# STEFAN BOSSE[1,2]*

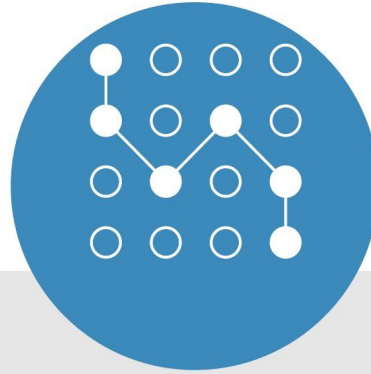[1] Institute of Computer Science
Researchgroup Practical Computer Science
University of Koblenz

[2] Department of Mechanical Engineering
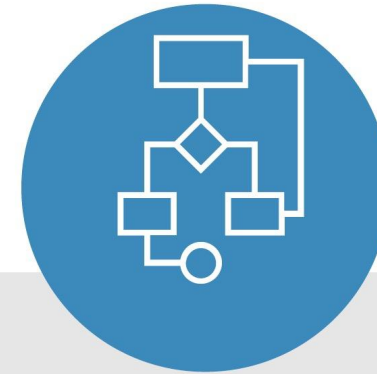Lehrstuhl für Materialkunde und Werkstoffprüfung
University of Siegen

**Explainable Data**

What data was used to train the model and why?

**Explainable Predictions**

What features and weights were used for this particular prediction?

**Explainable Algorithms**

What are the individual layers and the thresholds used for a prediction?

# WHY TRACEABILITY AND EXPLAINABILITY OF AI/ML MODELS IN MEASUREMENT AND TESTING TECHNOLOGY ARE MORE IMPORTANT THAN ACCURACY AND PRECISION

Stefan Bosse

LMW
Lehrstuhl für Materialkunde und Werkstoffprüfung

university of koblenz
Computer Science

# CONTENT

university
of koblenz
Computer Science

# BASICS

- **Data-driven Methods for**
  - Structural Health Monitoring
  - Damage Diagnostics
  - Material Testing
- **Data-driven Methods in**
  - Measurement - Sensors
  - Signal Processing
  - Aggregation
  - Transformation
  - Fusion
  - **Interpretation (Prediction)**
  - **Generation**
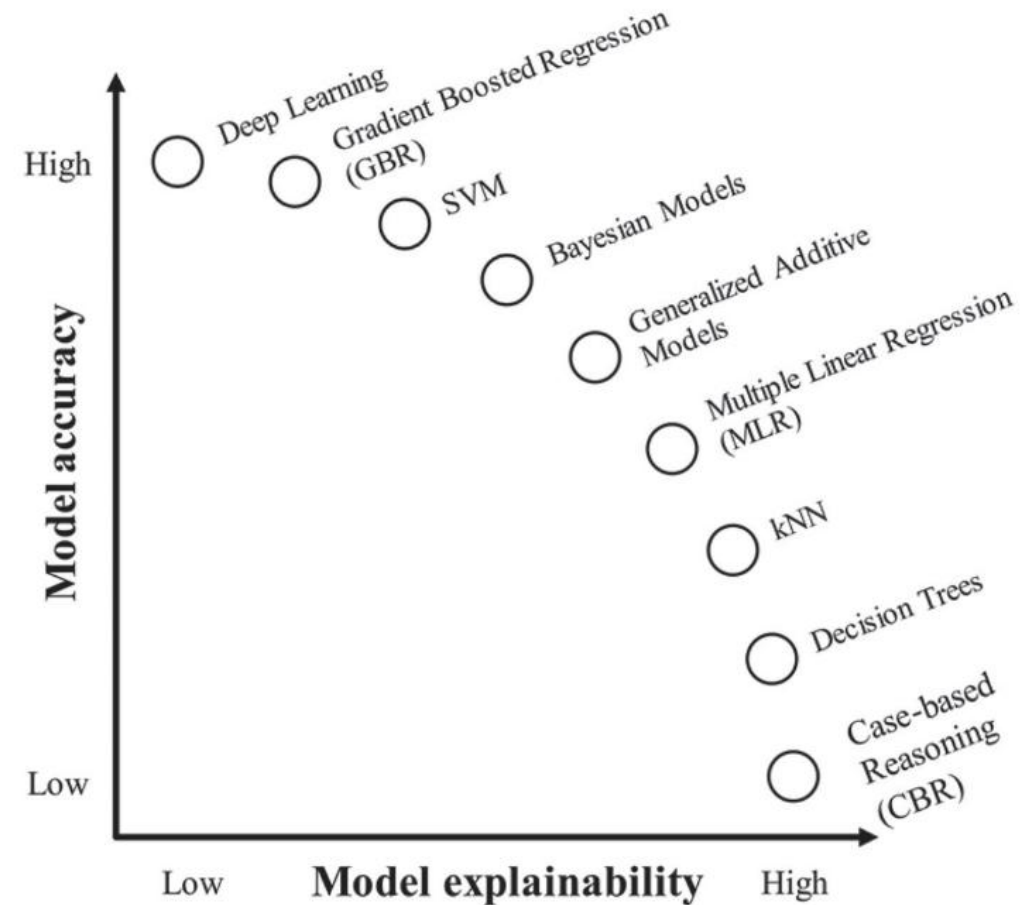
# BASICS

**Accuracy versa Precision**

- Strong. Accuracy: Close to expectation value (ground truth)

- Weak. Precision: Low variance

**Explainability versa Traceability or Tracktability**

- Strong. Explainability: An inductive model relationship between x and y based on knowledge

- Weak. Traceability: Which input contrbutes to output?

**Generalization: Specific or more general?**

- Interpolation versa Extrapolation



**ⓘ We have an opposite relationship between model accuracy and model explainability! And what's about model generalization?**

