

Application of Combinatorial Testing in Testing ML Systems

Workshop on Combinatorial Testing for AI-enabled Systems
September 4, 2024

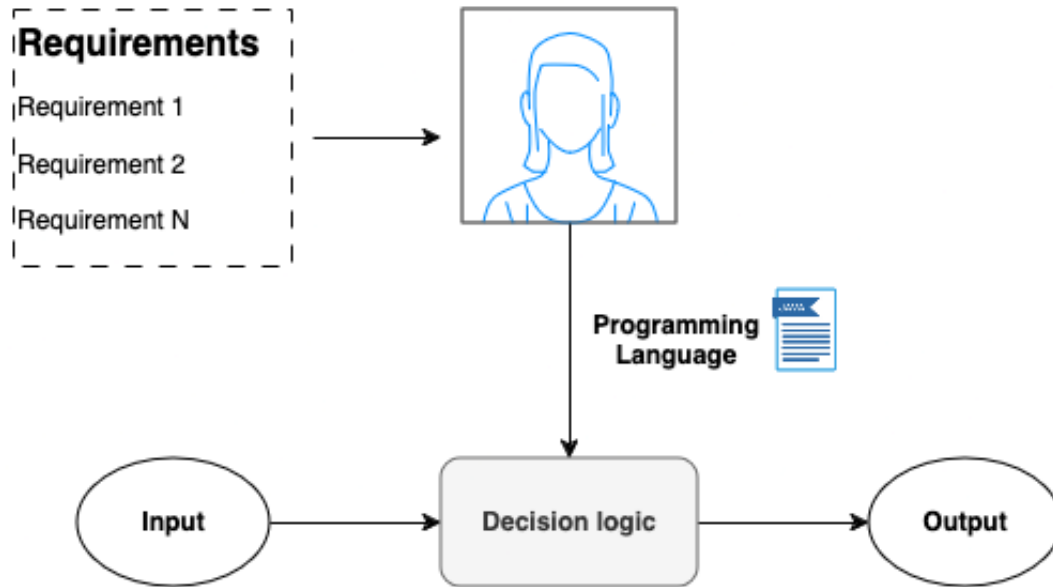
Jaganmohan Chandrasekaran, Ph.D.
Research Assistant Professor
Sanghani Center for AI & Data Analytics
Virginia Tech | jagan@vt.edu

Outline

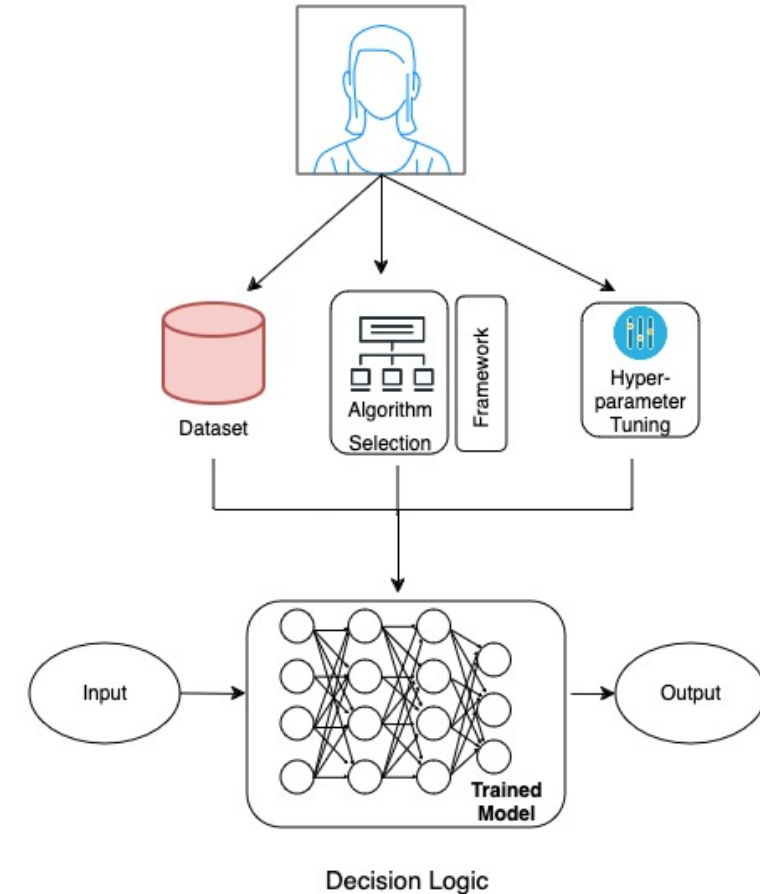
- Introduction
- Challenges in Testing ML Systems
- Combinatorial Testing for ML Systems
- Overview of Combinatorial Testing in Evaluating ML Systems
- Challenges and Future Directions

Traditional Software vs. ML-enabled Software Development

Traditional Software Development



ML-enabled Software Development

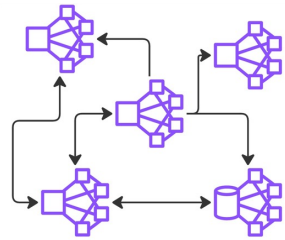


Unique Testing Challenges in ML



Data-Intensive

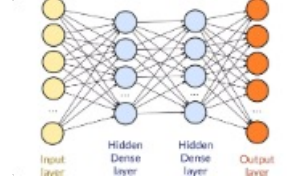
AI/ML systems rely heavily on data, significantly influencing behavior.



Large Input Space

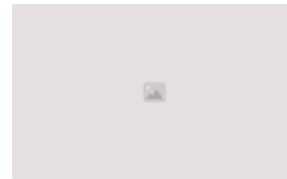
Decision logic is derived from the patterns underlying the training dataset, which is usually large in size.

9/6/24



Complex Decision Logic

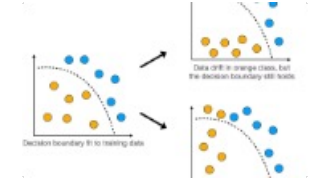
The internal decision logic is generally opaque and not easily interpretable by humans.



Lack of Test Oracle

In most cases, test instances will lack a reliable test oracle, and desired outcomes are unknown beforehand.

Application of CT in Testing ML Systems



Behavior

Non-deterministic behavior
Evolving behavior



Testing beyond Correctness

Fairness, safety, and trustworthiness are critical. Specialized testing techniques are required.

Outline

- Introduction
- Challenges in Testing ML Systems
- **Combinatorial Testing for ML Systems**
- Overview of Combinatorial Testing in Evaluating ML Systems
- Challenges and Future Directions

Why Combinatorial Testing for ML?

