

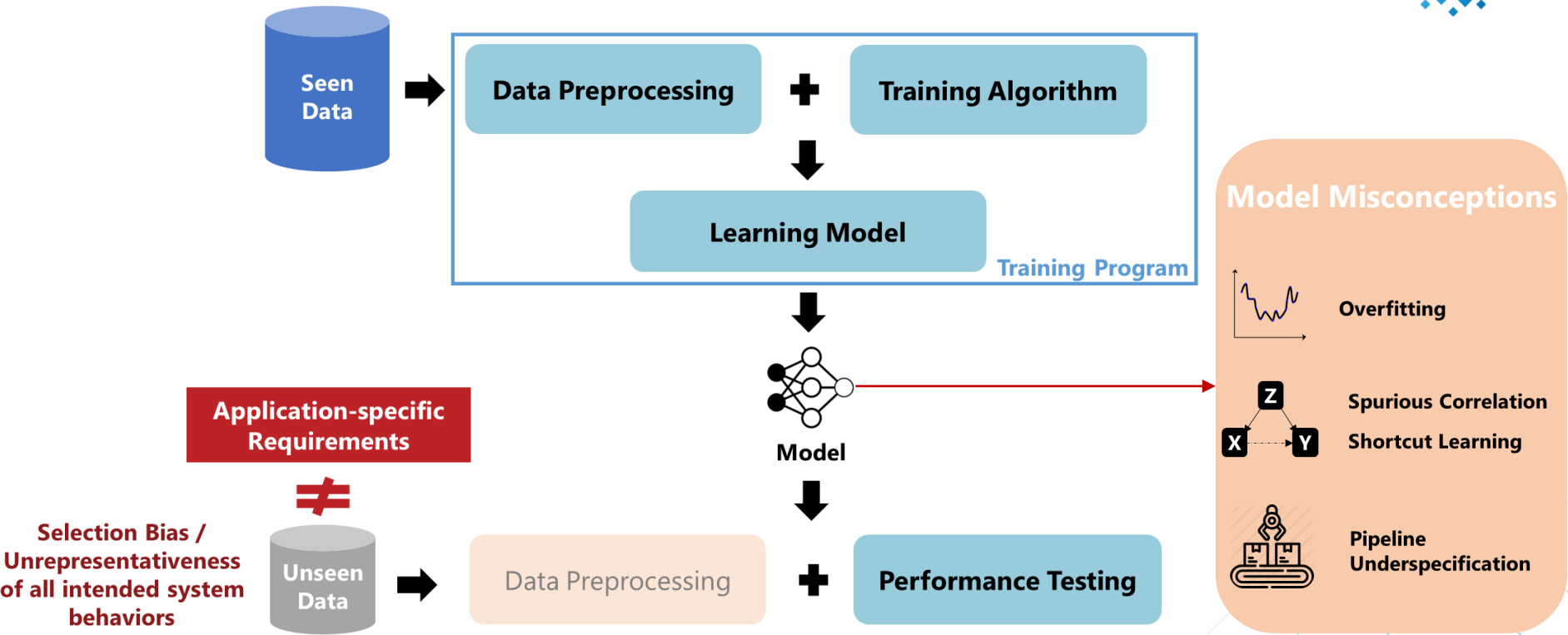


Domain-Aware DL Model Testing

Houssem Ben Braiek, Ph. D.