# BASICS
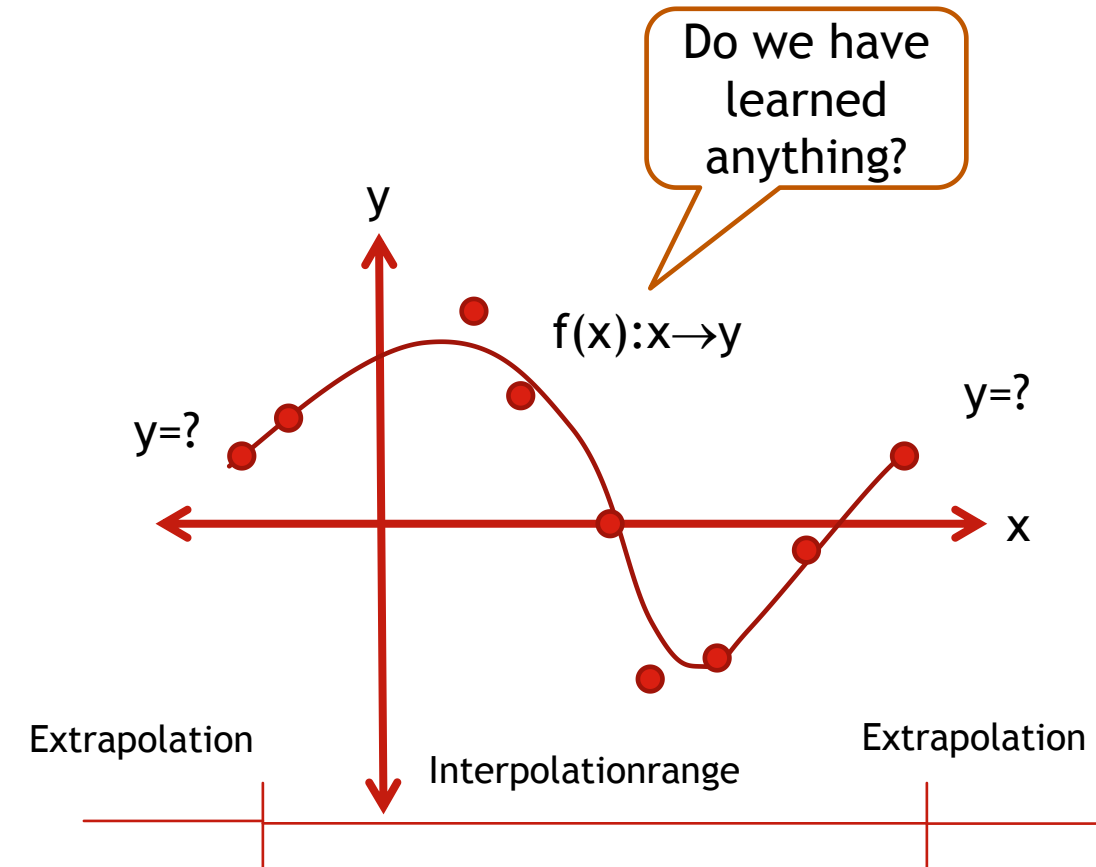
**Interpolation versa Extrapolation: Prediction Errors?**

Models

- Trees: I/E NA (partly by regression trees)

- Functions and Functional Graphs (SVM, ANN, CNN, ..)

- Training = Definition from Data ↔ Def. Parameter Space

Interpolation

- Input values between training points but inside parameter space

Extrapolation

- Input (and output) values outside the parameter space



Do we have learned anything?

f(x):x→y

y=?

y=?
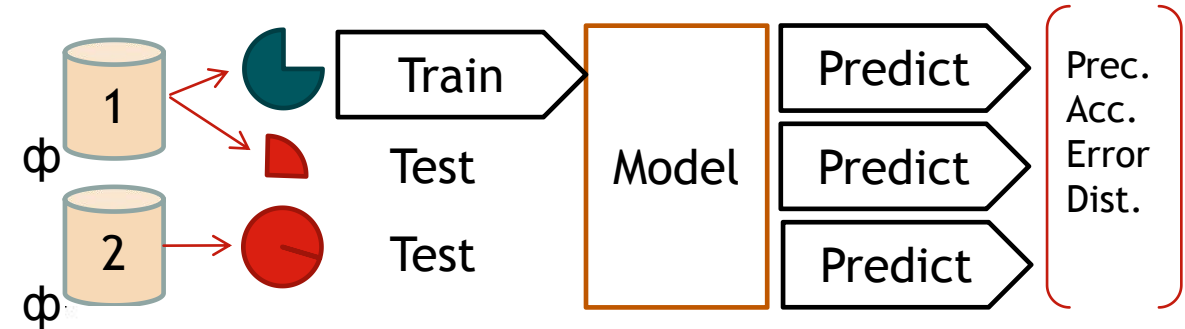
Extrapolation

Interpolationrange
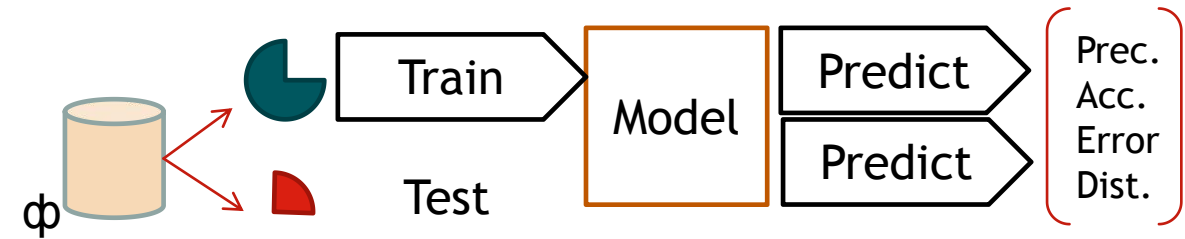
Extrapolation

ℹ️ **Interpolation should be possible with all functional models, extrapolation fails commonly! Generalization provides interpolation (degree 1) and extrapolation (degree 2). A specialized model will provide neither nor.**
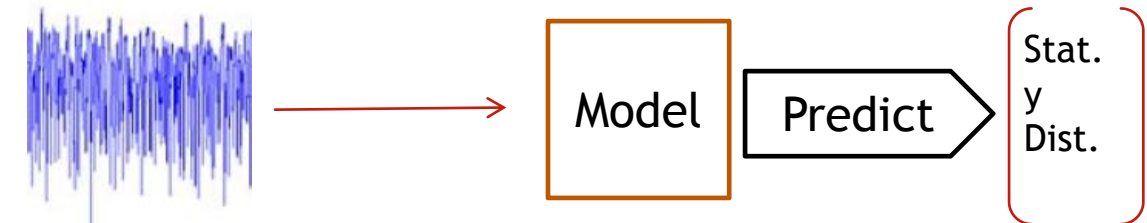
# BASICS

**Monitoring / Observation**

- Interpolation: Test with sample set not included in training (classical validation) but same parameters φ

- Extrapolation: Different sample distributions φ for test and training data

- Training Data versa Test Data: Accuracy, Precision, ...

**Robustness**

- What is the output if the input is ouside of the definition parameter space or non-sense data?

- Noise sensitivity?



φ → [ Train / Test ] → Model → Predict / Predict → Prec. Acc. Error Dist.

φ 1, φ 2 → [ Train / Test / Test ] → Model → Predict / Predict / Predict → Prec. Acc. Error Dist.

→ Model → Predict → Stat.y Dist.