Combinatorial testing, a black box test generation technique, focuses on systematically testing the interactions among the system's different parameters with a relatively smaller number of tests [1].

Enables the detection of interaction faults, including rare combinations that can lead to failures in software systems [2].

Why Combinatorial Testing for ML (2)?

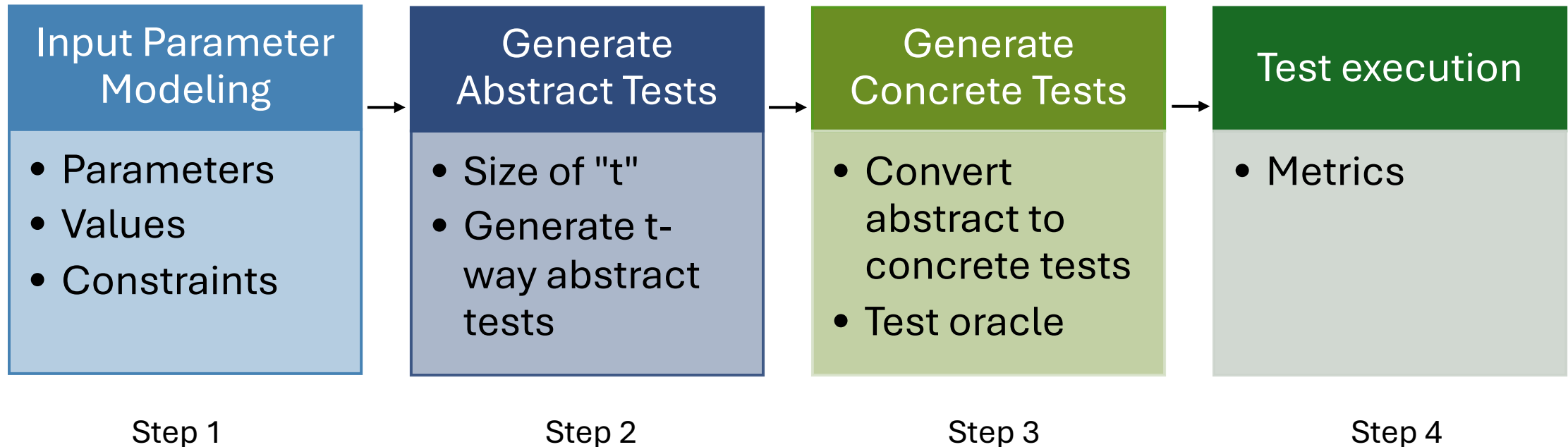
Combinatorial testing, a black box test generation technique, focuses on systematically testing the interactions among the system's different parameters with a relatively smaller number of tests.

Enables the detection of interaction faults, including rare combinations that can lead to failures in software systems.

Efficiently explore the large input space of ML models by systematically testing different feature interactions with a relatively small number of tests.

Help uncover rare feature combinations and edge case scenarios, thus reducing the risk of undetected faults.

Steps in Applying Combinatorial Testing to ML



Applying Combinatorial Testing - Image Dataset



1 Segmentation



Number of segments = 20
Parameter = Segment
Values = [True, False]

2

Value of t = 2

ACTS - ACTS Main Window

System Edit Operations Help

Algorithm: IPOG Strength: 2 Status : A new test set is built.

System View







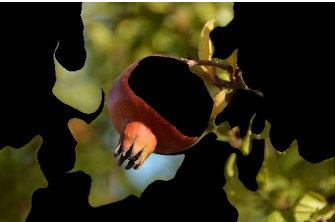

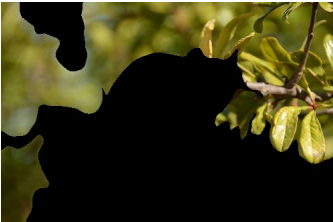


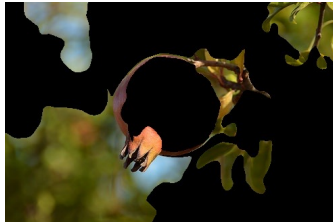
- [Root Node]
- [SYSTEM-Image_1]
- Segment_1
 - true
 - false
- Segment_2
 - true
 - false
- Segment_3
 - true
 - false
- Segment_4
 - true
 - false
- Segment_5
 - true
 - false
- Segment_6
 - true
 - false
- Segment_7
- Segment_8
- Segment_9
- Segment_10
- Segment_11
- Segment_12
- Segment_13
- Segment_14
- Segment_15
- Segment_16
- Segment_17
- Segment_18
- Segment_19
- Segment_20

Test Result Statistics

	Segment_1	Segment_2	Segment_3	Segment_4	Segment_5	Segment_6
1	true	true	false	false	false	false
2	true	false	true	true	true	true
3	false	true	true	false	true	false
4	false	false	false	true	false	true
5	false	true	false	true	true	false
6	false	false	true	false	false	true
7	true	true	false	false	false	true
8	false	false	true	true	true	false
9	false	true	false	false	false	false
10	true	false	true	true	true	true
11	false	true	true	false	false	true
12	true	false	true	true	true	true

Abstract test cases

Generating test cases - Image Dataset

Test case 1	Test case 2	Test case 3	Test case 4	Test case 5	Test case 6
					
Test case 7	Test case 8	Test case 9	Test case 10	Test case 11	Test case 12
					

Applying Combinatorial Testing - Tabular Dataset

- Parameters
 - Map all attributes, excluding the class label.
- Identify values for each parameter based on the attribute's data type
 - Categorical** attributes: Map unique values as representative values
 - Numerical** attributes: Identify representative values through discretization
- Identify constraints from the training dataset

