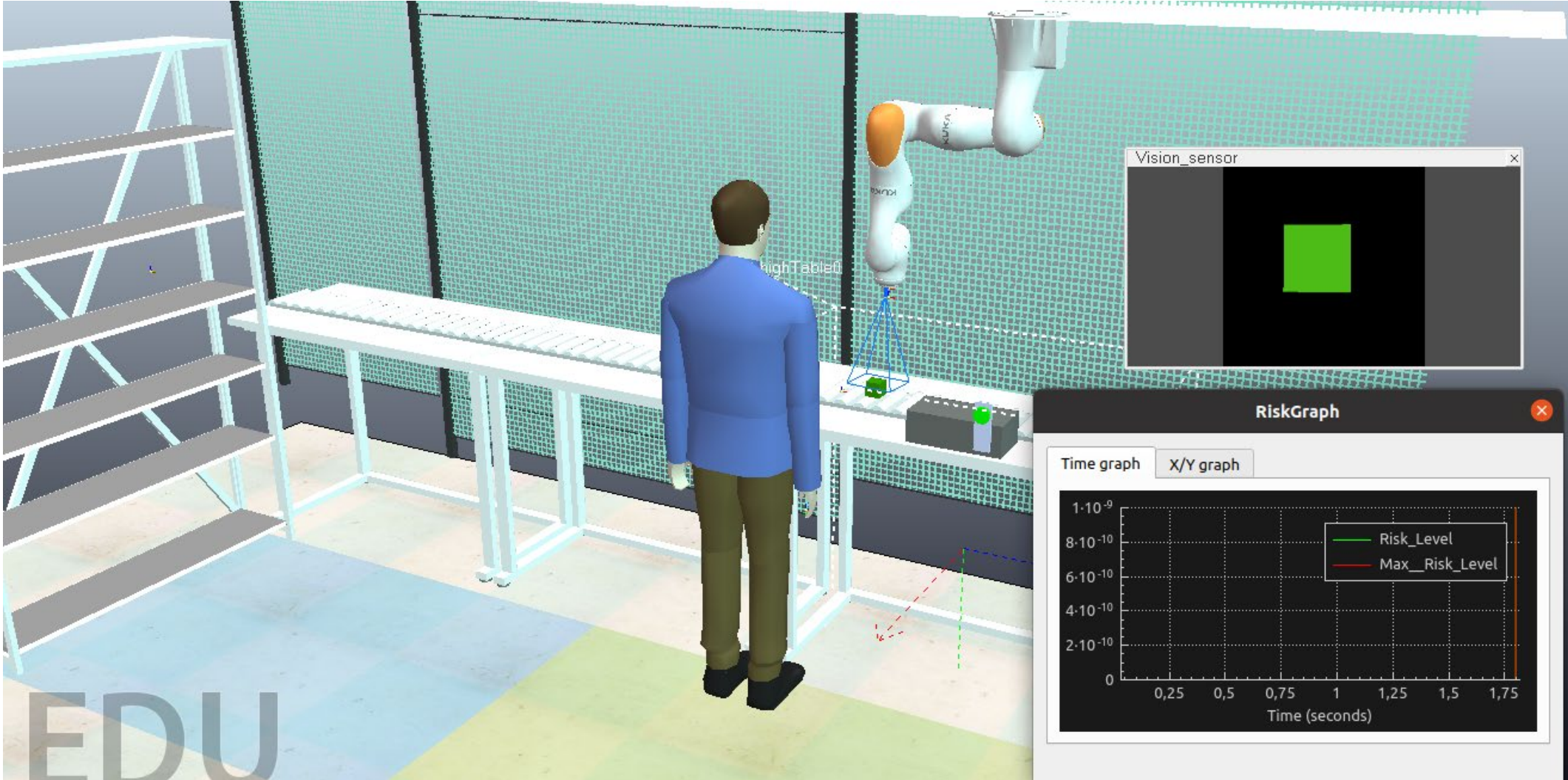


# Test Environment



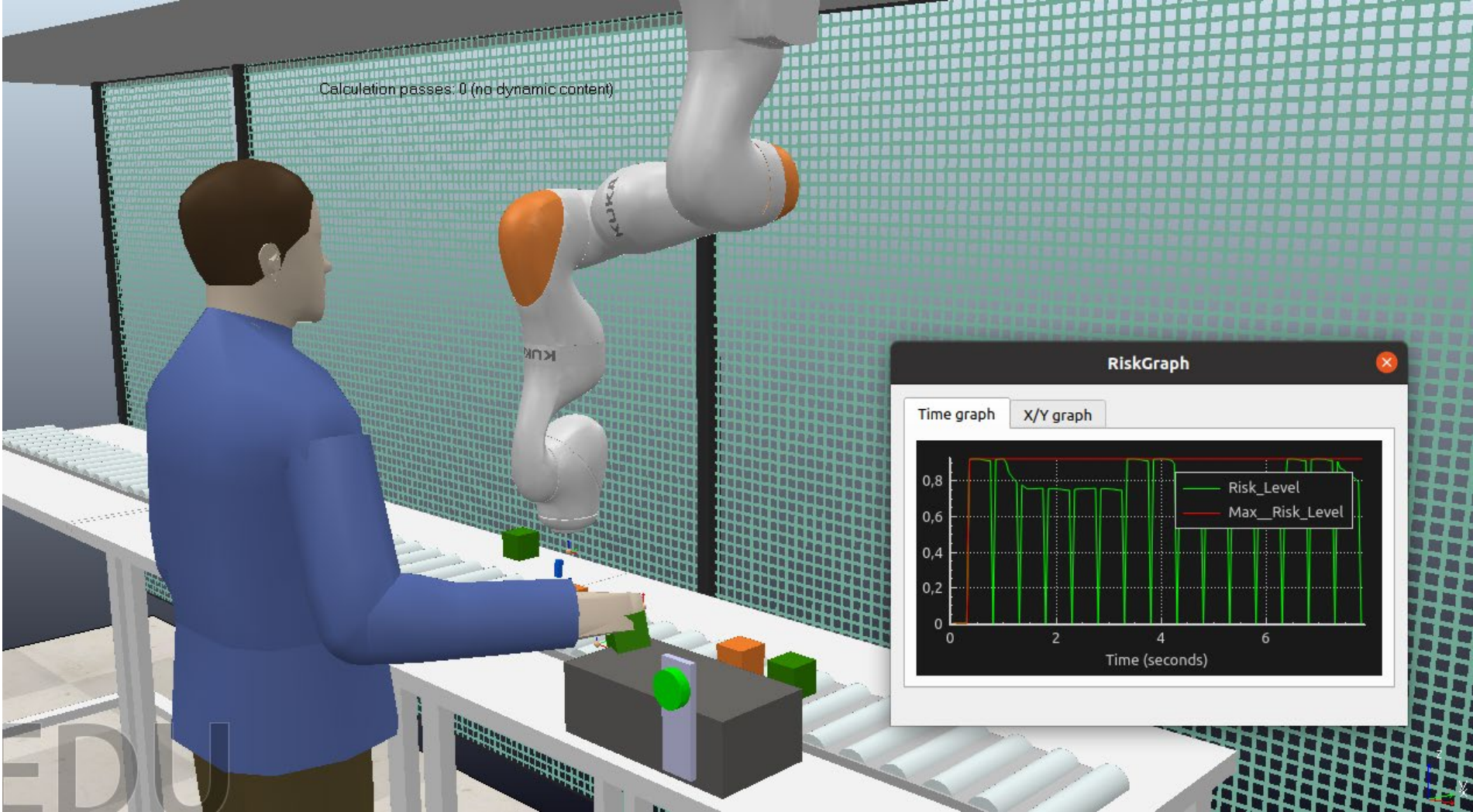


# Simulation in CoppeliaSim (1/2)

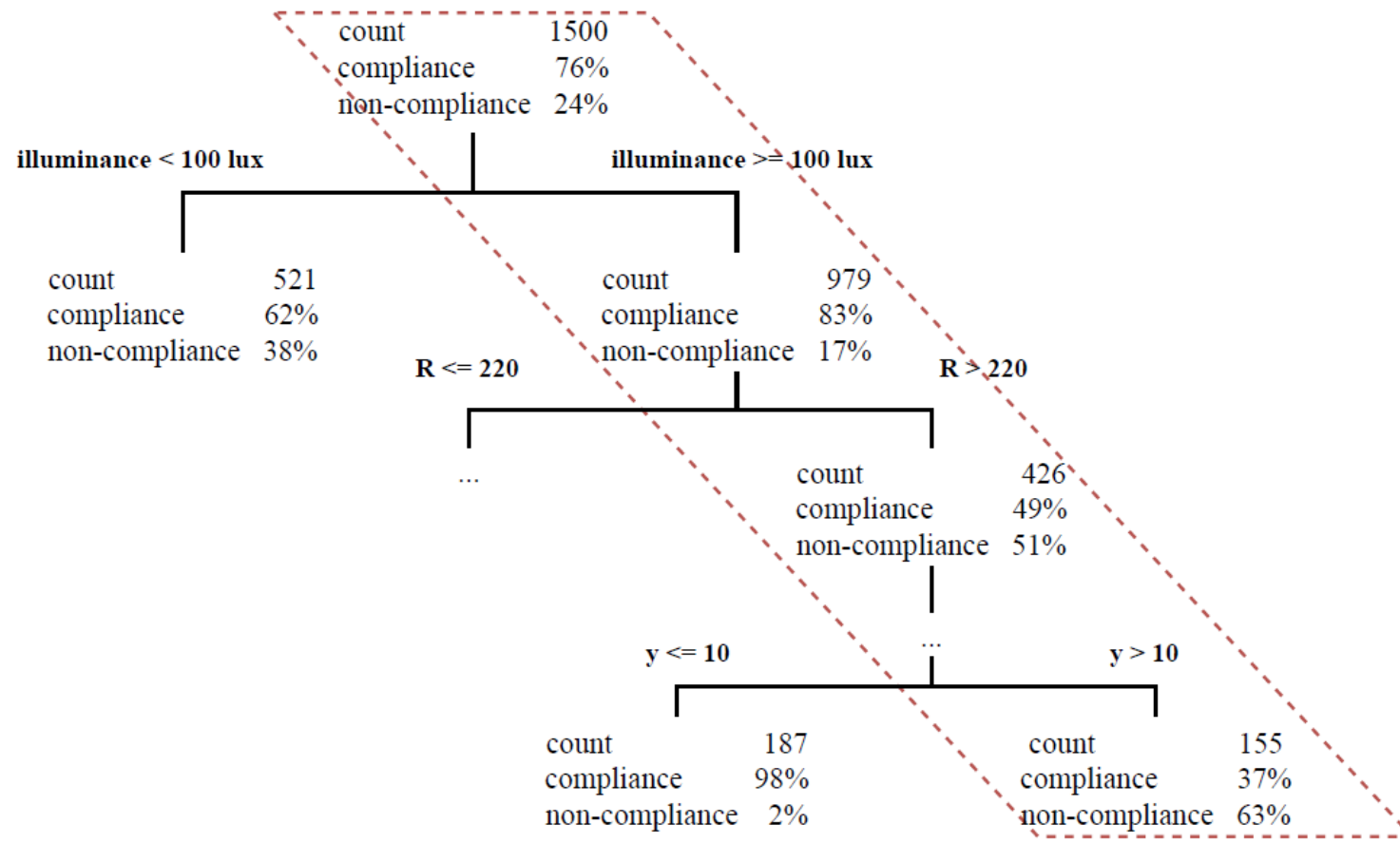


EDU

# Simulation in CoppeliaSim (2/2)



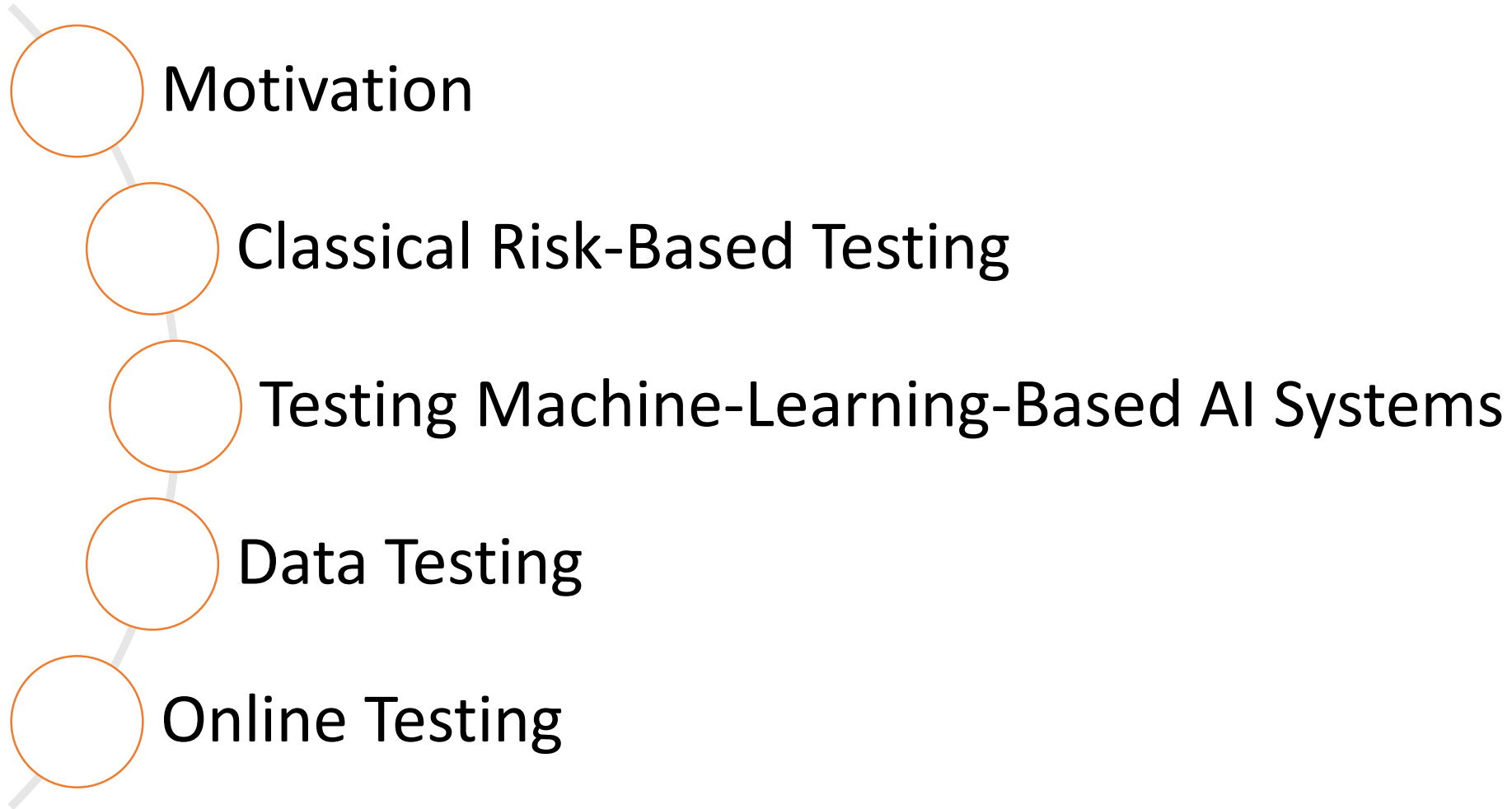
# Risk Quantification, Decision Tree and Rule Extraction



rule example #2

rule example #	illuminance (lux)	domain features operator arms color (R,G,B)	operator position (x,y)
1	<100	-	-
2	>100	R > 220, G > 236, B > 200	x > 180, y > 10