**Training error low, test error high: Specialized Model ► Nice to have, but useless!**
**Training and test error low: Interpolating Model ► Degree 1 of Generalization**
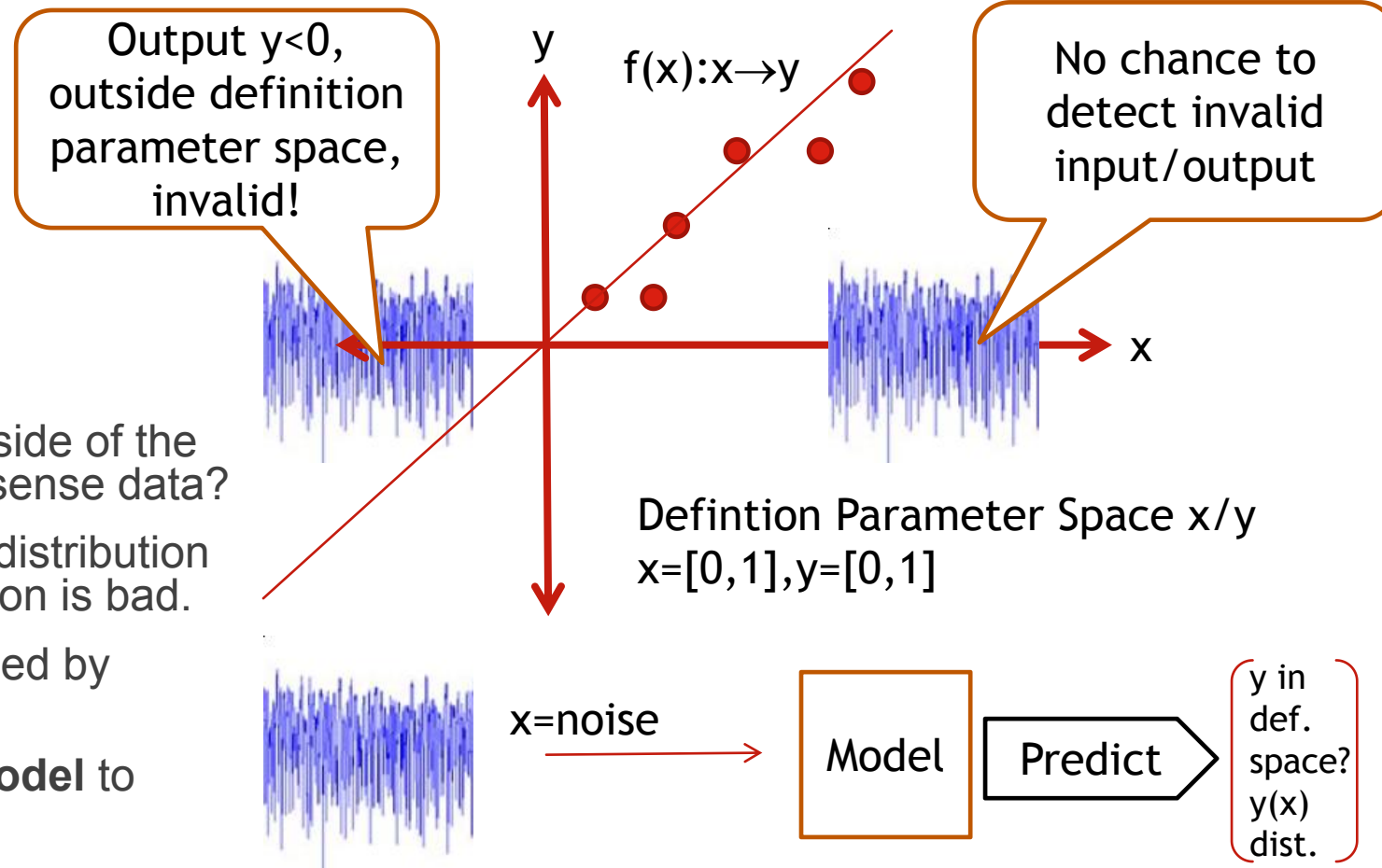**Different sample distributions covering different parameter spaces:**
**Extrapolation Test, low error outside training distribution ► Degree 2 of Generalization**

**university of koblenz**
Computer Science

# BASICS



**Robustness**

- What is the output if the input is ouside of the definition parameter space or non-sense data?

- Noise sensitivity? Make a y(noise) distribution test - high variant y(noise) distribution is bad.

- Can we detect **invalid output** caused by **invalid input**?

- Do we need a separate **scoring model** to classify input data validity? Or the „noise" class as additional output?

> **ℹ** With simple models we can maybe recognize and detect invalid input and invalid output. With any complex, highly non-linear, multi-variate and deep (nested) models this is mostly impossible!
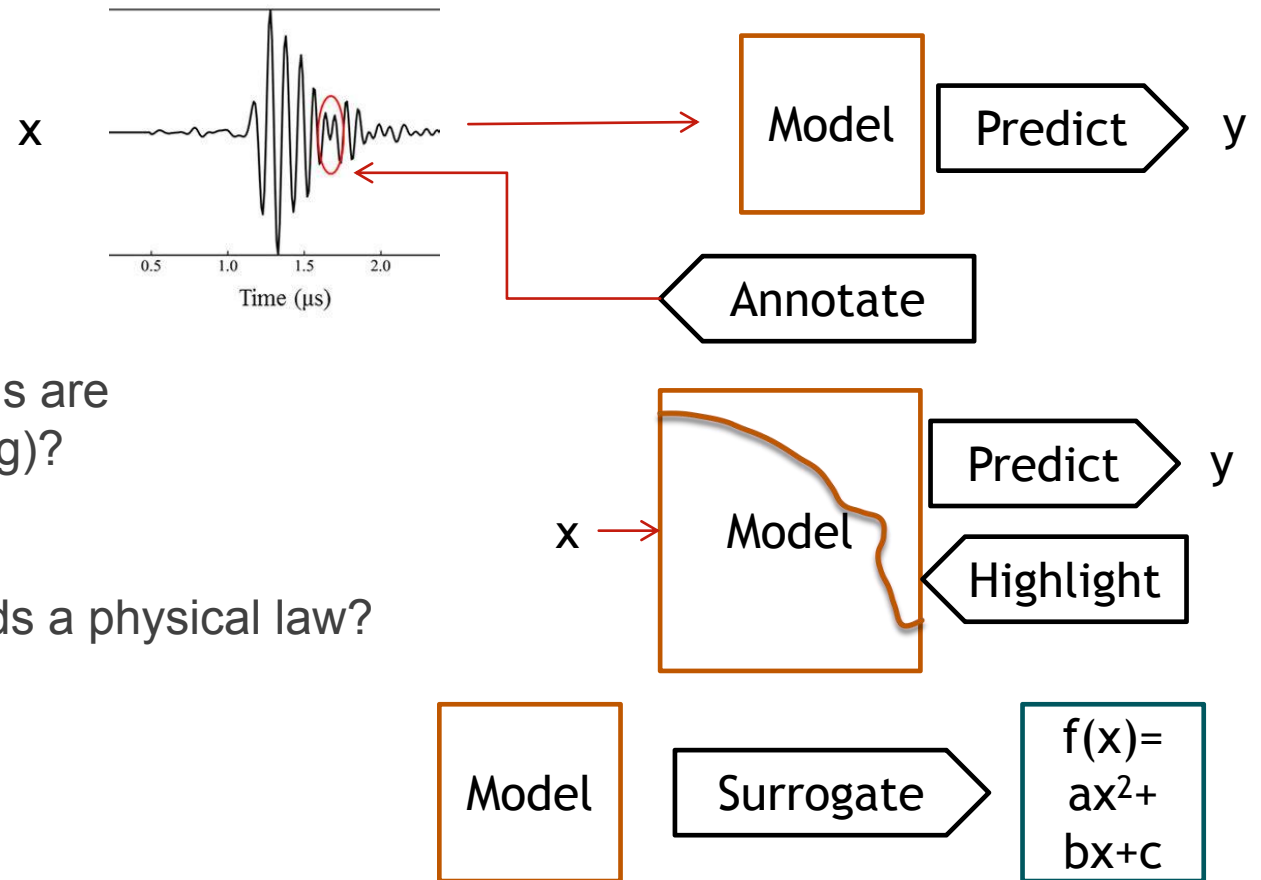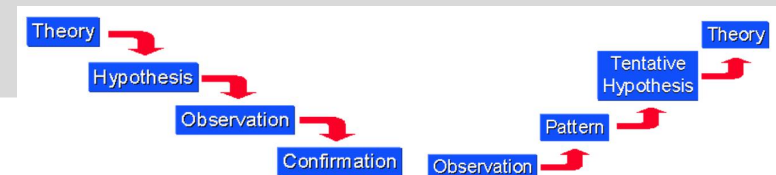
# BASICS

**Traceability**

- Which part of input is relevant?