Underspecification Issues of Unseen Datasets



Selection Bias /
Unrepresentativeness
of all intended system
behaviors

Application-specific
Requirements

Unseen
Data

Data Preprocessing + Training Algorithm
↓
Learning Model
Training Program

Model

Data Preprocessing + Performance Testing

Model Misconceptions

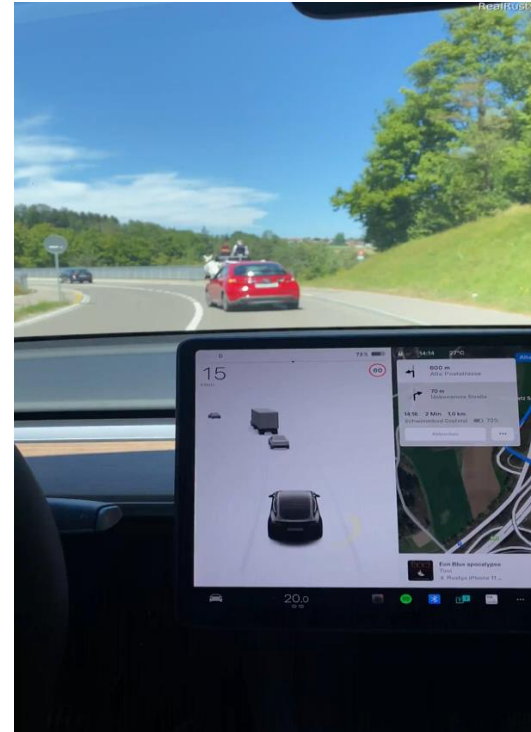
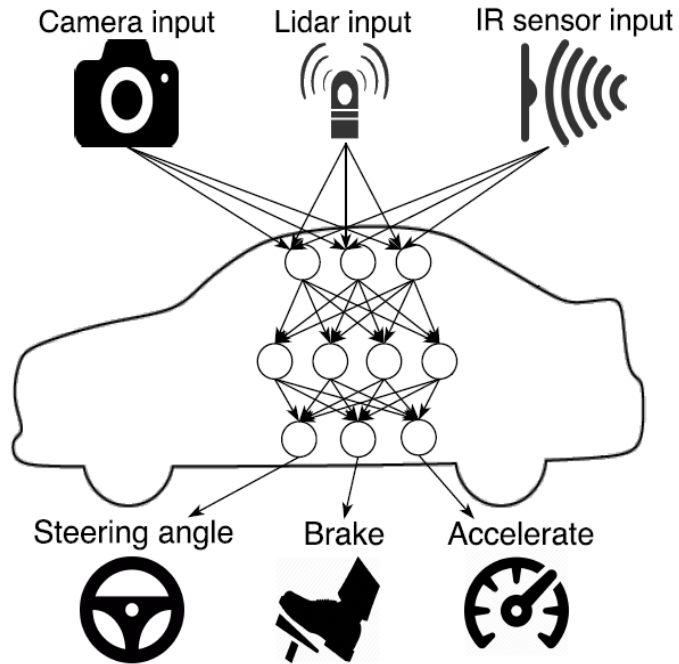
- Overfitting
- Spurious Correlation
Shortcut Learning
- Pipeline Underspecification

Why do DL practitioners perceive the value of DL testing differently?



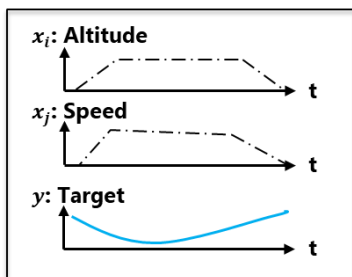
	Low Risk	High Risk
Quantifiable Performance	Outperform the state-of-the-art on testing benchmarks , e.g., ImageNet, Coco, etc.	Maintain an acceptable performance for a critical function under carefully controlled conditions , e.g., a custom-made cobot that performs repetitive tasks in a manufacturing facility.
Non-Quantifiable Performance	Provide added value over legacy baselines or fill a gap , e.g., filtering ads, recommending movies, etc.	Guarantee an acceptable performance for a critical function under all foreseeable operational conditions , e.g., a generic-purpose cobot that assists the elderly with household duties.

High Risk, Non-Quantifiable Performance ...





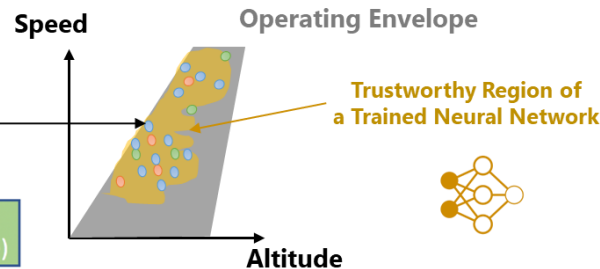
The Case of Aircraft System Performance Models



Timeseries Data Flights

(i) Extraction of **steady-state flight** data points

(ii) Preprocessing & Splitting data points into :

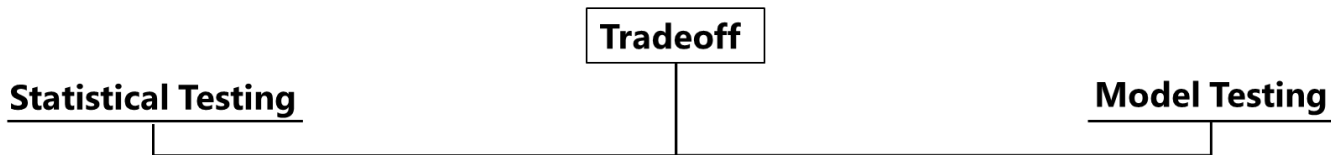


→ A trained NN could illustrate the system performance over the range of included-or-close operational conditions.