1	age	workclass	education	education_n	marital_stat	occupation	relationship	race	sex
2	39	State-gov	Bachelors	13	Never-marri	Adm-clerical	Not-in-famil	White	Male
3	50	Self-emp-no	Bachelors	13	Married-civ-	Exec-managi	Husband	White	Male
4	38	Private	HS-grad	9	Divorced	Handlers-cle	Not-in-famil	White	Male
5	53	Private	11th	7	Married-civ-	Handlers-cle	Husband	Black	Male
6	28	Private	Bachelors	13	Married-civ-	Prof-specialt	Wife	Black	Female
7	37	Private	Masters	14	Married-civ-	Exec-managi	Wife	White	Female
8	49	Private	9th	5	Married-spo	Other-servic	Not-in-famil	Black	Female
9	52	Self-emp-no	HS-grad	9	Married-civ-	Exec-managi	Husband	White	Male
10	31	Private	Masters	14	Never-marri	Prof-specialt	Not-in-famil	White	Female
11	42	Private	Bachelors	13	Married-civ-	Exec-managi	Husband	White	Male
12	37	Private	Some-colleg	10	Married-civ-	Exec-managi	Husband	Black	Male
13	30	State-gov	Bachelors	13	Married-civ-	Prof-specialt	Husband	Asian-Pac-Isl	Male
14	23	Private	Bachelors	13	Never-marri	Adm-clerical	Own-child	White	Female
15	32	Private	Assoc-acdm	12	Never-marri	Sales	Not-in-famil	Black	Male
16	40	Private	Assoc-voc	11	Married-civ-	Craft-repair	Husband	Asian-Pac-Isl	Male
17	34	Private	7th-8th	4	Married-civ-	Transport-m	Husband	Amer-Indian	Male
18	25	Self-emp-no	HS-grad	9	Never-marri	Farming-fish	Own-child	White	Male
19	32	Private	HS-grad	9	Never-marri	Machine-op-	Unmarried	White	Male

A snippet of the Adult Income Dataset [12, 13]

Designing an Input Parameter Model - Tabular Dataset

Discretized Numerical Values

1	age	workclass	education	education_num	marital_status	occupation	relationship	race	sex
2	39	State-gov	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male
3	50	Self-emp-not-inc	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male
4	38	Private	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male
5	53	Private	11th	7	Married-civ-spouse	Handlers-cleaners	Husband	Black	Male
6	28	Private	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female
7	37	Private	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female
8	49	Private	9th	5	Married-spouse-absent	Other-service	Not-in-family	Black	Female
9	52	Self-emp-not-inc	HS-grad	9	Married-civ-spouse	Exec-managerial	Husband	White	Male
10	31	Private	Masters	14	Never-married	Prof-specialty	Not-in-family	White	Female
11	42	Private	Bachelors	13	Married-civ-spouse	Exec-managerial	Husband	White	Male
12	37	Private	Some-college	10	Married-civ-spouse	Exec-managerial	Husband	Black	Male
13	30	State-gov	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband	Asian-Pac-Islander	Male
14	23	Private	Bachelors	13	Never-married	Adm-clerical	Own-child	White	Female
15	32	Private	Assoc-acdm	12	Never-married	Sales	Not-in-family	Black	Male
16	40	Private	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband	Asian-Pac-Islander	Male
17	34	Private	7th-8th	4	Married-civ-spouse	Transport-moving	Husband	Amer-Indian-Eskimo	Male
18	25	Self-emp-not-inc	HS-grad	9	Never-married	Farming-fishing	Own-child	White	Male
19	32	Private	HS-grad	9	Never-married	Machine-operating	Unmarried	White	Male