- Which functional terms or paths in graphs are contributing to output (or decision making)?

**Explainibility**

- Anything learned? Generalization towards a physical law?

- Analytical explanations $x \rightarrow y$?

- Surrogate Modeling

- Explainable Models (e.g. XANN)

x

Time (μs)

Model | Predict | y

Annotate

x → Model | Predict | y

Highlight

Model | Surrogate | f(x)= ax$^2$+ bx+c

$$f(x)= ax^2 + bx+c$$

ⓘ **The most important methodology: Explain the model behavior or extraxt knowledge from a data-driven model. Induction is better than Deduction.**
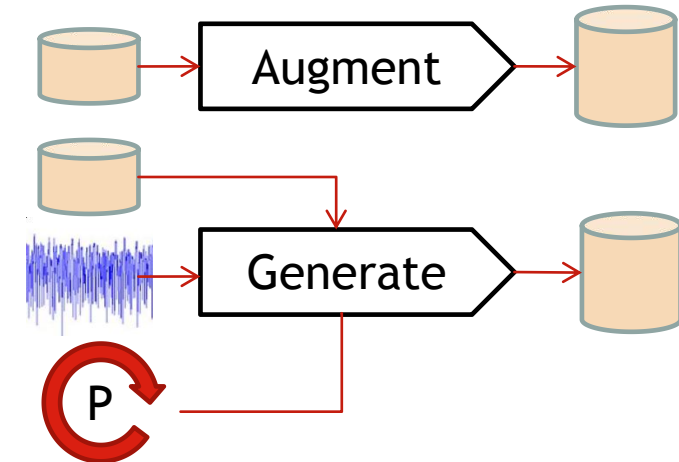
Theory
Hypothesis
Observation
Confirmation

Theory
Tentative Hypothesis
Pattern
Observation

**university of koblenz**
Computer Science

# PREDICTION VERSA GENERATION



**Motivation: Lack of parameter variance and sparse paramater space coverage of experimental data.**

**Signal Data Generation: Output is Signal**

- Model-driven or random Augmentation from real measured data (linear independent data?), additive and multiplicative noise, super-position (real-synthetic)

- Model-based or Model-driven Simulation Methods (Physics-drivem but Reality gap!)

- Model-free Random-process Generative Models (data-driven), e.g. Generative Adversarial Networks or Variational Autoencode Models

- Model-free (?) or Model-driven Parameterizable Generative Models (data-driven)



**The most important quesion of signal data generation: Produces the generator what we want - is it physically correct, and is the parameter space broadly covered?**

# EXAMPLE 1: USELESS DATA-DRIVEN MODEL

## Ultrasonic High-frequncy Pulse-Echo Measurements for Porosity Detection in Die-casted Aluminum Plates
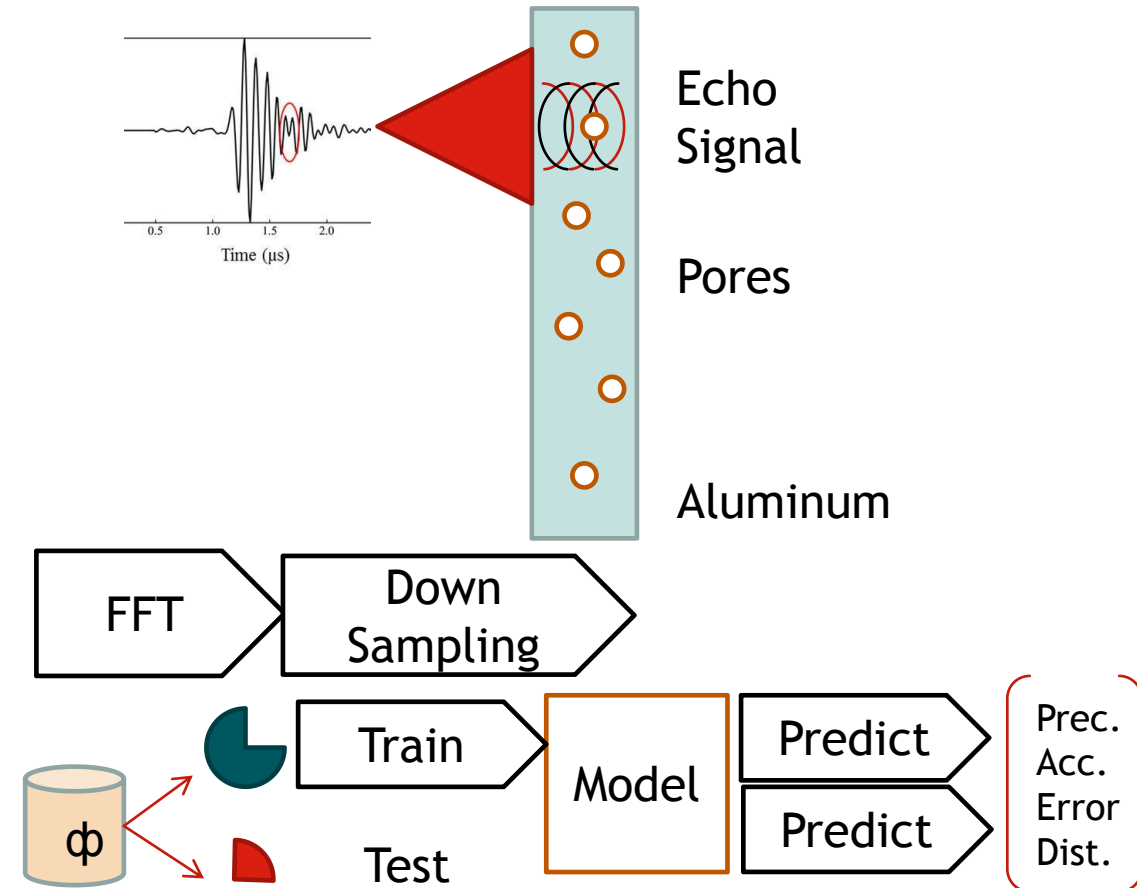
**Input**

- Time-resolved US signal (stimuls: bipolar pulse, broad frequency spectrum) , 50 specimens, 3 measuring locations

- Transducer: Dual-piezo-crystal, 5 MHz, 10 mm Dia.

