

Data

*Altruism*

*by Design*

Machine Learning, honestly.



"I want AI to do my laundry and dishes so that I can do art and writing, not for AI to do my art and writing so that I can do my laundry and dishes."

Author and videogame enthusiast **Joanna Maciejewska**  
(although bathroom cleaning goes)

# Models know What you did last summer

## A Survey of Privacy Attacks in Machine Learning

MARIA RIGAKI and SEBASTIAN GARCIA, Czech Technical University in Prague, Czech Republic

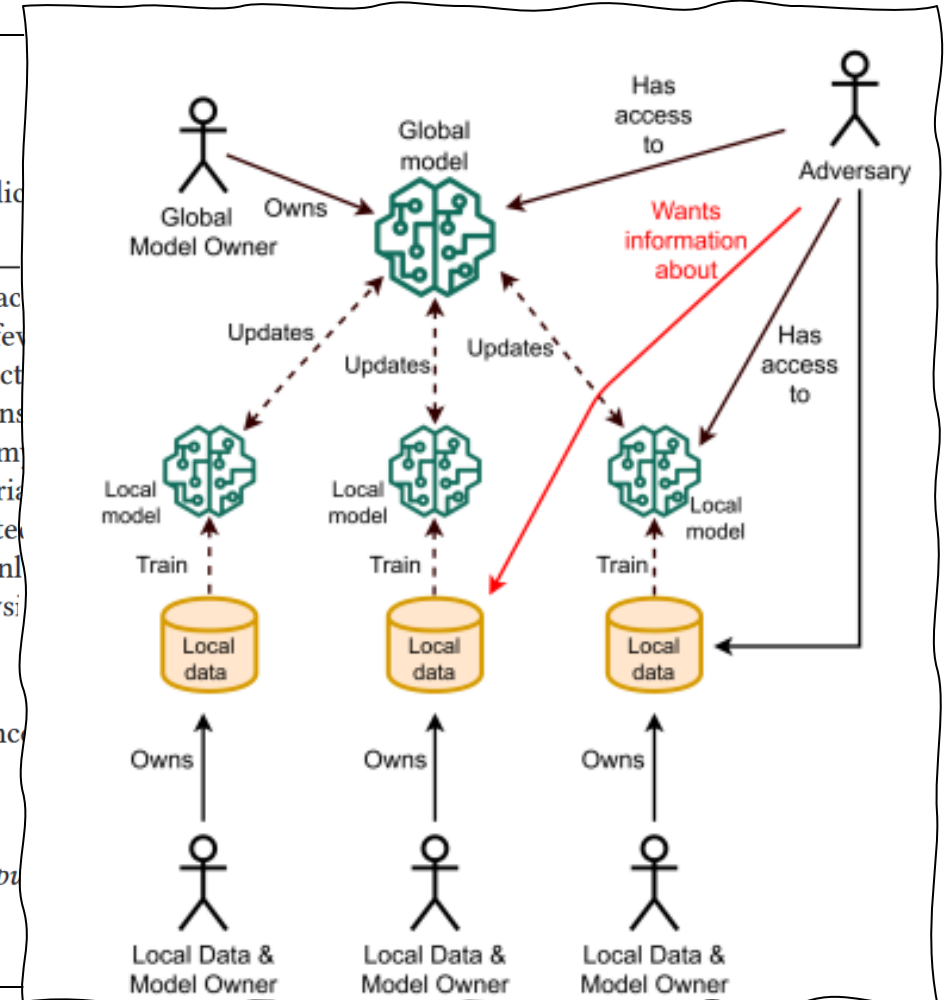
As machine learning becomes more widely used, the need to study its implications in security and privacy becomes more urgent. Although the body of work in privacy has been steadily growing over the past few years, research on the privacy aspects of machine learning has received less focus than the security aspect. Our contribution in this research is an analysis of more than 45 papers related to privacy attacks against machine learning that have been published during the past seven years. We propose an attack taxonomy together with a threat model that allows the categorization of different attacks based on the adversary's knowledge, and the assets under attack. An initial exploration of the causes of privacy leaks is presented as well as a detailed analysis of the different attacks. Finally, we present an overview of the most commonly proposed defenses and a discussion of the open problems and future directions identified during our analysis.

CCS Concepts: • **Computing methodologies** → **Machine learning**; • **Security and privacy**;

Additional Key Words and Phrases: Privacy, machine learning, membership inference, property inference, model extraction, reconstruction, model inversion