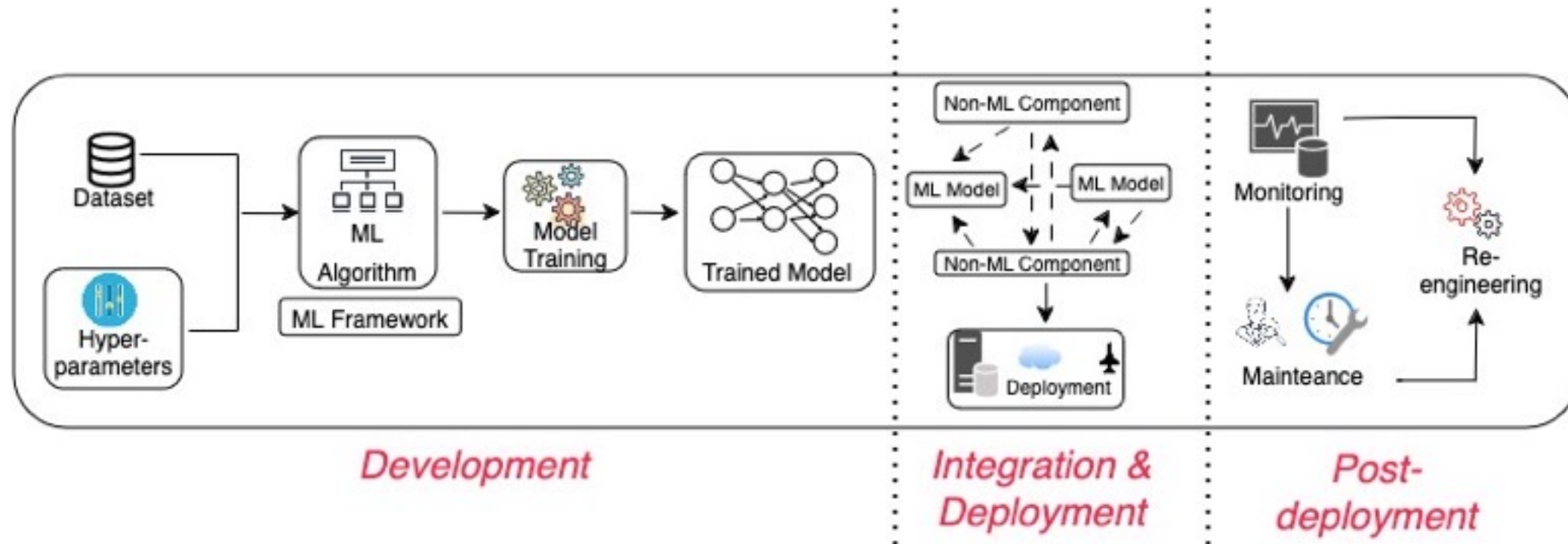
Categorical Values

Parameters

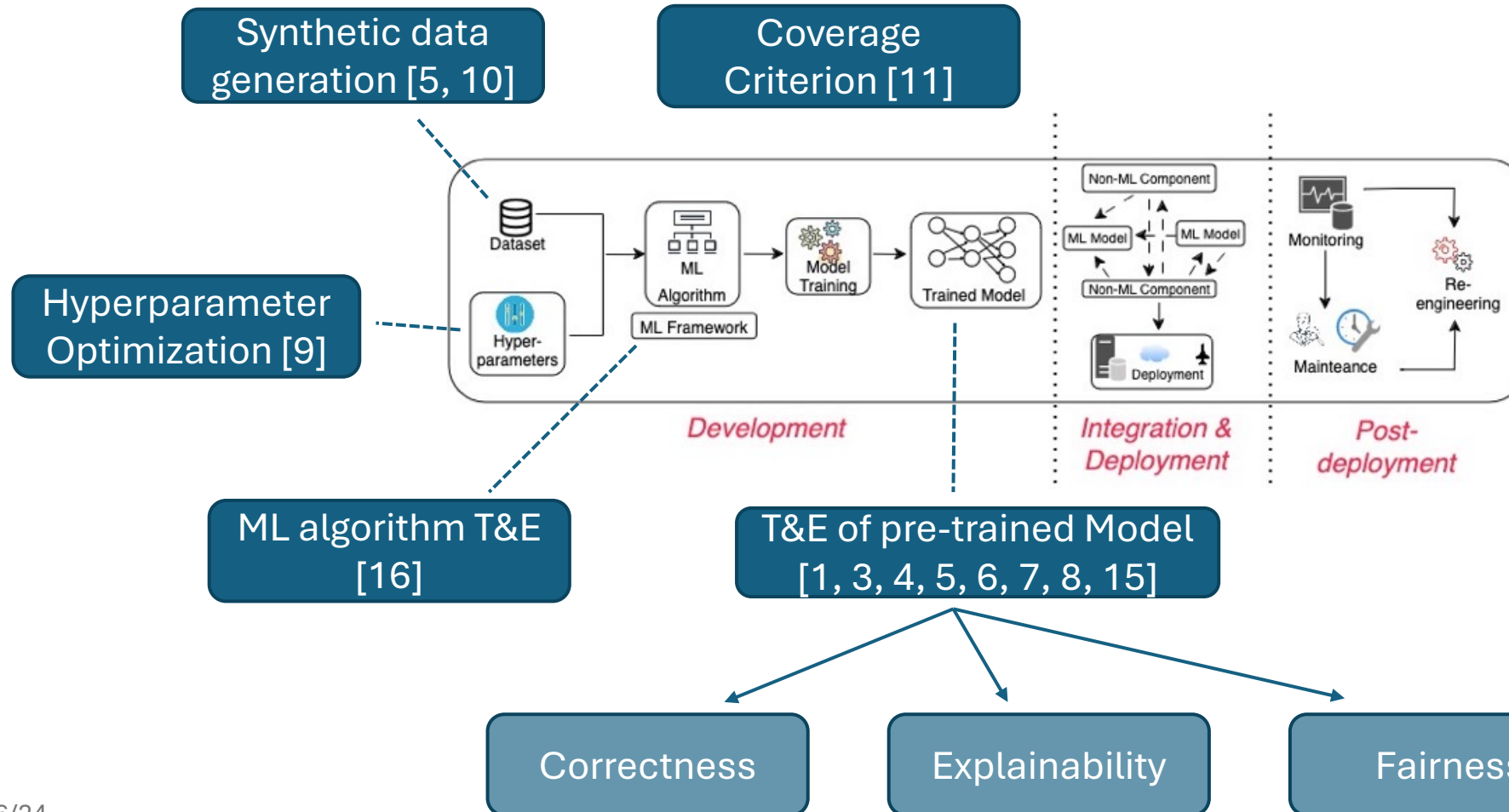
Outline

- Introduction
- Challenges in Testing ML Systems
- Combinatorial Testing for ML Systems
- **Overview of Combinatorial Testing in Evaluating ML Systems**
- Challenges and Future Directions

Lifecycle - ML-enabled software system



Application of Combinatorial Testing in the ML Lifecycle

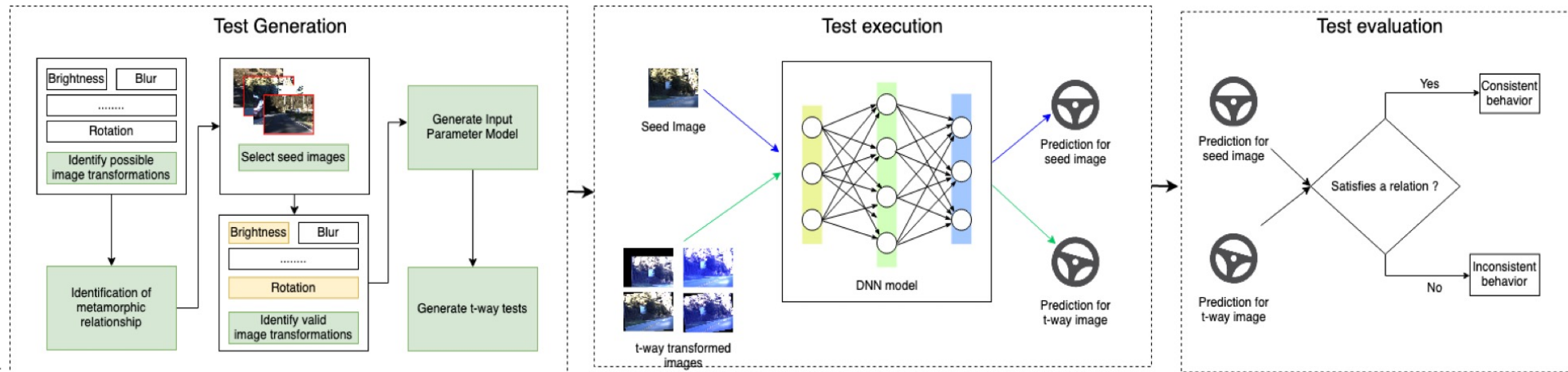


Testing pre-trained Deep Neural Network models

- Generate t-way synthetic images to test Deep Neural Networks (DNNs) used in self-driving cars
- Each **synthetic image** is a **combination of image transformations**.