'the equipment, systems, and installations must be designed and installed to ensure they perform their intended functions under all foreseeable operating conditions.' U.S Code of Federal Regulations, parts 23, 25, 27, 29

A trustworthy performance model must be qualified to be representative of system behavior throughout the range of foreseeable operational conditions.

The Need for Domain-Aware DL Testing Methods



Estimate the **iid performance** of the model for completely **new inputs**.

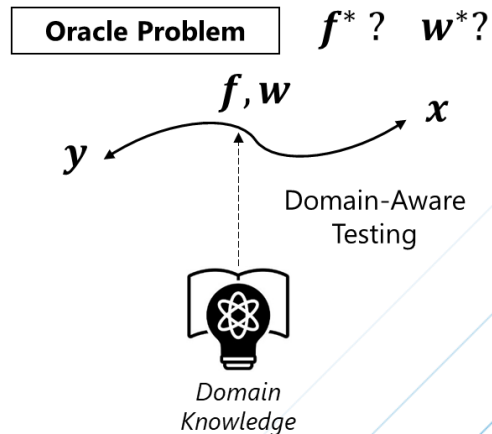
$$Err = \sum_{i \in D_{test}} (\hat{y}^{(i)} - y^{(i)})^2$$

Use unseen test data D_{test} as a proxy for future entries (x_{new}).

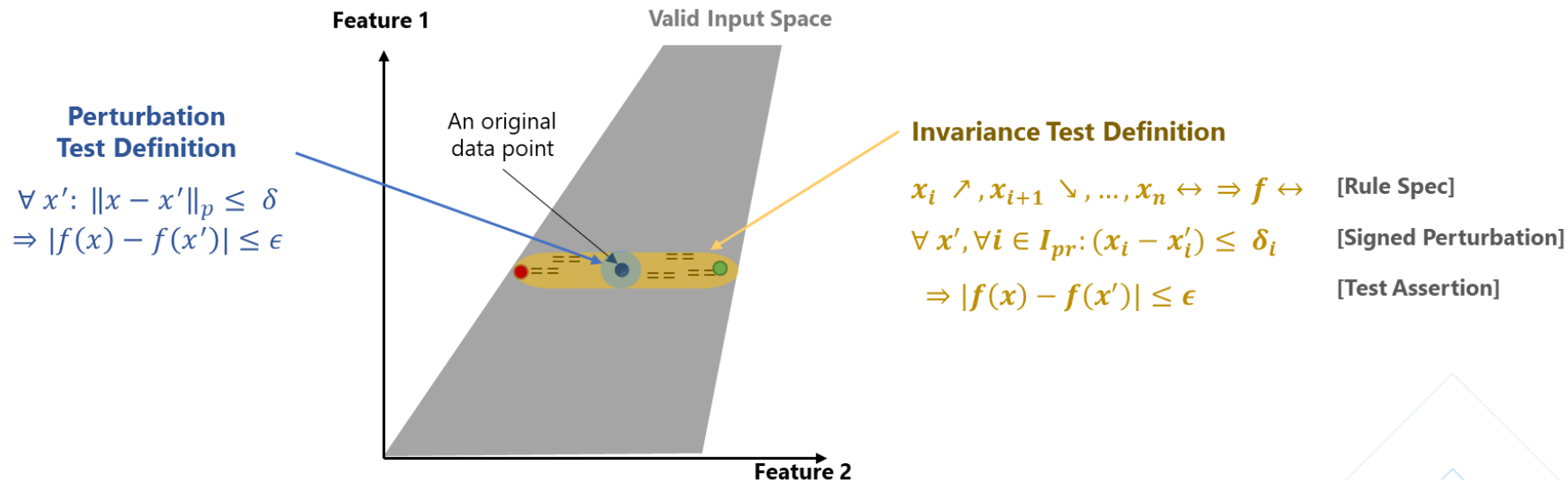
$$D_{test} = \{(x^{(i)}, y^{(i)})\}_{i \in [1, N]}$$

Collection of D_{test} is costly in aircraft industry

Test the **internal logic/mappings** of the model against the prior knowledge on the nature of the relation between x and y .