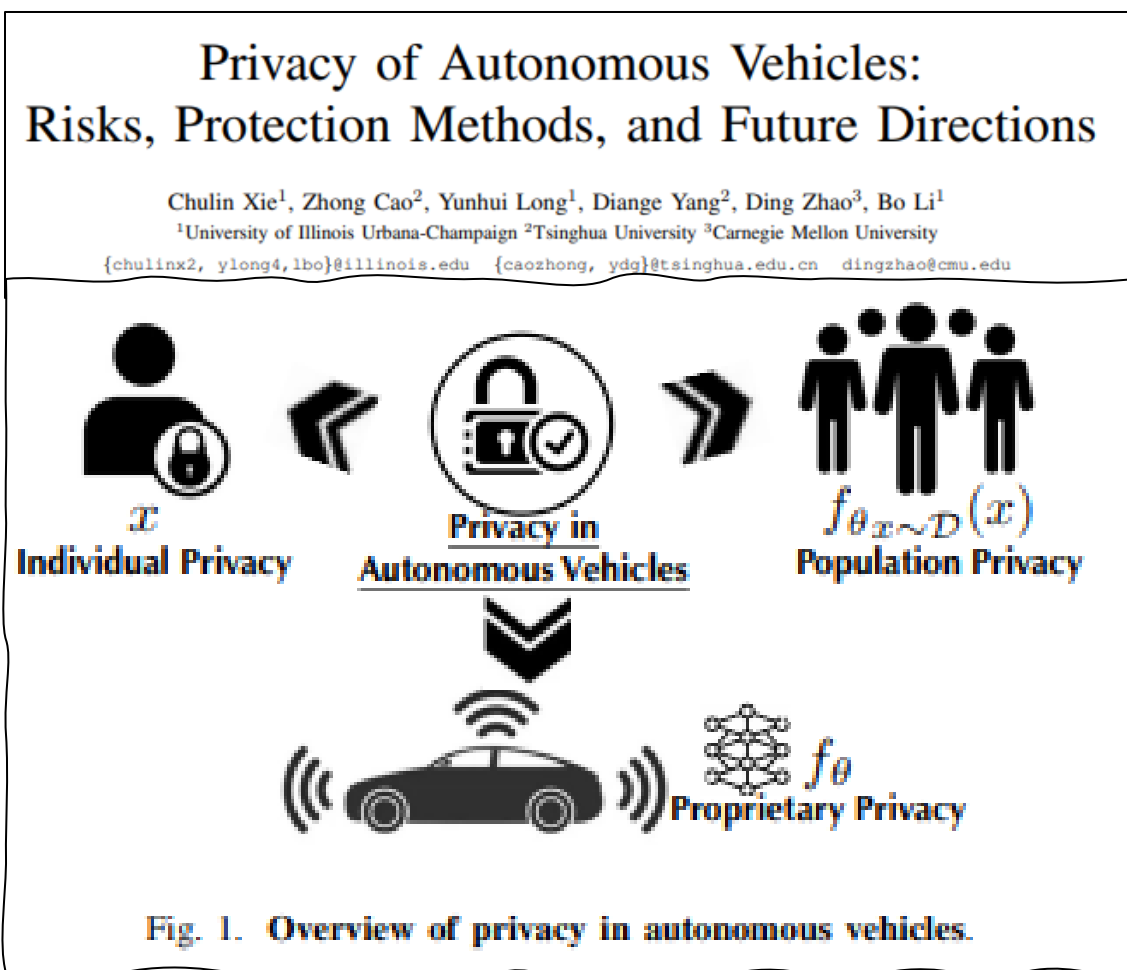
### ACM Reference format:

Maria Rigaki and Sebastian Garcia. 2023. A Survey of Privacy Attacks in Machine Learning. *ACM Comput. Surv.* 56, 4, Article 101 (November 2023), 34 pages.  
<https://doi.org/10.1145/3624010>



# Connected cars

## Eyes unseen, yet always near



**TABLE II  
INDIVIDUAL PRIVACY: RISKS AND PRIVACY ATTACKS**

Attack Category	Attack target	Adversary's Knowledge	Attack Method	Data Type
Membership Inference	Classification models	Black-box model, input distribution, model type	shadow model	location, image, tabular data [78]
		Black-box model	confidence-thresholding	location, image, text, tabular data [79]
		Black-box model	shadow dataset generation	image, tabular data [80]
		Black-box model, average training loss	threshold-based	tabular data [81]
		Black-box model, input distribution, model type	shadow model, robustness evaluation	images, location [82]
Generative models	Classification models	White/Black-box model	confidence-thresholding, GAN	images [79]
		Black-box model	latent encoding	images [80]
Aggregated statistics	Classification models	Black-box model, prior observation	game-based procedure	location [81]
Model Inversion	Classification models	White/Black-box model	optimization-based	image, tabular data [82]
		Black-box model, public data	inversion model	image [82]
		White-box model, public data, corrupted target input	GAN	image [83]
Model Memorization	Classification models	White/Black-box model	encoding	image, text [84]
	Generative models	White-box model	random sequences insertion	text [13]

Information technology — Artificial  
intelligence — AI system life cycle  
processes

*Technologies de l'information — Intelligence artificielle — Processus  
de cycle de vie des systèmes d'IA*