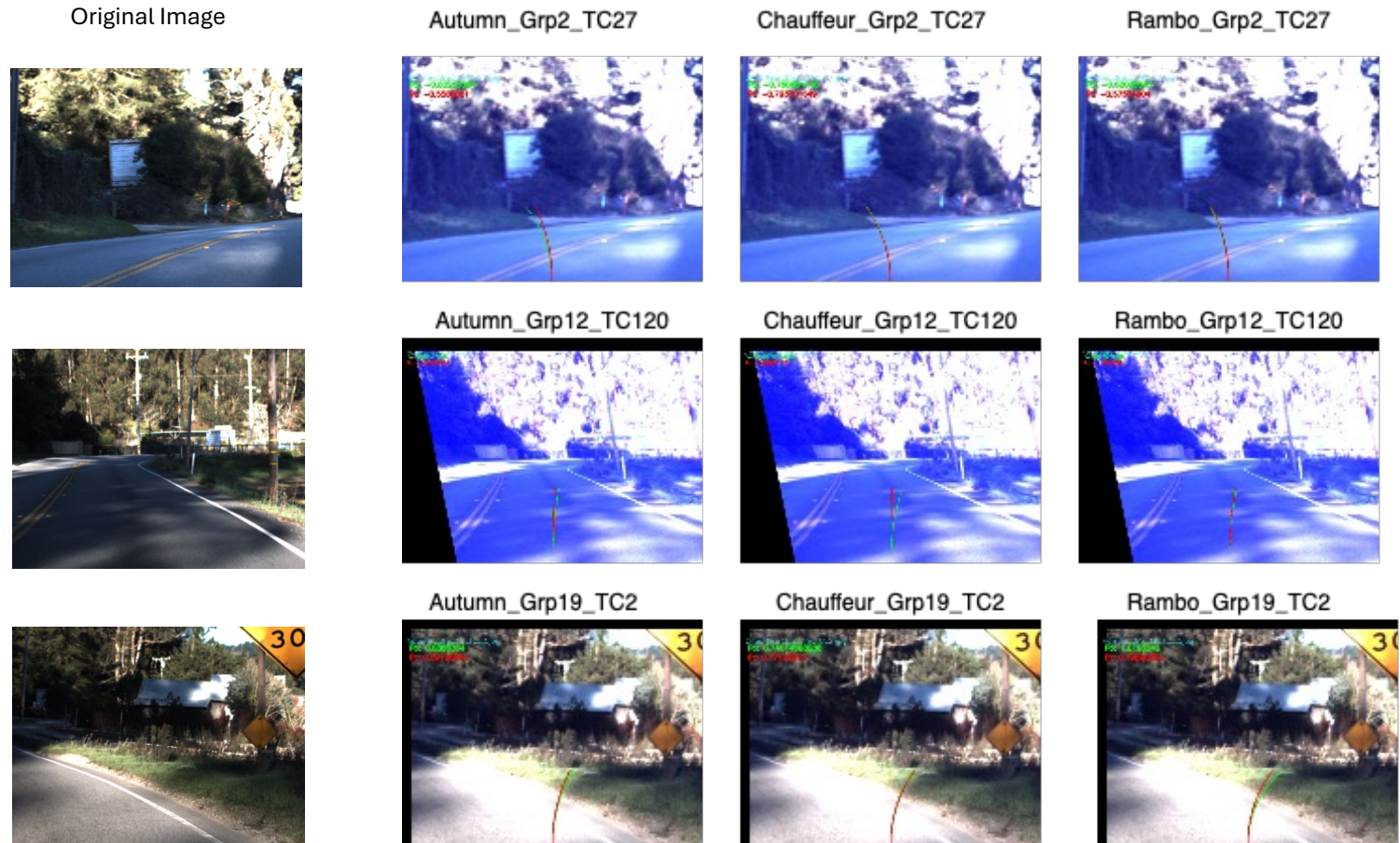
A set of image transformations

Transformations		Values
Blur	Averaging	3x3, 4x4, 5x5, 6x6
	Gaussian	3x3, 5x5, 7x7
	Median	3, 5
	Bilateral	(9, 75, 75)
Brightness		10, 20, 30, 40, 50, 60, 70, 80, 90, 100
Contrast		1.2, 1.4, 1.6, 1.8, 2.0, 2.2, 2.4, 2.6, 2.8, 3.0
Rotation		3, 6, 9, 12, 15, 18, 21, 24, 27, 30
Scale		1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0, 5.5, 6.0
Shear (Horizontal)		0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0
Translation		(10,10), (20,20), (30,30), (40,40), (50,50), (60,60), (70,70), (80,80), (90,90), (100,100)



Testing pre-trained Deep Neural Network models (2)

- Synthetic images generated by combining a set of image transformations (t-way tests) can **successfully identify inconsistent behavior**.
- T-way synthetic images **can significantly increase neuron coverage**, a measure of the proportion of neurons activated in a DNN model.



T-way tests successfully uncover behavior inconsistencies among pre-trained models

Explaining Model Decisions

Explaining a model's decision is similar to the fault localization

- Fault localization: Given a failing scenario, a software engineer *identifies which part of the input that causes the failure*
- Explainable AI: Given a decision (made by an ML model), the objective is to *identify features that influence the model's decision.*