- Material: Aluminum alloys (primary, secodary 58/89%)

**Features**

- Frequency spectrum of response signal (FFT)

- Assumption: Attenuation is frequency dependent and frequency spectrum is depending on pores (size, density)

**Output**

- Pore grade classification (A,B,C,D)

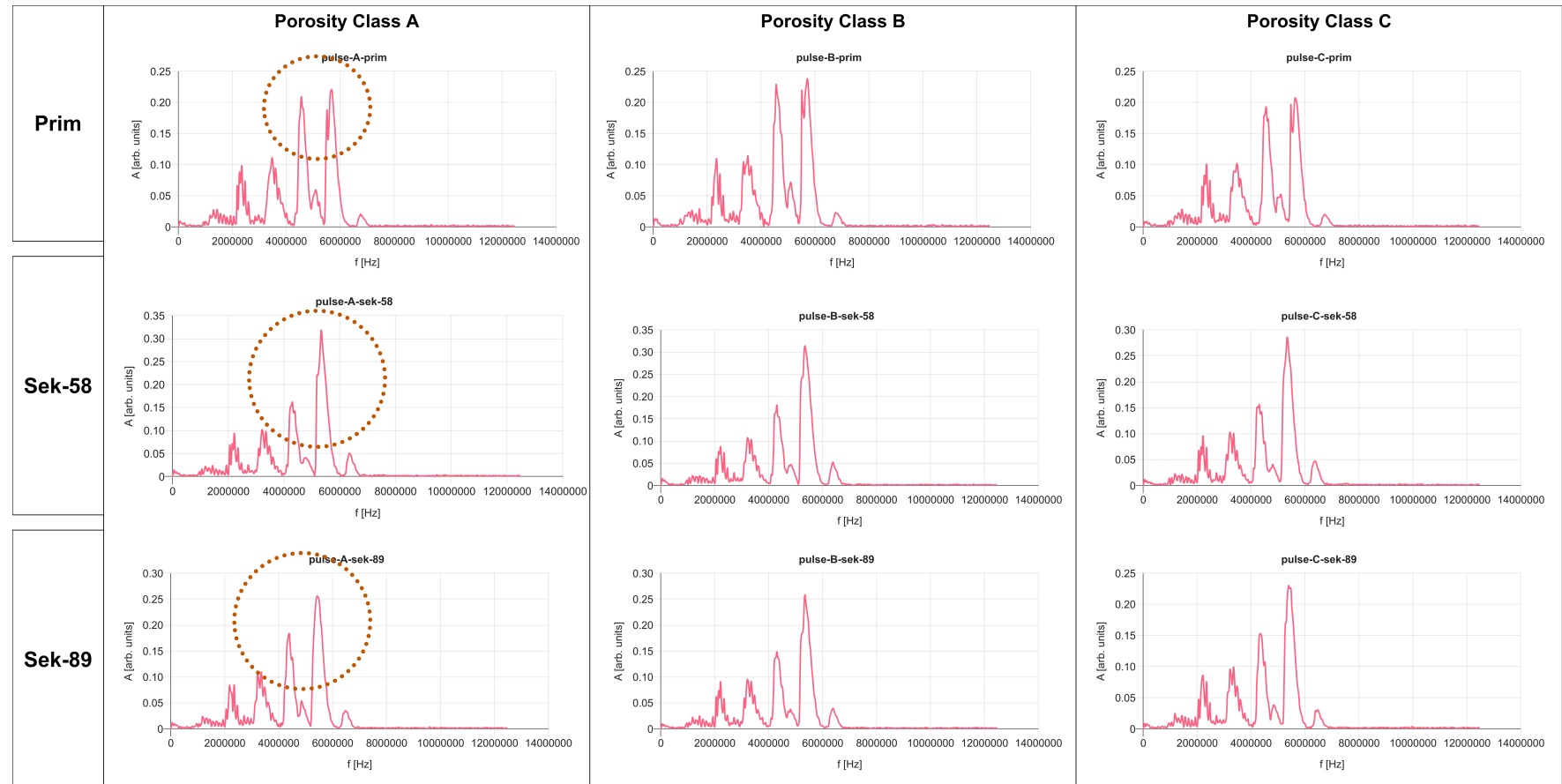- Simple ANN (two layers [10,5], sigmoid activation function softmax output layer)

# EXAMPLE 1: USELESS DATA-DRIVEN MODEL

Ultrasonic High-frequncy Pulse-Echo Measurements for Porosity Detection in Die-casted Aluminum Plates

*No porosity class correlation visible. Hidden?*

*Classification of alloy class is directly visible! Strong feature correlation*

# EXAMPLE 1: USELESS DATA-DRIVEN MODEL

Ultrasonic High-frequncy Pulse-Echo Measurements for Porosity Detection in Die-casted Aluminum Plates

**Data**

- φ: 50 specimen x 3 positions x 5 augmentation (multiplciative normal distrubuted noise 10%), mixed 3 alloys
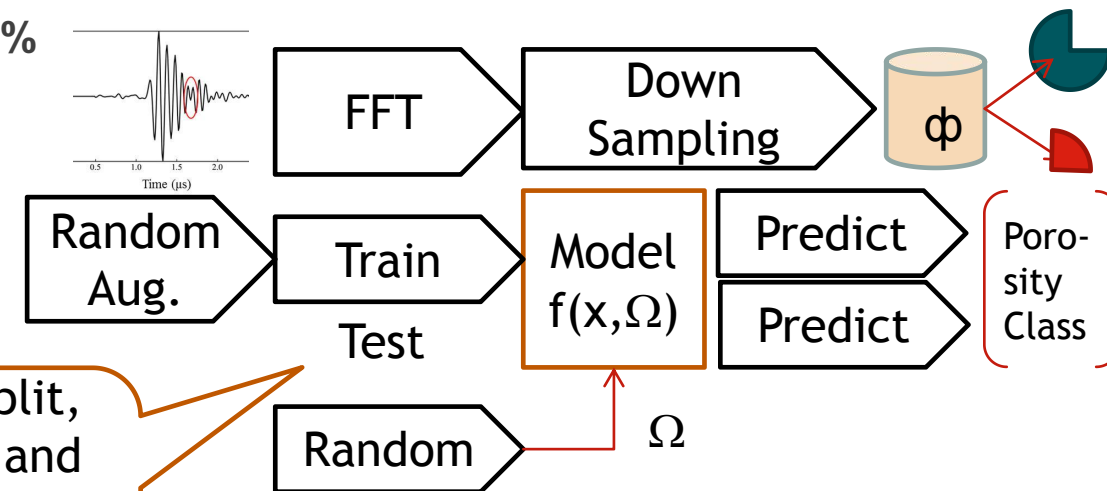
- Training/Test split: 80/20% (random)

**Results**

- Training of ANN (adam optimizer) with augmented data: Smooth and convergent!