# Data is the heartbeat of everything

## *Artificial Intelligence Life Cycle*

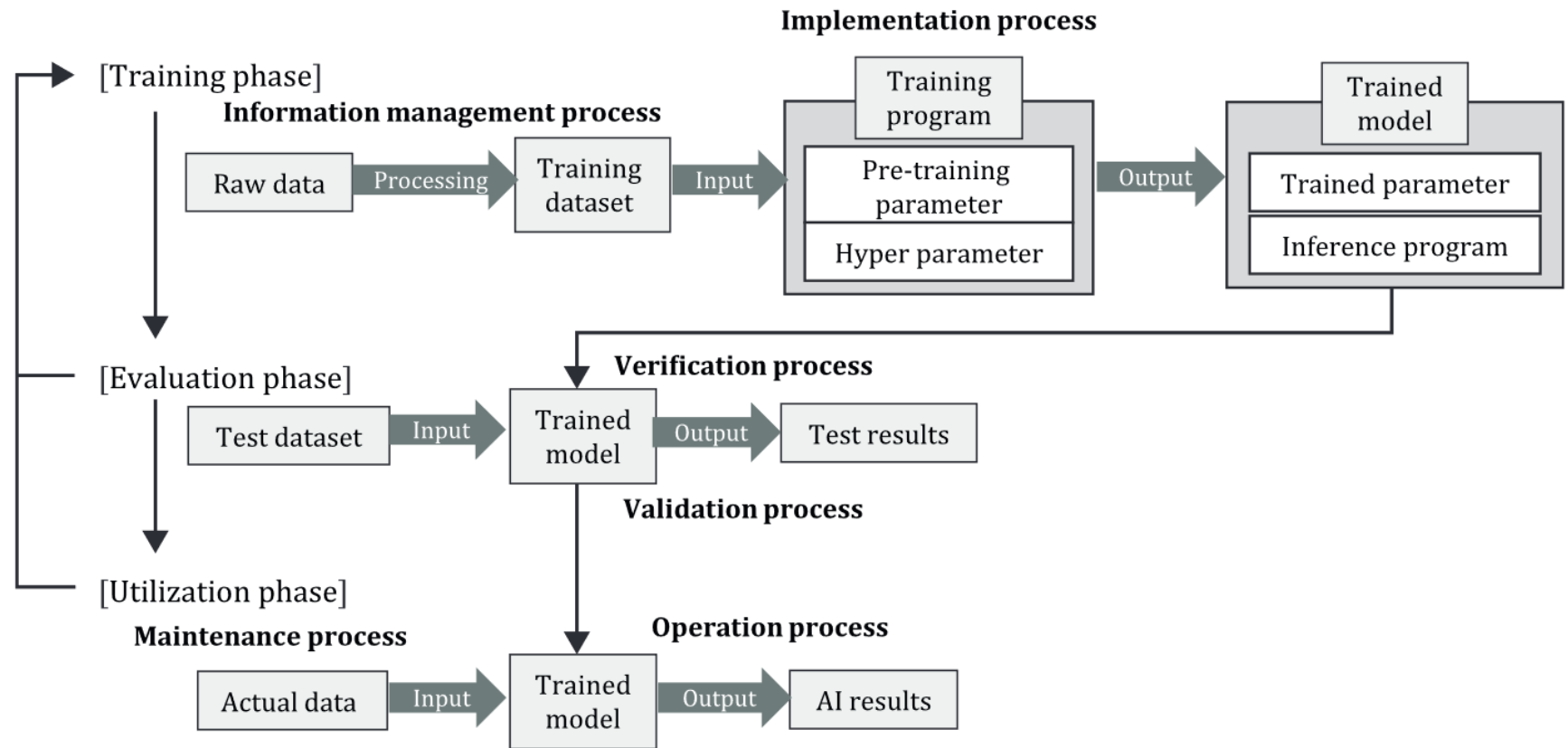


Figure A.1 — The flow of AI-specific processes

Information technology — Artificial  
intelligence — AI system life cycle  
processes

*Technologies de l'information — Intelligence artificielle — Processus  
de cycle de vie des systèmes d'IA*