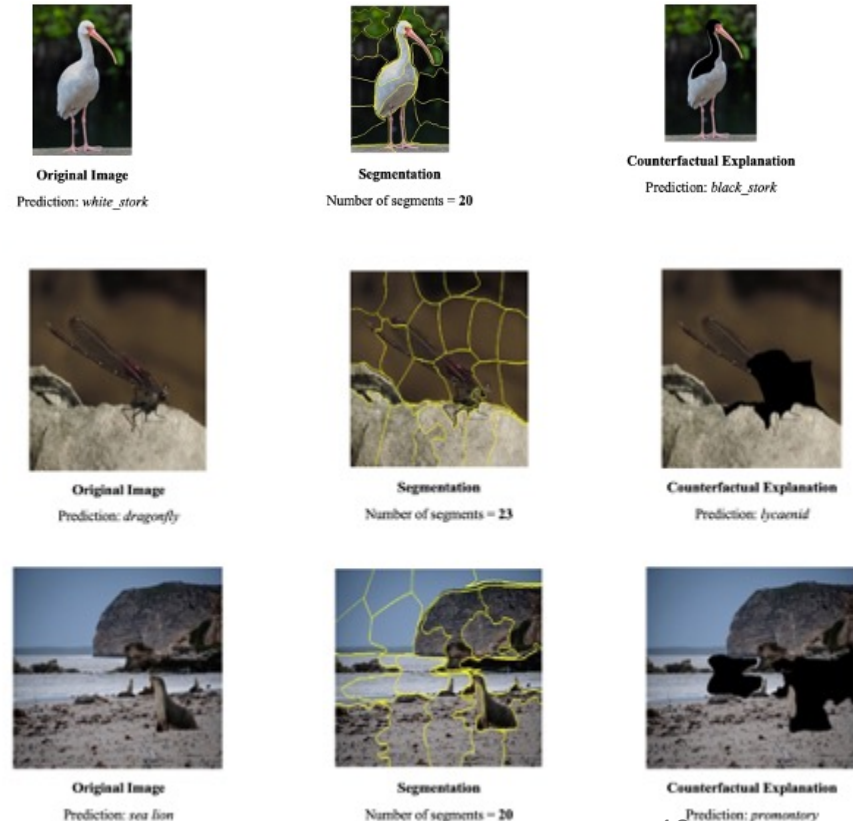
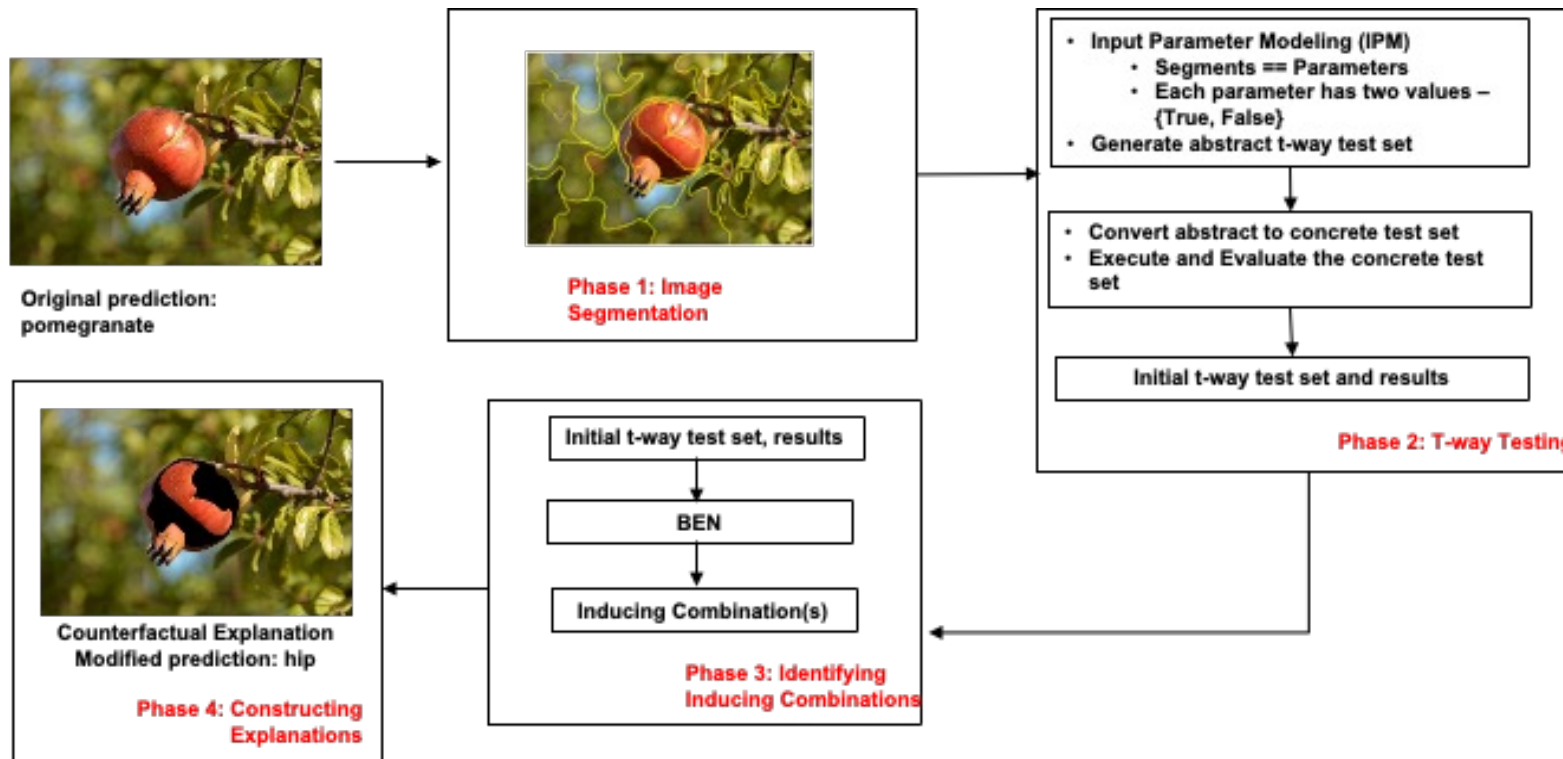
Can we use BEN, a combinatorial testing-based software fault localization tool, to generate an explanation?

Explaining Model Decisions (2)

Explaining the behavior of image classifiers through counterfactual explanations

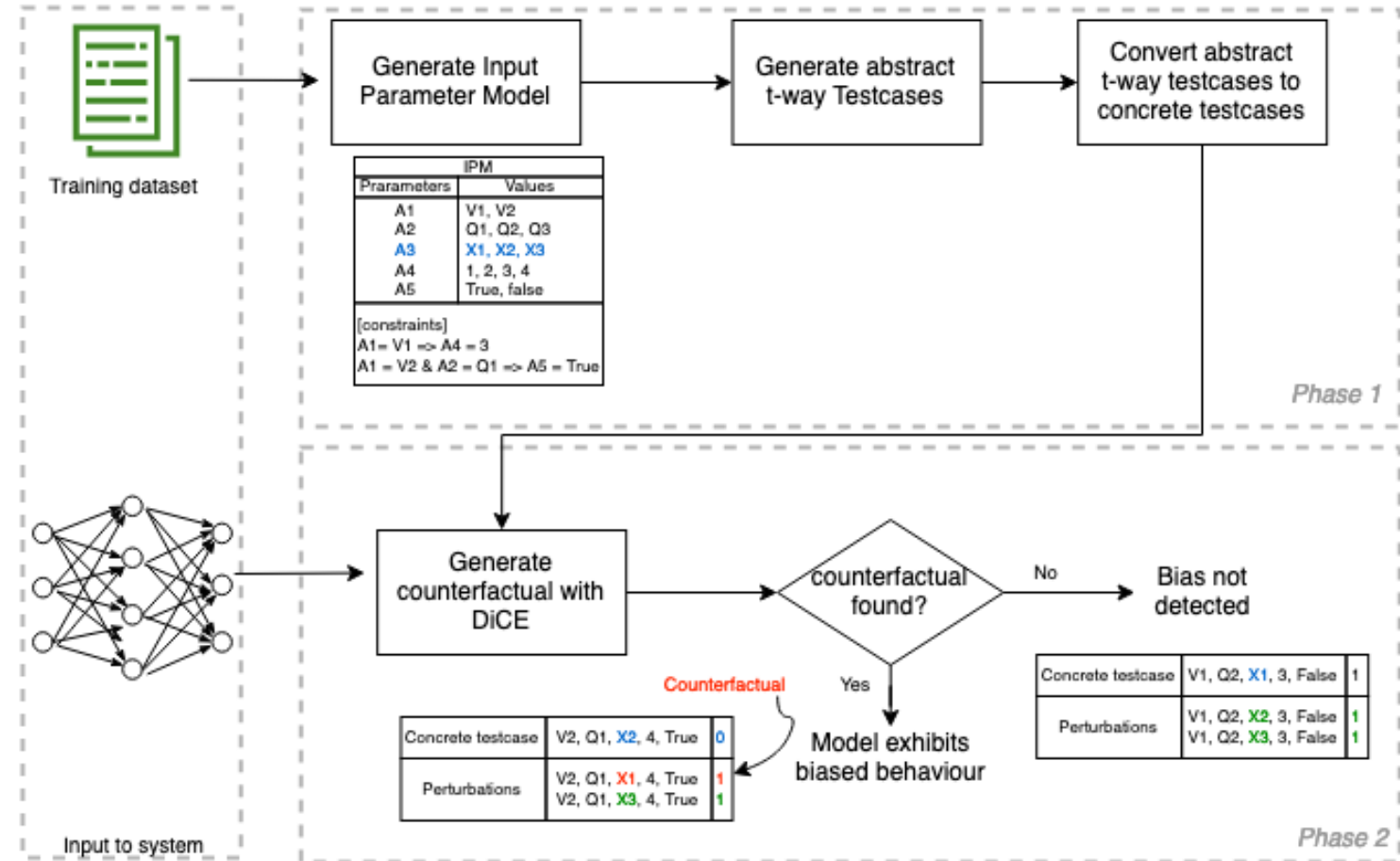
Determine a **minimal subset of features** in the input image that, **if removed, would result in a different classification.**