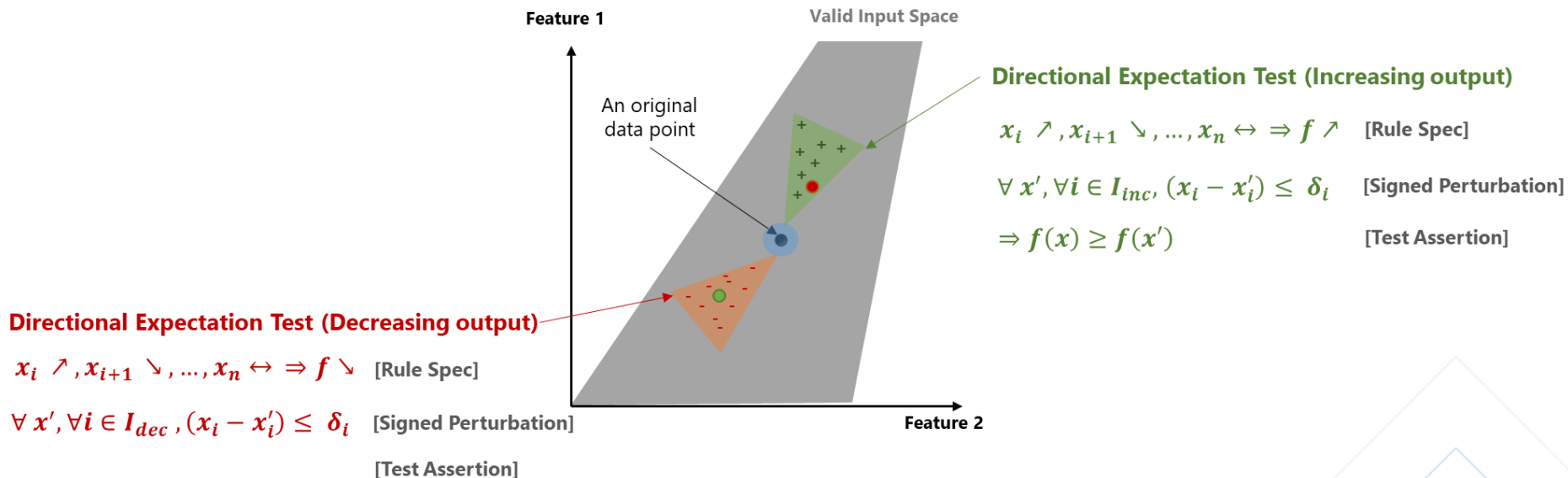


Invariance Tests



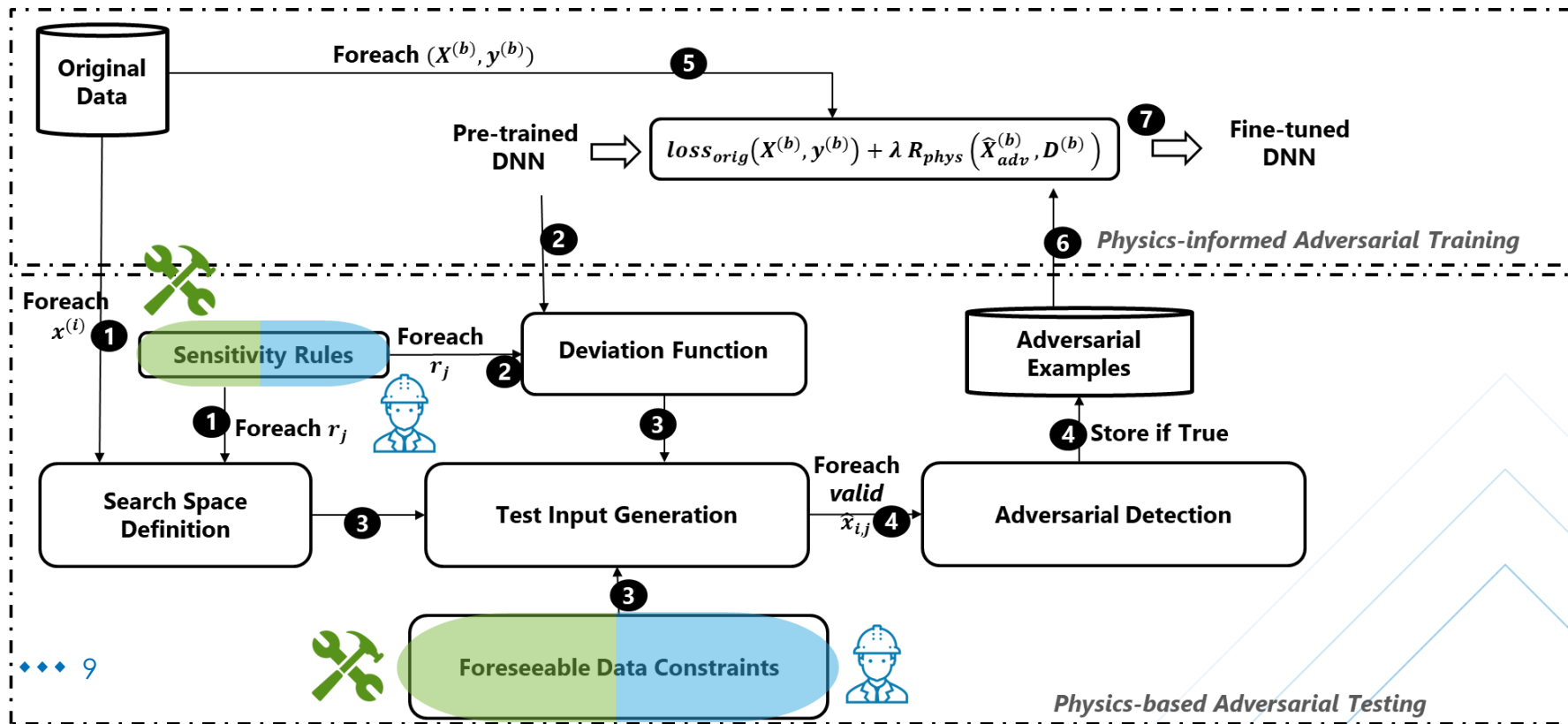
 These represent the failed inputs x for which the predictions are not consistent with the derived invariance tests.

Directional Expectation Tests



These represent the failed inputs x for which the predictions are not consistent with the derived directional expectation tests.

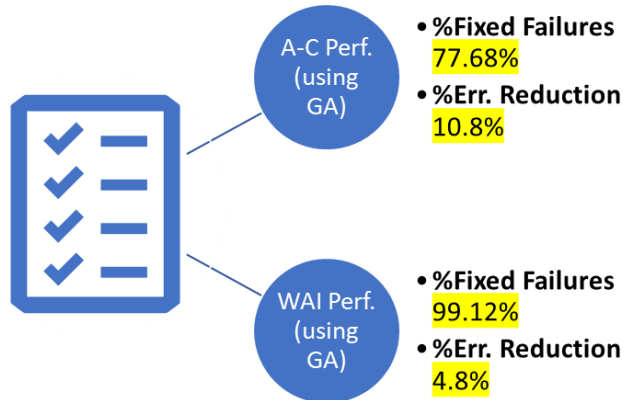
End-to-End Workflow of the Proposed Method



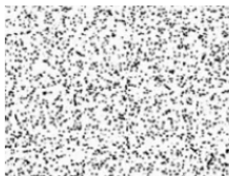
Evaluation Models & Results