# Data is the heartbeat of everything

## *Artificial Intelligence Life Cycle*

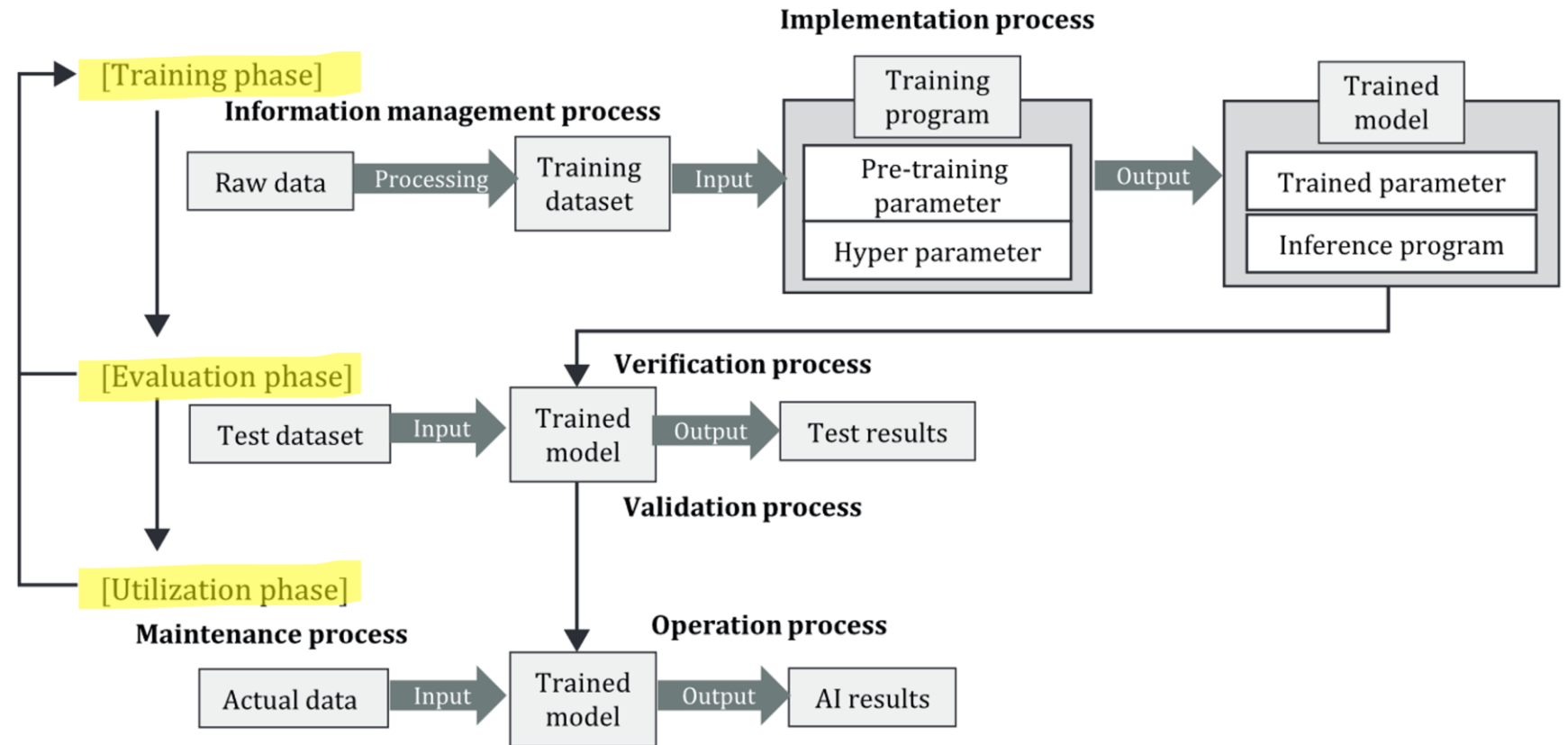
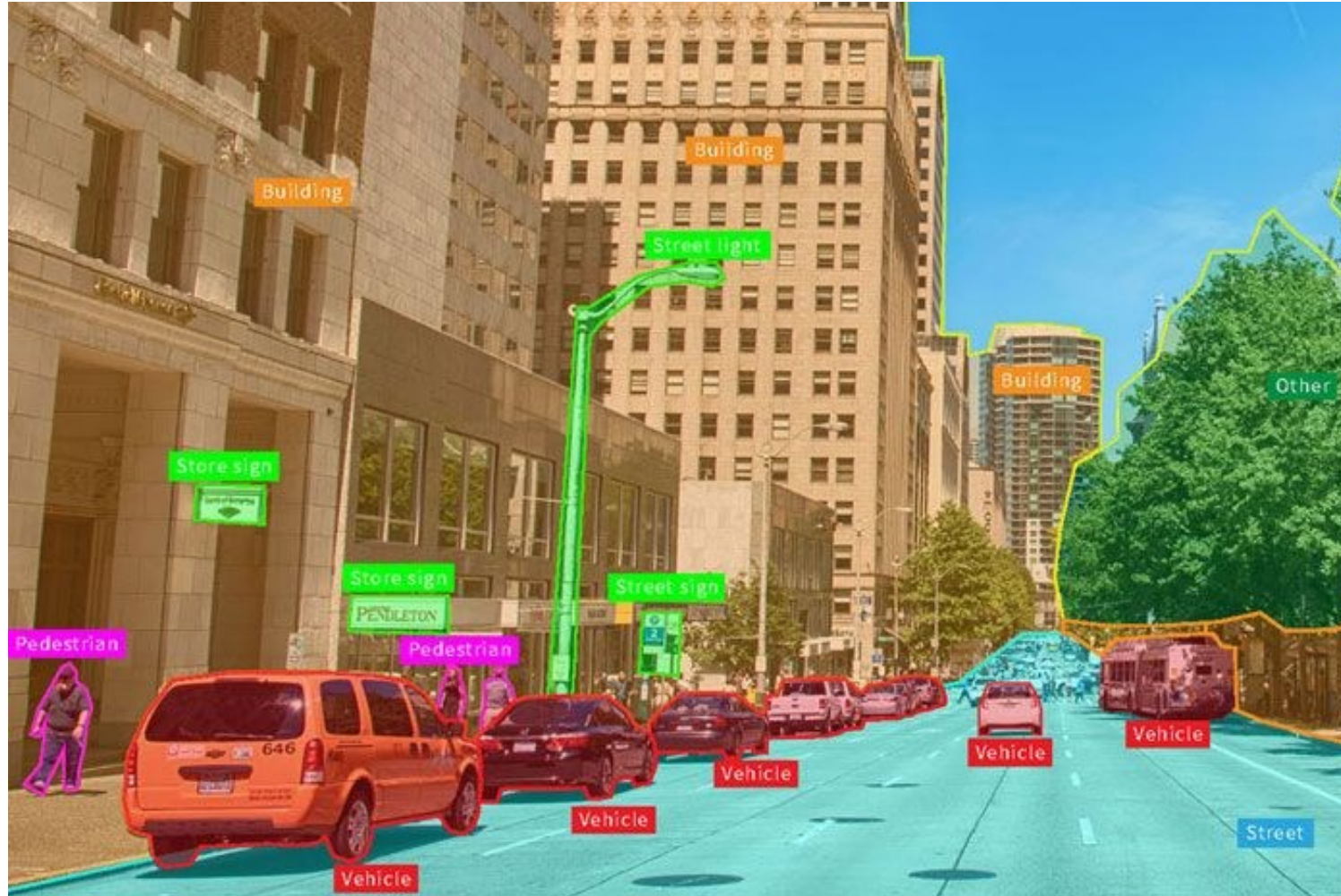