The CT-based approach effectively generates a counterfactual explanation for 44 out of 50 (~90%) images.



Fairness Testing of ML models

- Evaluating the discriminatory behavior of pre-trained ML models using combinatorial testing
- A model-agnostic approach that can successfully identify fairness violations in ML models



Fairness Testing of ML models (2)

Results suggest that t-way tests effectively identify biases introduced by the training dataset and the learning algorithm.

Datasets	# of t-way tests	Protect ed Attribu tesx	Number of fairness violations			
			Logisti c Regres sion (LR)	Rando m Forest (RF)	Support Vector Machin es (SVM)	Deep Neural Network (DNN)
Adult Income	676	Sex, Race	58	25	23	11
German Credit	81	Sex, Age	5	10	9	26
COMPAS	64	Sex, Race	27	37	24	4

Fairness violations identified by t-way tests

Dataset	Attribute Modified	# CF found			
		LR	RF	SVM	DNN
Adult Income	Only Race	16	10	5	4
	Only Sex	11	1	3	5
	Both Race and Sex	31	14	16	2
	Total	58	25	23	11

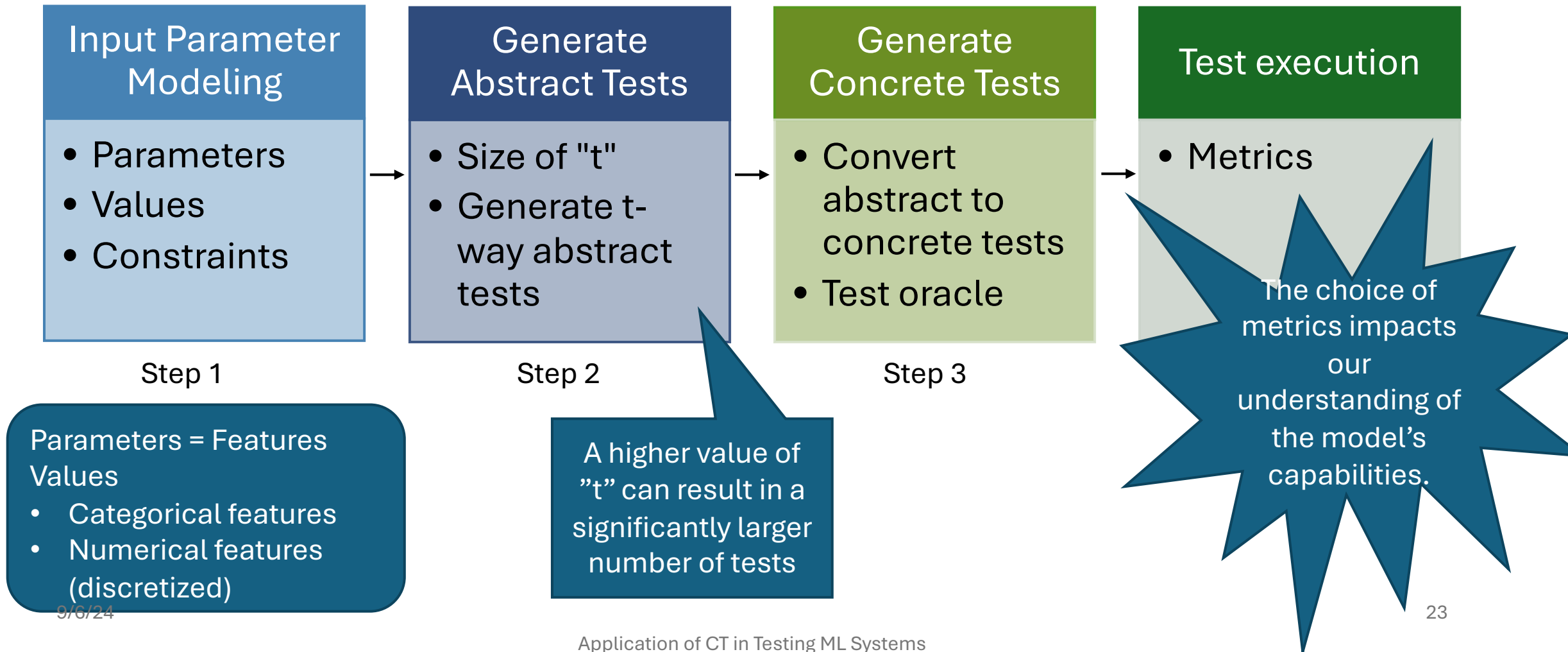
Dataset	Attribute Modified	# CF found			
		LR	RF	SVM	DNN
German Credit	Only Race	0	0	0	4
	Only Sex	0	0	0	10
	Both Race and Sex	5	10	9	12
	Total	5	10	9	26

Dataset	Attribute Modified	# CF found			
		LR	RF	SVM	DNN
COMPAS	Only Race	11	19	5	2
	Only Sex	1	3	1	1
	Both Race and Sex	15	15	18	1
	Total	27	37	24	4

Outline

- Introduction
- Challenges in Testing AI/ML Systems
- Combinatorial Testing for AI/ML Systems
- Overview of Combinatorial Testing in Evaluating AI/ML Systems
- **Challenges and Future Directions**