- Classification error: **Training Data=0%..3% (!), Test=40%±20%**

**Explainibility**

- Large Training/Test error ratio: **Specialized Model**! No I/E

- Very small/weak input features were amplified resulting in a practical unusable, unpredictable and instable model

- **Unknown functional x-y relation**



Interpolation (I)/ Extrapolation (E) Tests not possible

Random data split, augmentation, and model training was repeared N times (Statistiscs)

# EXAMPLE 2: GENERATIVE MODEL

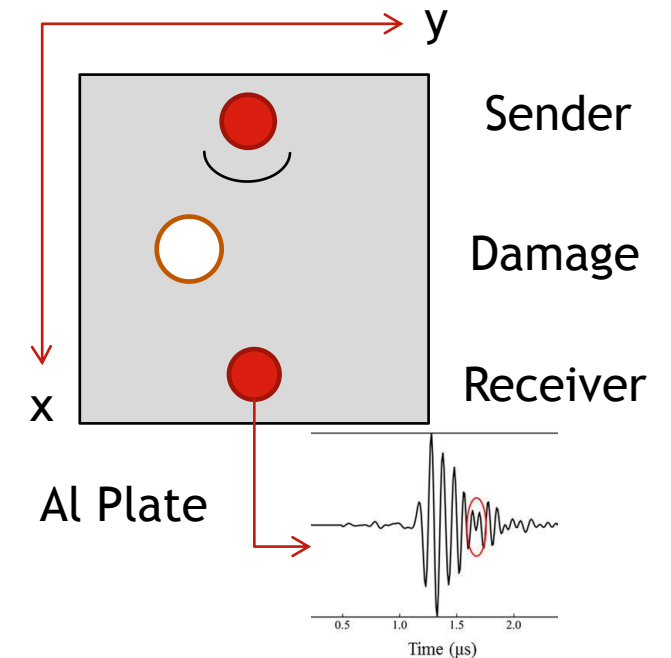Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data

**Signals**

- Time-resolved US signal (stimuls: gaussian-masked sind pulse, narrow frequency spectrum)

- Assumed measuring set-up: On sender transduce, one receiver transducer, single and straight path

- Material: Solid (e.g., Aluminum), Damage: Air, e.g. a circular hole

**Simulation and Ground-truth (GT) Data**

- Parameter space is limited - simulation using visco-elastic wave equiation is used

- A large set of data can be generated with exact labelling (GT)

**Generative Model**

- Generative Adversarial Network (GAN)

- Input: Random vector, Output: GUW signal, Training: GUW signals from simulation



$$\phi=\{f_{wave}, dim., pos_{dam}, size_{dam}, T, ..\}$$
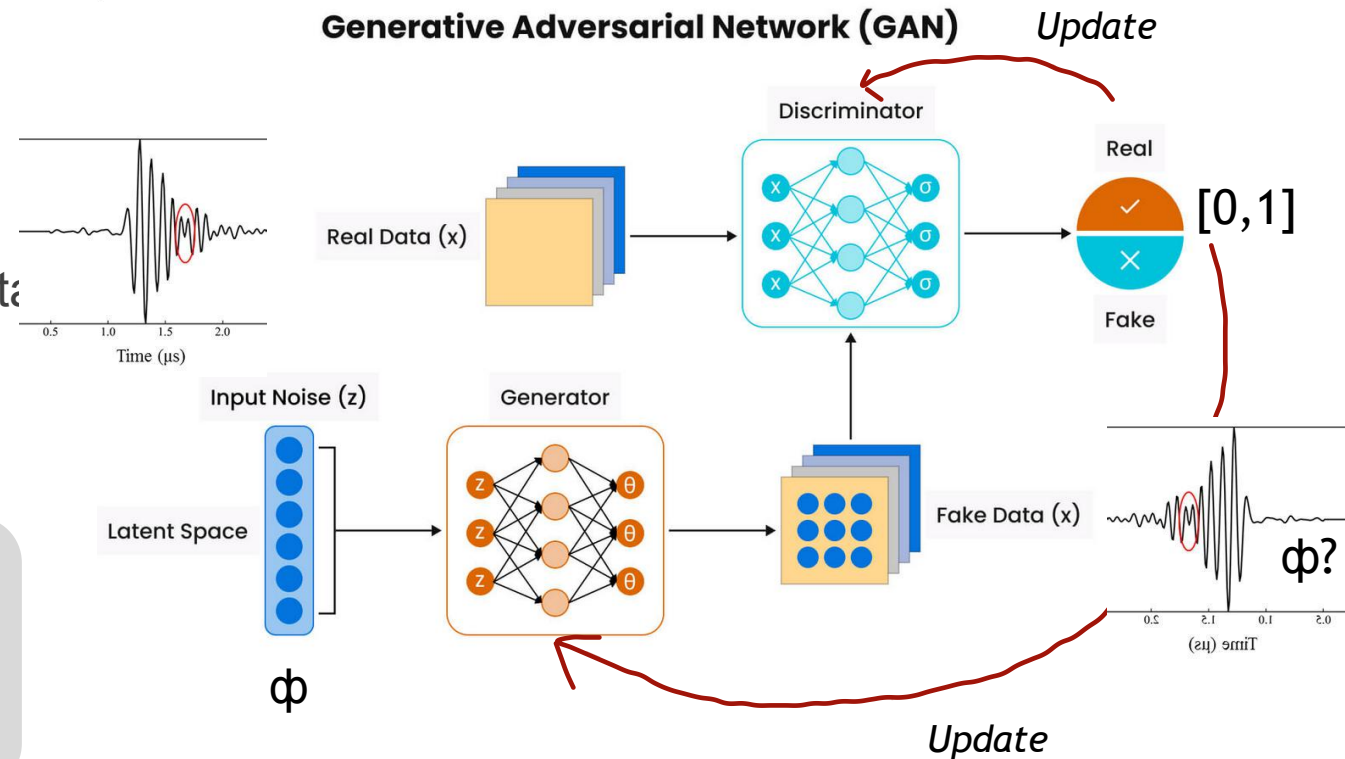
# EXAMPLE 2: GENERATIVE MODEL

Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data

**Generative Adversarial Network Model**

- Generator

- Discriminator (only Training)

- The generator nevers sees the orginal data

- Feedback only from discriminator which predicts a fake score [0,1]

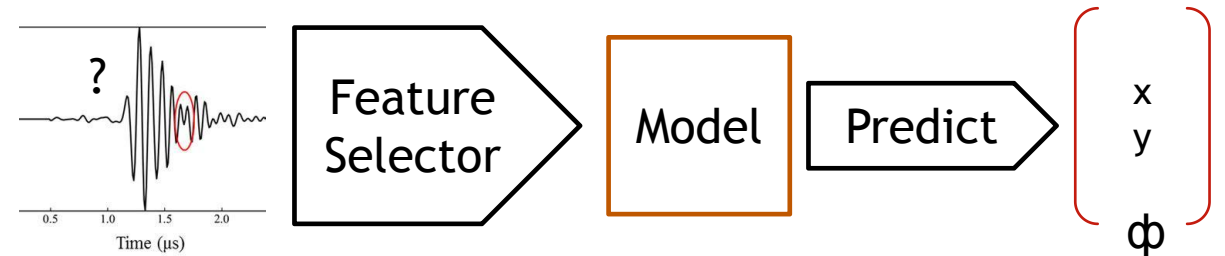⚠️ **How can we evaluate the signal generator quality? Physical correctness? Covered Generator parameter space?**

# EXAMPLE 2: PARAMETER PREDICTOR MODEL

Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data

**Parameter Predictor Model**

- Input: GUW signal (real/simu./synth.)