Figure A.1 — The flow of AI-specific processes

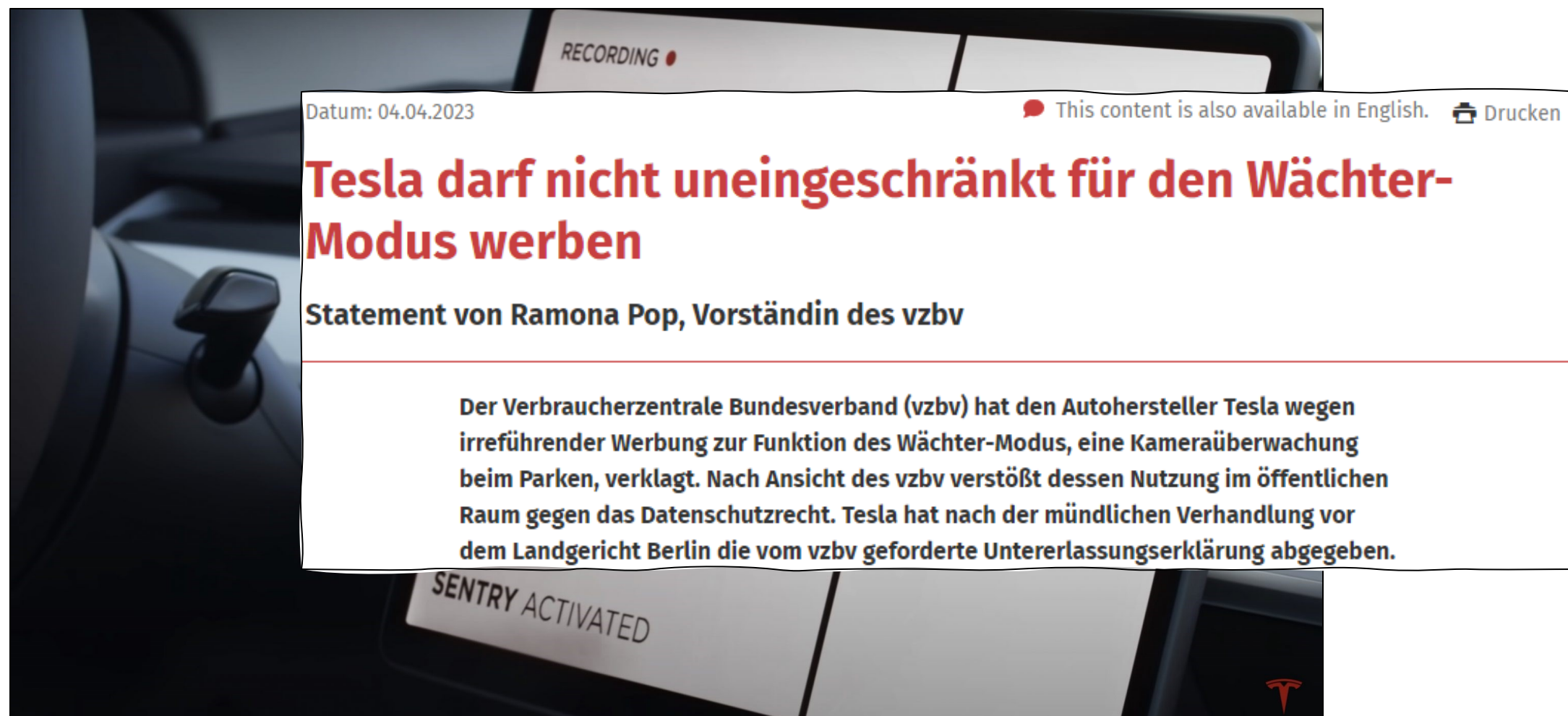
# From life's chaos, data brings order

## Through bounding boxes & segmentation




# Data moves at the speed of light

The law struggles to catch it



The image shows a screenshot of a news article overlaid on a background of a car's interior. The background features a digital display with the word "RECORDING" and a red dot, and another display below it with the text "SENTRY ACTIVATED" and the Tesla logo. The news article text is contained within a white box with a black border.