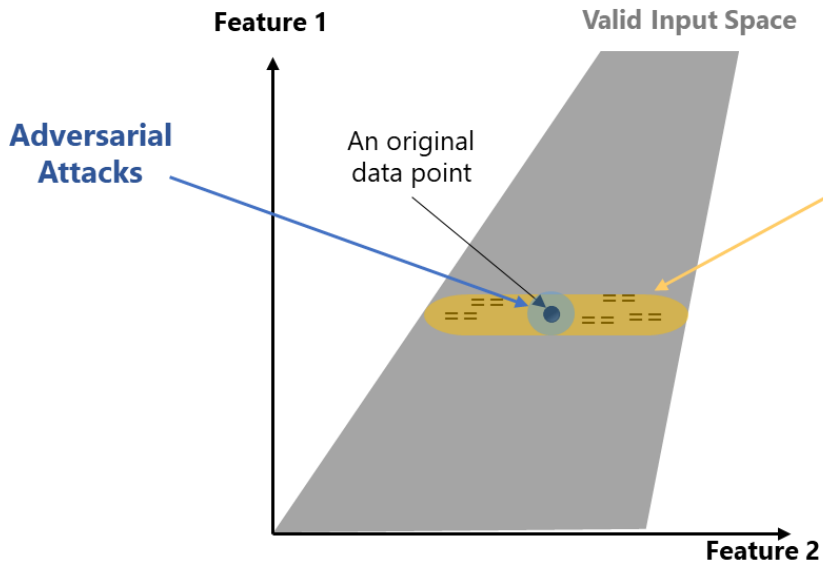
Model	Predicted Target	Description
Aircraft(A-C) Performance Model	α : angle of attack	The model maps steady-state angle of attack (α) to features related to flight conditions and wing configurations.
Wing Anti-Icing (WAI) Performance Model	T_{skin}^b : A-wing leading-edge skin temperature	The model maps the states of skin temperature sensors to features related to flight conditions, wing configurations, and high-pressure pneumatic system conditions at the wing root.
	T_{skin}^b : B-wing leading-edge skin temperature	