- Output: Parameter vector

- Here: Damage position $\mathbf{p}_{dam}$=(x,y), trained with simulation data (GT)
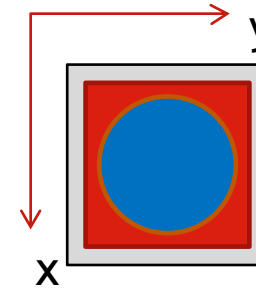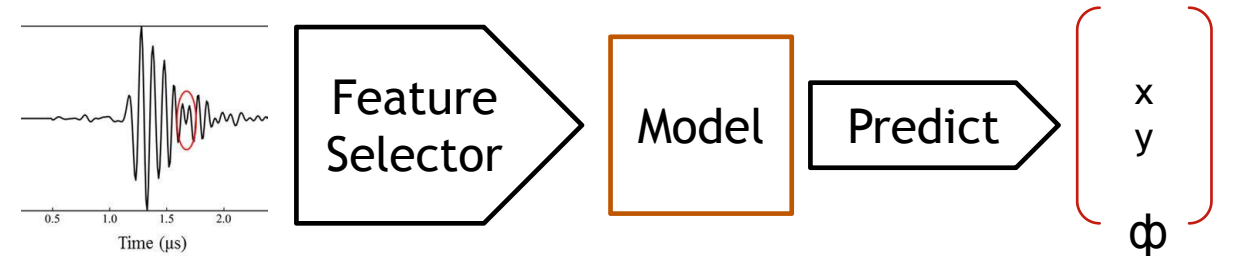
- Interpolation and Extrapolation required



⚠️ **Use another data-driven predictor model for the data parameter space! Good idea? How good is the predictor? Can we trust the model?**

# EXAMPLE 2: PARAMETER PREDICTOR MODEL
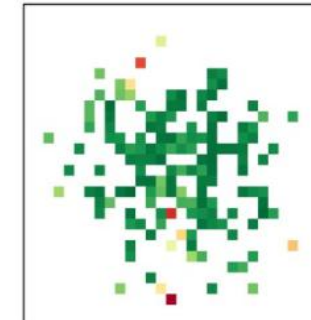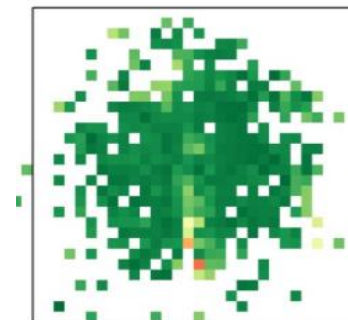
Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data
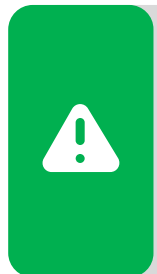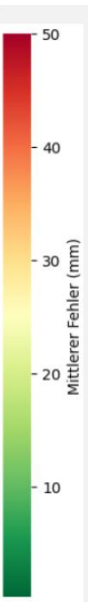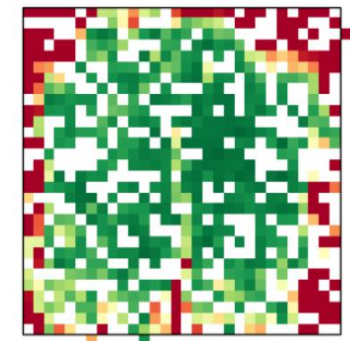
**Parameter Predictor Model: IE Test**

- Two sample sets (φ is damage location): Normal (N) and uniform (U) random distribution of damage location (x,y)

- Training with normal distributed data (80%), Test with uniform distributed data (100%)

⚠️ **It seems the model is I and partly E capable, at least we can hope!**



Regr. Error Training Data (N)   Test Data (N)          Test Data (U)

[Sidar Kilinc, 2025]

# EXAMPLE 2: PARAMETER PREDICTOR MODEL

Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data
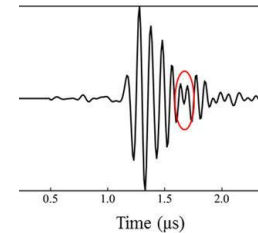
**Parameter Predictor Model: IO Analysis**

- Identify major input elements / regions contributing to the output (strong feature analysis), Input: signal *s*, Output: location coordinates (*x,y*)

- Gradient is a measure of activation:

$$\frac{\partial x}{\partial s_i} \approx \frac{\Delta x}{\Delta s_i}, \ \frac{\partial y}{\partial s_i} \approx \frac{\Delta y}{\Delta s_i}$$

⚠️ **If the prediction / regression error is low then stimulus signal regions can be identified. If the error is high, no clear correlation is visible!**
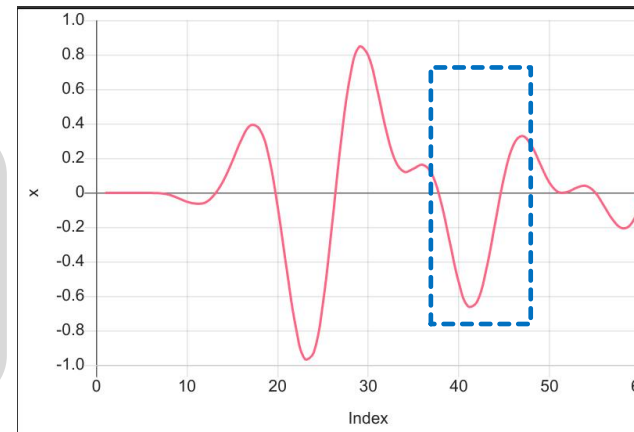


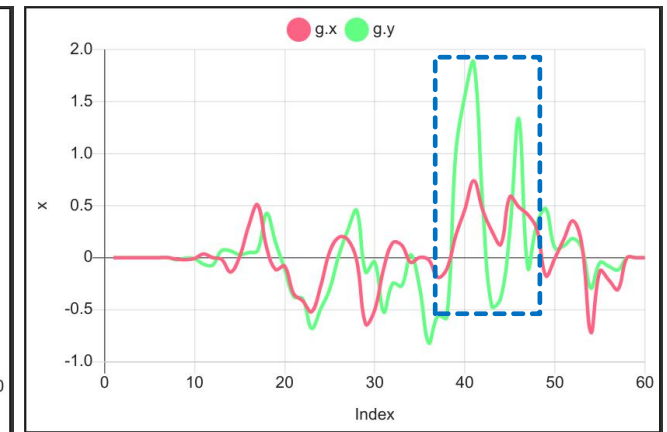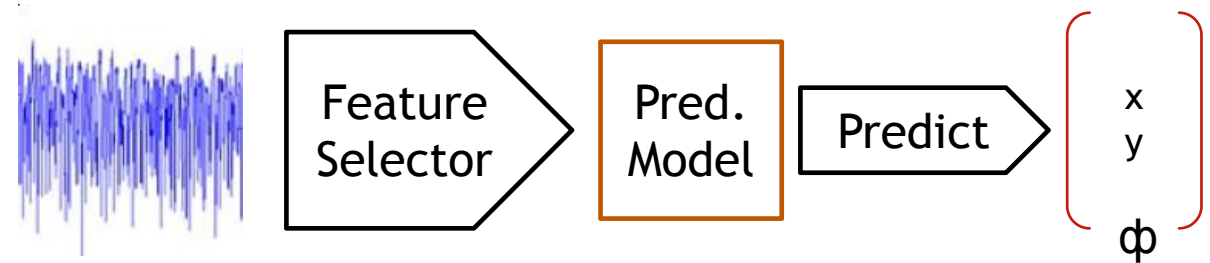| Ground Truth | x=352 | y=307 |
|---|---|---|
| Predicted | x=226 (-25%) | y=308 (0.3%) |

(CNN Model)



Signal

Gradients

# EXAMPLE 2: SCORER MODEL

Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data

**Parameter Predictor Model: Noise Test**

- Invalid and noise data test: Feed model with pure random noise and random sine waves.