Datum: 04.04.2023 This content is also available in English.  Drucken

## Tesla darf nicht uneingeschränkt für den Wächter-Modus werben

Statement von Ramona Pop, Vorsitzin des vzbv

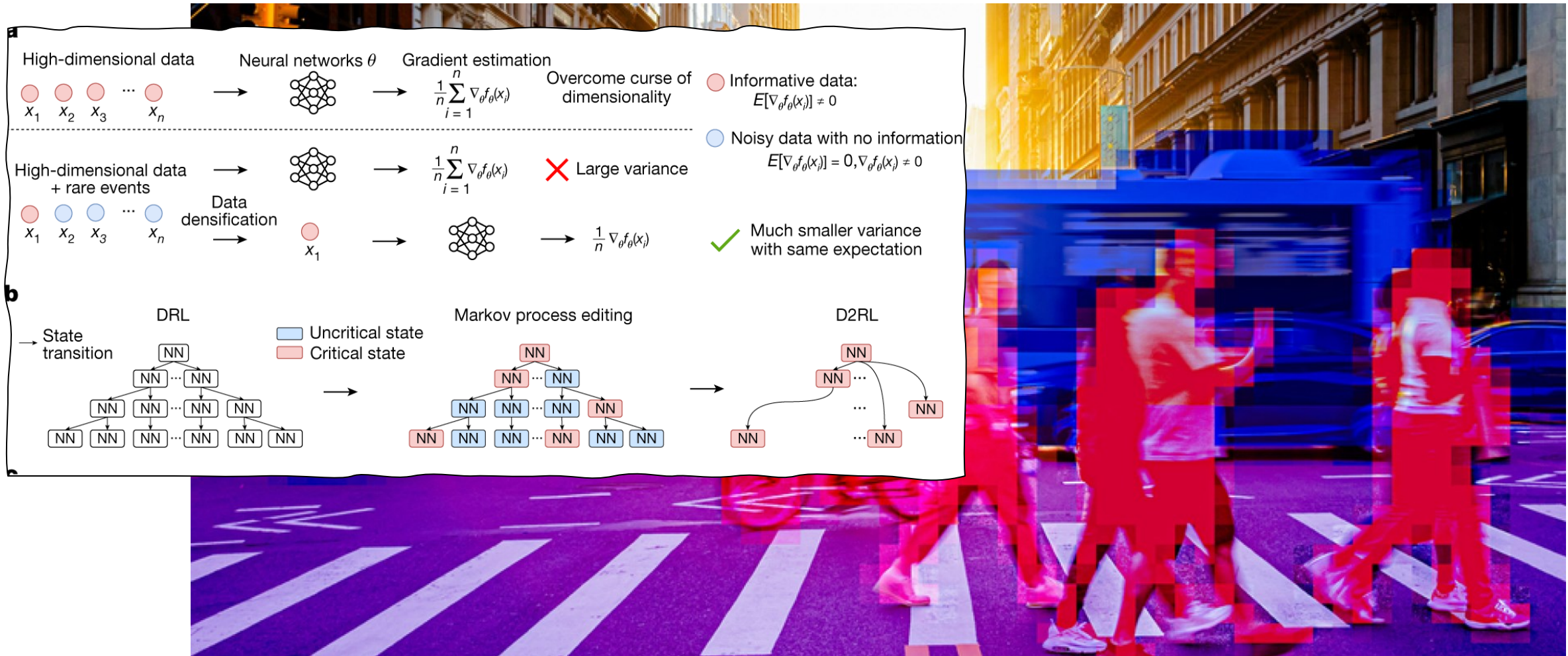
---

Der Verbraucherzentrale Bundesverband (vzbv) hat den Autohersteller Tesla wegen irreführender Werbung zur Funktion des Wächter-Modus, eine Kameraüberwachung beim Parken, verklagt. Nach Ansicht des vzbv verstößt dessen Nutzung im öffentlichen Raum gegen das Datenschutzrecht. Tesla hat nach der mündlichen Verhandlung vor dem Landgericht Berlin die vom vzbv geforderte Untererlassungserklärung abgegeben.



# Neural networks feast on data

## To reinforce stability and safety



# Autonomous driving's fairness

## Depends on inclusive data

### Bias Behind the Wheel: Fairness Analysis of Autonomous Driving Systems

XINYUE LI, Peking University, China

ZHENPENG CHEN, Nanyang Technological University, Singapore

JIE M. ZHANG, King's College London, United Kingdom

FEDERICO  
YIN  
XUAN

This p  
tonom  
demog  
extens  
and 3,  
undet  
variou  
strate  
female  
both e  
fairne  
code, d