Analogies with other DL Applications



Noise
Perturbation



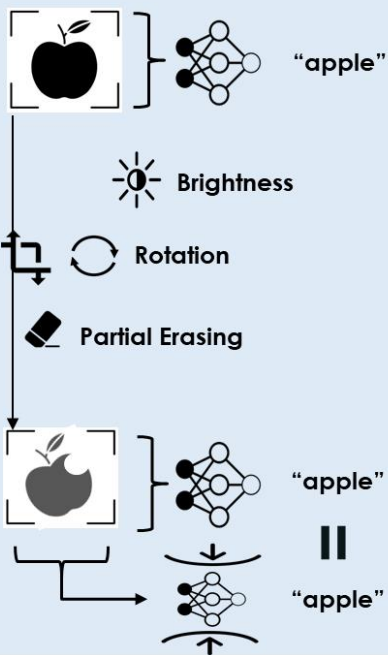
Semantically-preserving
Data Transformations



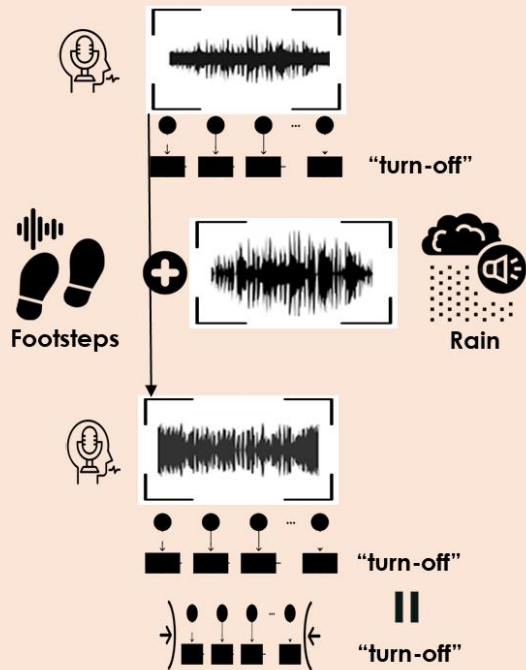
Semantically-preserving Data Transformations



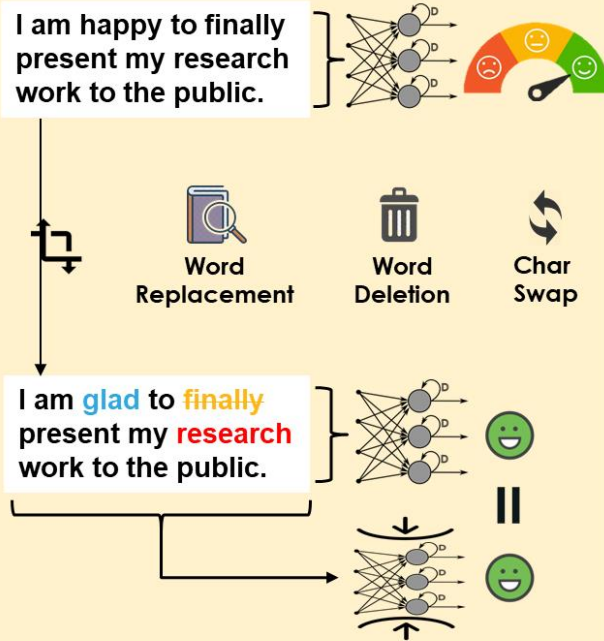
Object Classification



Speech Recognition



A= NLP Sentiment Detection



Semantically-preserving Data Transformations



For Images:

Pixel-value Transformations:



Brightness



Contrast



Blurring



Partial Erasing



Pixel Perturbation

Affine Transformations:



Translation



Shearing



Scaling



Rotation



For Audio Speeches:

Signal-wise Conversions:



Speed



Pitch



Loudness

Additive Noise Signals:



Random noisy perturbations



Colored noises: white, pink, brown.