Challenges in Adapting CT for ML

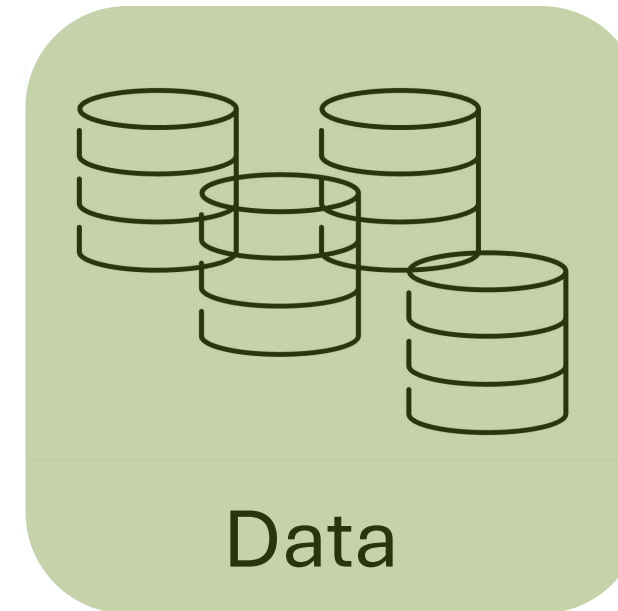
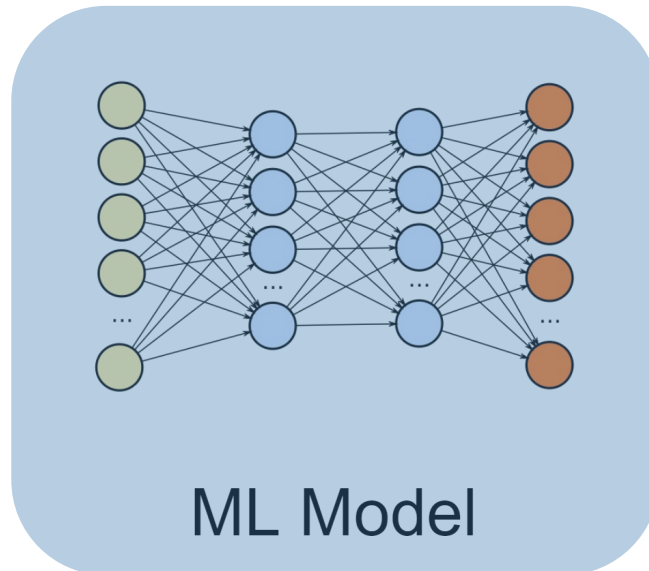


Future Directions

- Model debugging
- Post deployment
 - Monitoring
 - Regression testing
- Other data modalities
 - Natural Language Processing Datasets

Next Session...

The performance of ML systems is directly influenced by the quality of the underlying ML model and the data used to train it.



This session

Next session

Thank you!!!!

Collaborators*

Tyler Cody

Laura Freeman

Erin Lanus

Jeff Lei

Brian Lee

Raghu Kacker

Krishna Khadka

Rick Kuhn

M S Raunak

Ankita Patel

& all the anonymous reviewers
from the community



9/6/24
* Listed in alphabetical order

References

1. Chandrasekaran, J., Lei, Y., Kacker, R., & Kuhn, D. R. (2021, April). A combinatorial approach to testing deep neural network-based autonomous driving systems. In 2021 IEEE international conference on software testing, verification and validation workshops (ICSTW) (pp. 57-66). IEEE.
2. Kuhn, R., Lei, Y., & Kacker, R. Practical Combinatorial Testing: Beyond Pairwise. Kuhn, R., Lei, Y., & Kacker, R. Practical Combinatorial Testing: Beyond Pairwise.
3. Ma, L., Juefei-Xu, F., Xue, M., Li, B., Li, L., Liu, Y., & Zhao, J. (2019, February). Deepct: Tomographic combinatorial testing for deep learning systems. In 2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER) (pp. 614-618). IEEE.
4. Gladisch, C., Heinzemann, C., Herrmann, M., & Woehrle, M. (2020). Leveraging combinatorial testing for safety-critical computer vision datasets. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (pp. 324-325).
5. Dola, S., McDaniel, R., Dwyer, M. B., & Soffa, M. L. (2024, April). CIT4DNN: Generating Diverse and Rare Inputs for Neural Networks Using Latent Space Combinatorial Testing. In Proceedings of the IEEE/ACM 46th International Conference on Software Engineering (pp. 1-13)
6. Kuhn, D. R., Kacker, R. N., Lei, Y., & Simos, D. E. (2020, October). Combinatorial methods for explainable AI. In 2020 IEEE international conference on software testing, verification and validation workshops (ICSTW) (pp. 167-170). IEEE.
7. Chandrasekaran, J., Lei, Y., Kacker, R., & Kuhn, D. R. (2021, April). A combinatorial approach to explaining image classifiers. In 2021 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW) (pp. 35-43). IEEE.
8. Patel, A. R., Chandrasekaran, J., Lei, Y., Kacker, R. N., & Kuhn, D. R. (2022, April). A combinatorial approach to fairness testing of machine learning models. In 2022 IEEE international conference on software testing, verification and validation workshops (ICSTW) (pp. 94-101). IEEE.
9. Khadka, K., Chandrasekaran, J., Lei, Y., Kacker, R. N., & Kuhn, D. R. (2024, April). A Combinatorial Approach to Hyperparameter Optimization. In Proceedings of the IEEE/ACM 3rd International Conference on AI Engineering-Software Engineering for AI (pp. 140-149).
10. Khadka, K., Chandrasekaran, J., Lei, Y., Kacker, R. N., & Kuhn, D. R. (2023, April). Synthetic data generation using combinatorial testing and variational autoencoder. In 2023 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW) (pp. 228-236). IEEE.
11. Sekhon, J., & Fleming, C. (2019, May). Towards improved testing for deep learning. In 2019 IEEE/ACM 41st International Conference on Software Engineering: New Ideas and Emerging Results (ICSE-NIER) (pp. 85-88). IEEE.
12. Becker, Barry and Kohavi, Ronny. (1996). Adult. UCI Machine Learning Repository. <https://doi.org/10.24432/C5XW20>.
13. UCI Machine Learning Repository: Adult Data Set | <https://archive.ics.uci.edu/dataset/2/adult>
14. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115, 211-252.
15. Cody, T., Lanus, E., Doyle, D. D., & Freeman, L. (2022, April). Systematic training and testing for machine learning using combinatorial interaction testing. In 2022 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW) (pp. 102-109). IEEE.
16. Chandrasekaran, J., Feng, H., Lei, Y., Kuhn, D. R., & Kacker, R. (2017, March). Applying combinatorial testing to data mining algorithms. In 2017 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW) (pp. 253-261). IEEE