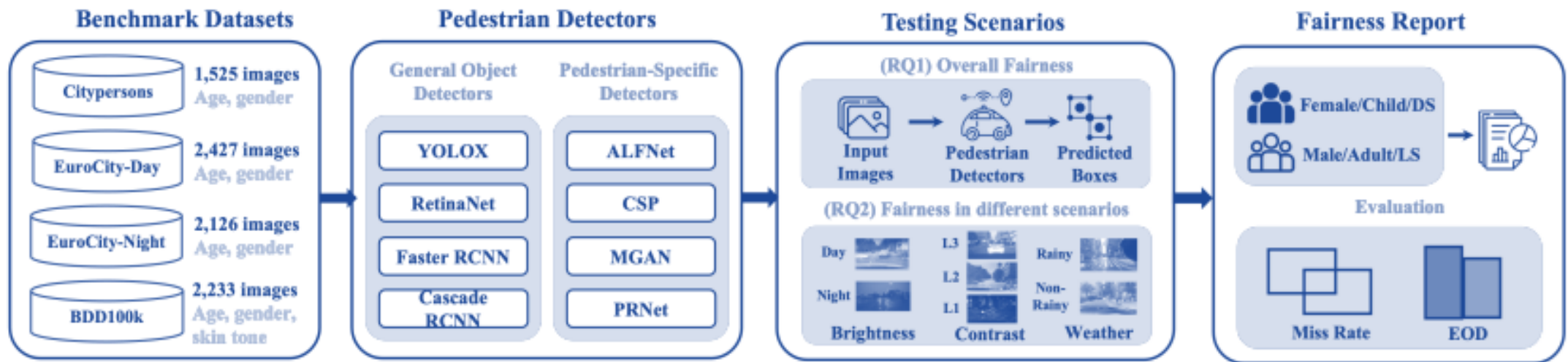


Fig. 1. Overview of our experimental settings.

# Diverse data saves lives

## Without it, systems fail

Table 4. Number of labeled pedestrian instances per dataset.

Dataset	Gender		Age		Skin tone	
	Male	Female	Adult	Child	Light	Dark
CityPersons	2,357	1,822	4,568	233	-	-
EuroCity-Day	1,726	1,646	4,498	100	-	-
EuroCity-Night	1,265	1,318	4,165	68	-	-
BDD100k	3,457	2,479	6,293	190	2,724	789
<b>Overall</b>	<b>8,805</b>	<b>7,265</b>	<b>19,524</b>	<b>591</b>	<b>2,724</b>	<b>789</b>

Detectors	Age					
	Day-time			Night-time		
	MR Adult	MR Child	EOD (Age)	MR Adult	MR Child	EOD (Age)
YOLOX	10.80%	40.77%	-29.97%	17.81%	54.93%	-37.12%
RetinaNet	12.73%	41.92%	-29.19%	18.97%	61.97%	-43.01%
Faster RCNN	4.50%	25.00%	-20.50%	7.34%	33.80%	-26.46%
Cascade RCNN	4.52%	25.58%	-21.05%	6.81%	33.80%	-26.99%
ALFNet	26.52%	49.42%	-22.90%	65.10%	83.10%	-18.00%
CSP	28.64%	46.92%	-18.28%	63.42%	76.06%	-12.64%
MGAN	27.87%	43.46%	-15.59%	49.72%	69.01%	-19.30%
PRNet	34.80%	56.92%	-22.13%	72.60%	92.96%	-20.36%
<b>Average</b>	<b>18.80%</b>	<b>41.25%</b>	<b>-22.45%</b>	<b>37.72%</b>	<b>63.20%</b>	<b>-25.48%</b>

# Diverse data saves lives

## Without it, systems fail

Table 4. Number of labeled pedestrian instances per dataset.

Dataset	Gender		Age		Skin tone
	Male	Female	Adult	Child	
CityPersons	2,357	1,822			
EuroCity-Day	1,726	1,646			
EuroCity-Night	1,265	1,318			
BDD100k	3,457	2,479			
<b>Overall</b>	<b>8,805</b>	<b>7,269</b>			

Age	
Adult	Child
4,568	233
4,498	100
4,165	68
6,293	190
<b>19,524</b>	<b>591</b>

Detectors	Day-time			EOD (Age)
	MR Adult	MR Child	EOD (Age)	
	YOLOX	10.80%	40.77%	
RetinaNet	12.73%	41.92%	-29.19%	-43.01%
MaskRCNN	4.50%	25.00%	-20.50%	-26.46%
MaskRCNN	4.52%	25.58%	-21.05%	-26.99%
ALFNet	26.52%	49.42%	-22.90%	-18.00%
CSP	28.64%	46.92%	-18.28%	-12.64%
MGAN	27.87%	43.46%	-15.59%	-19.30%
PRNet	34.80%	56.92%	-22.13%	-20.36%
<b>Average</b>	<b>18.80%</b>	<b>41.25%</b>	<b>-22.45%</b>	<b>-25.48%</b>

EOD (Age)
-29.97%
-29.19%
-20.50%
-21.05%
-22.90%
-18.28%
-15.59%
-22.13%
<b>-22.45%</b>