- Train an additional scorer model that tests the input signal data for validity.

⚠️ **What happens if the predictor gets noise or random data? Make a test... The predictor model outputs a broad parameter range..**



Feature Selector → Pred. Model → Predict → [ x  y   φ ]

Random Sine Waves



Predict →

*Output x Distribution*

*Output y Distribution*

university of koblenz
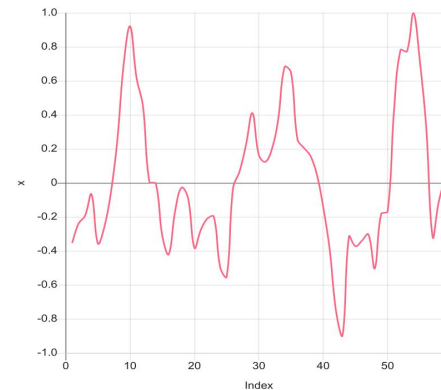Computer Science

# EXAMPLE 2: SCORER MODEL

Random-process and Data-driven Generative Model for Guided Ultrasonic Wave Data

**Parameter Predictor Model: Noise Test**

- Invalid and noise data test: Feed model with pure random noise and random sine waves.

- Train an additional scorer model that tests the input signal data for validity.
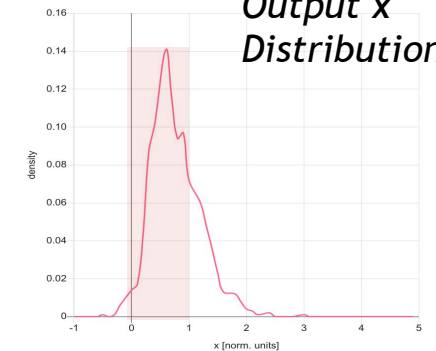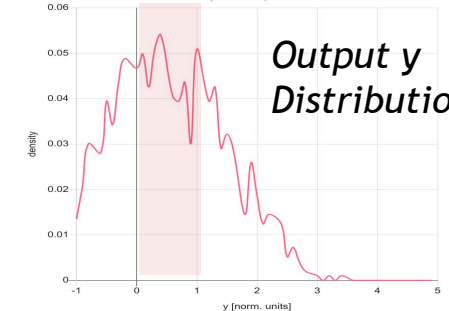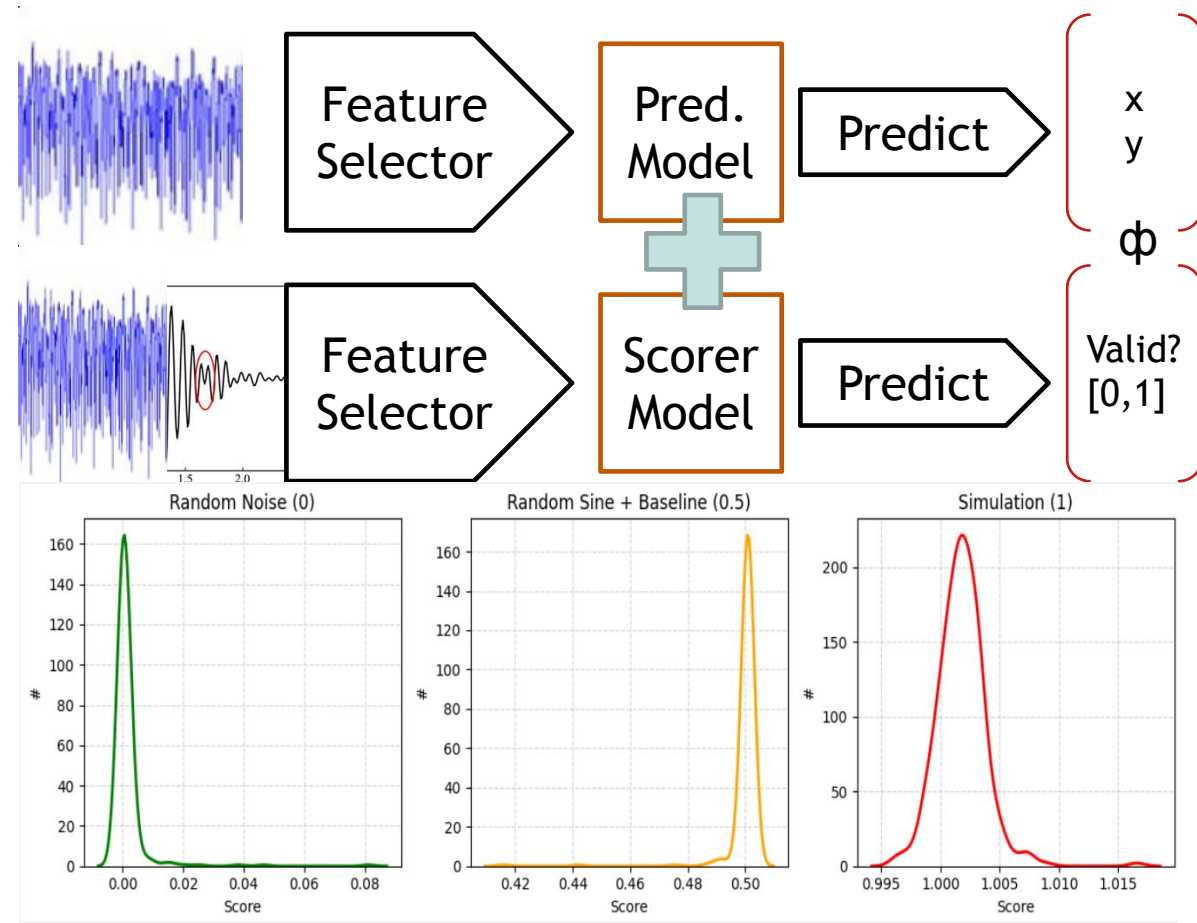
- Hierarchical model: 1. Scorer 2. Predictor

> ⚠️ **What happens if the predictor gets noise or random data? Make a test... The predictor model outputs a broad parameter range. A scorer is required.**



| Noise | Random Sine Waves | GUW |
|-------|-------------------|-----|
| 0 | 0.3 /0.5 (+ GUW) | 1.0 |

# CONCLUSIONS

| Generalization | Traceability | Explainability |
|---|---|---|
| ▪ Degree 1: Model can be used for Interpolation w/o large deviations within defined parameter Space<br><br>▪ Degree 2: Model can be used outside (trained) parameter space w/o large deviation<br><br>▪ Degree 3: The model is stable against noise and invalid data<br><br>▪ If we want to test and evaluate generative models (random process) we need degree 3! | ▪ Effect of input on model activation (paths inside a model) / Activation paths<br><br>▪ Which part or region of model input is relevant for a specific output?<br><br>▪ Model behavior with invalid or pure noise data, model selection | ▪ Do we have learned something from the data-driven model? Induction versa deduction?<br><br>▪ Why is the model giving a specific output for a specific input?<br><br>▪ Surrogate models can help to reduce a complex model to simplified and explainable well known functional laws<br><br>▪ Can we solve the inverse problem with a specific model?<br><br>▪ Can we explain generative models? |

university
of koblenz
Computer Science

# THANK YOU

Stefan Bosse

sbosse@uni-koblenz.de

www.edu-9.de