# Summary



# References

- [1] Zhang, Jie M., et al. *Machine learning testing: Survey, landscapes and horizons*. IEEE Transactions on Software Engineering, 2020
- [2] Riccio, Vincenzo, et al. *Testing machine learning based systems: a systematic mapping*. Empirical Software Engineering 25.6 (2020): 5193-5254.
- [3] Felderer, M., Schieferdecker, I. *A taxonomy of risk-based testing*. Software Tools for Technology Transfer, 16(5), 559-568, Springer, 2014
- [4] Foidl, H., Felderer, M.: *Integrating software quality models into risk-based testing*. Software Quality Journal, 26(2), 809-847, 2018
- [5] Foidl, H., Felderer, M., Ramler, R. *Data Smells: Categories, Causes and Consequences, and Detection of Suspicious Data in AI-based Systems*. 1st International Conference on AI Engineering (CAIN 2022), ACM, 2022
- [6] Haq, Fitash Ul, et al. *Can Offline Testing of Deep Neural Networks Replace Their Online Testing?*. Empirical Software Engineering 26:5, 1-30, 2021
- [7] Adigun, Jubril Gbolahan, et al. *Collaborative Artificial Intelligence Needs Stronger Assurances Driven by Risks*. Computer 55:3, 52-63, 2022
- [8] Harel-Canada, F., et al. *Is neuron coverage a meaningful measure for testing deep neural networks?*. Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 2020



## Testing the Untestable: Risikobasierte Qualitätssicherung für Machine-Learning Systeme

Prof. Dr. Michael Felderer

Department of Computer Science

Universität Innsbruck

Austria

 [michael.felderer@uibk.ac.at](mailto:michael.felderer@uibk.ac.at)

 [@mfelderer](https://twitter.com/mfelderer)