# The model's manifold grows

## With each layer of data diversity



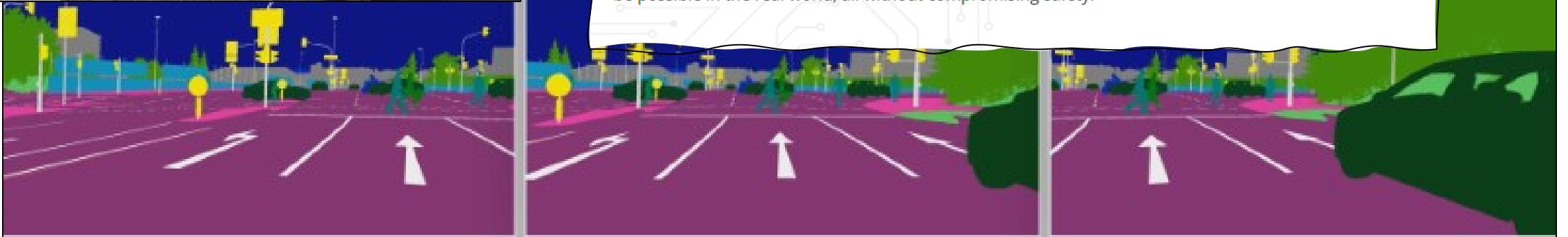
## SYNTHETIC TRAINING DATA

SYNTHETIC TRAINING DATA

SOFTWARE IN-THE-LOOP

// QUICKLY AND THOROUGHLY TRAIN YOUR ARTIFICIAL INTELLIGENCE (AI)

rFpro enables the training, testing and validation of a vehicle's AI through the mass creation of highly accurate training data. It subjects the AI to more edge cases and life-threatening situations than would be possible in the real world, all without compromising safety.



# Comprehensive, diverse data

## The lifeblood of accurate medical AI

Article | Published: 18 May 2020

### A deep learning system for differential diagnosis of skin diseases

[Yuan Liu](#), [Ayush Jain](#), [Clara Eng](#), [David H. Way](#), [Kang Lee](#), [Peggy Bui](#), [Kimberly Kanada](#), [Oliveira Marinho](#), [Jessica Gallegos](#), [Sara Gabriele](#), [Vishakha Gupta](#), [Nalini Singh](#), [Vivek Hofmann-Wellenhof](#), [Greg S. Corrado](#), [Lily H. Peng](#), [Dale R. Webster](#), [Dennis Ai](#), [Susan R. Carter Dunn](#) & [David Coz](#)

*Nature Medicine* **26**, 900–908

23k Accesses | 377 Citations

Original Investigation | Innovations in Health Care Delivery

December 13, 2016

### Development and Validation of a Deep Learning Algorithm for Detection of Diabetic Retinopathy in Retinal Fundus Photographs


[Varun Gulshan](#), PhD<sup>1</sup>; [Lily Peng](#), MD, PhD<sup>1</sup>; [Marc Coram](#), PhD<sup>1</sup>; [et al](#)

[» Author Affiliations](#) | [Article Information](#)

*JAMA*. 2016;316(22):2402-2410. doi:10.1001/jama.2016.17216

Letter | Published: 20 May 2019

### End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography

[Diego Ardila](#), [Atilla P. Kiraly](#), [Sujeeth Bharadwaj](#), [Bokyung Choi](#), [Joshua J. Reicher](#), [Lily Peng](#), [Daniel Tse](#) , [Jo](#), [David P. Naidich](#) & [Shravya Shetty](#)

[this article](#)

[tmetric](#) | [Metrics](#)

# Protectionism

## Will not last forever



Der Hamburgische Beauftragte für  
Datenschutz und Informationsfreiheit

Diskussionspapier: Large Language Models  
und personenbezogene Daten

1. **Die bloße Speicherung eines LLMs stellt keine Verarbeitung** im Sinne des Art. 4 Nr. 2 DSGVO dar. **Denn in LLMs werden keine personenbezogenen Daten gespeichert.** Soweit in einem LLM-gestützten KI-System personenbezogene Daten verarbeitet werden, müssen die Verarbeitungsvorgänge den Anforderungen der DSGVO entsprechen. Dies gilt insbesondere für den Output eines solchen KI-Systems.
2. **Mangels Speicherung** personenbezogener Daten im LLM **können die Betroffenenrechte der DSGVO nicht das Modell selbst zum Gegenstand haben.** Ansprüche auf Auskunft, Löschung oder Berichtigung **können sich jedoch zumindest auf Input und Output eines KI-Systems der verantwortlichen Anbieter:in oder Betreiber:in beziehen.**
3. **Das Training von LLMs mit personenbezogenen Daten muss datenschutzkonform erfolgen.** Dabei sind auch die Betroffenenrechte zu beachten. Ein ggf. datenschutzwidriges Training wirkt sich aber nicht auf die Rechtmäßigkeit des Einsatzes eines solchen Modells in einem KI-System aus.

# *Data Altruism By Design*

## Carving out a GDPR exception for Machine Learning

1. Lawful access

2. The storing of the data is not the purpose

3. The systems aren't used to the detriment of the data subjects or in a way that

4. Reversed burden of proof

5. Point 4 includes implementing appropriate TOMs throughout the AI lifecycle

An insurance fund to cover risks.



# Kontakt

**Tea Mustać**

[tea.mustac@spiritlegal.com](mailto:tea.mustac@spiritlegal.com)

[www.spiritlegal.com](http://www.spiritlegal.com)