Indoor Noises: breathing, footsteps, laughing, clock-tick, etc.



Outdoor Noises: Engine, Fireworks, Rain, Train, etc.



For Natural Language Texts :

Char-level Transformations:



Random Insertion



Random Swap



Random Deletion

Word-level Transformations:



Synonym/Embedding Replacement



Random Insertion

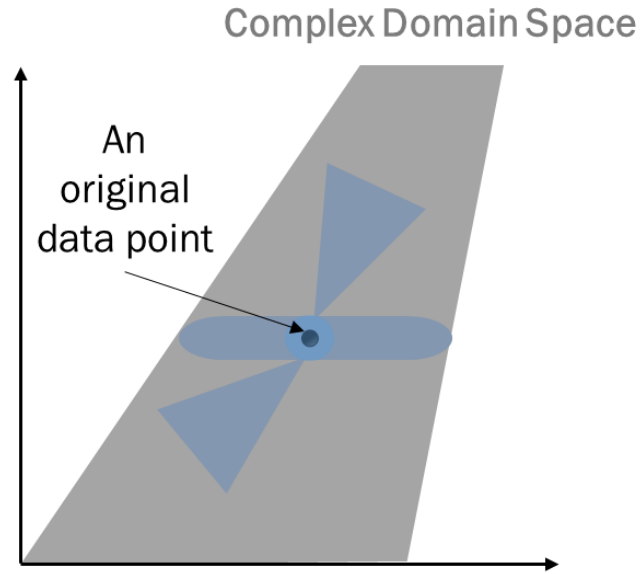


Random Swap



Random Deletion

How can we generate valid inputs from complex domains?

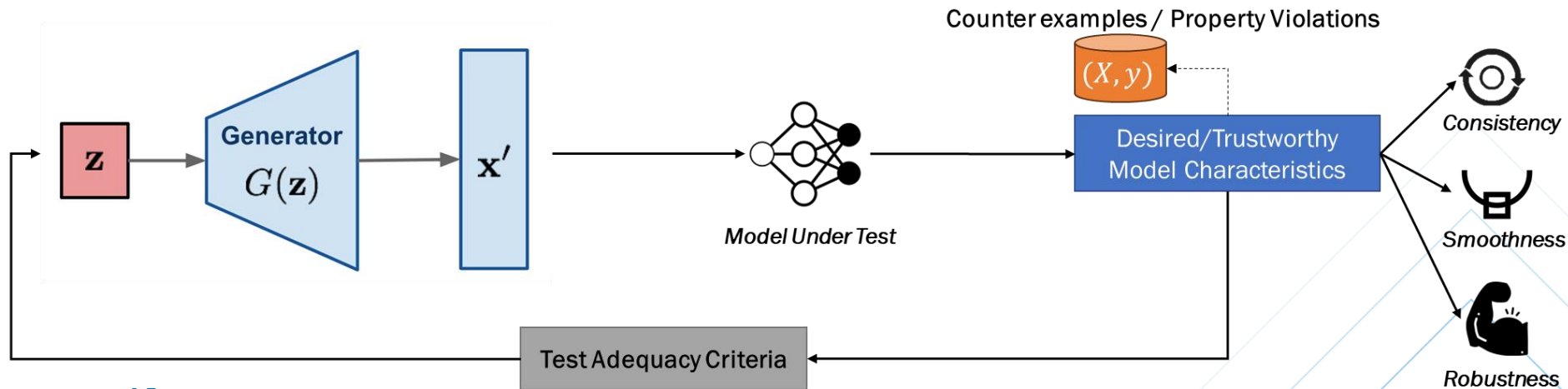


As software tests are written in code, DL tests can be produced by DL models !



DeepRoad [1] use GANs:

- map image from source domain to latent domain.
- generate image in the new domain from latent domain.



◆◆◆ 15

DEEL

DEpendable & Explainable Learning

