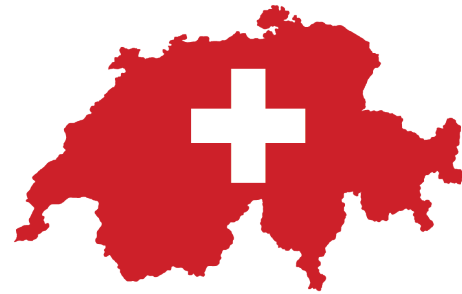# A COLLABORATION WITH



**TAHEREH ZOHDINASAB**

**UNIVERSITÀ DELLA SVIZZERA ITALIANA, SWITZERLAND**



**PAOLO TONELLA**

**UNIVERSITÀ DELLA SVIZZERA ITALIANA, SWITZERLAND**



**ALESSIO GAMBI**

**UNIVERSITY OF PASSAU, GERMANY**



**ACM SIGSOFT INTERNATIONAL SYMPOSIUM ON SOFTWARE TESTING AND ANALYSIS (ISSTA) 2021**

2

# THE AGE OF DEEP LEARNING (DL)
## CAN WE TRUST SOFTWARE THAT LEARNS?

# RESEARCH SUBJECTS



**MNIST**

**HAND-WRITTEN DIGIT CLASSIFICATION**



**BEAMNG**

**LANE-KEEPING FOR AUTONOMOUS CARS**

# TRADITIONAL DL ASSESSMENT



TRAINING SET

VALIDATION SET

TEST SET

DNN CLASSIFIER

ACC = 0.95

PERFORMANCE METRIC

# HOW WILL THE DL BASED SYSTEM BEHAVE FOR DATA BEYOND THE ORIGINAL DATASET?

# AUTOMATED TEST GENERATION FOR DL BASED SYSTEMS



**DLFUZZ**
ADVERSARIAL
IMAGES

**ASFAULT**
ADVERSARIAL
ROAD NETWORKS

**DEEPJANUS**
FRONTIER OF
BEHAVIOURS

**NEED FOR A HUMAN-INTERPRETABLE CHARACTERISATION OF THE MISBEHAVIOURS**

8

# FEATURE MAPS

# FEATURE SELECTION

# FEATURE SELECTION METHODOLOGY

# FEATURE SELECTION METHODOLOGY: OPEN CODING

# FEATURE SELECTION METHODOLOGY: OPEN CODING

## Open Coding

### Pilot Study



Labelled Inputs



| Feature |
|---|
| Boldness |
| Smoothness |
| Discontinuity |

### Assessor 1
Boldness:[0],
Smoothness:[-2],
Discontinuity:[1]

### Assessor 2
Boldness:[2],
Smoothness:[-2],
Discontinuity:[1]

# FEATURE SELECTION METHODOLOGY: OPEN CODING

# FEATURE SELECTION METHODOLOGY: OPEN CODING

# FEATURE SELECTION METHODOLOGY: METRICS IDENTIFICATION

## Metrics Identification

Candidate
Metrics
Design

Candidate Metrics

| Case study | Feature | Metric |
|---|---|---|
| | Boldness | Lum |
| | Smoothness | AvgAng |
| MNIST | Discontinuity | Mov |
| | Rotation | Or |

# FEATURE SELECTION METHODOLOGY: METRICS IDENTIFICATION



| Case study | Feature | Metric | Correlation | P-value |
|---|---|---|---|---|
| MNIST | Boldness | Lum | 0.67 | < 0.002 |
| | Smoothness | AvgAng | 0.05 | 0.241 |
| | Discontinuity | Mov | 0.9 | < 0.002 |
| | Rotation | Or | 0.43 | < 0.002 |

# FEATURE SELECTION METHODOLOGY: METRICS IDENTIFICATION



Metrics Identification

Candidate Metrics Design → Candidate Metrics → Metrics Validation and Correlation

| Case study | Feature | Metric | Correlation | P-value |
|---|---|---|---|---|
| MNIST | Boldness | Lum | 0.67 | < 0.002 |
| | ~~Smoothness~~ | ~~AvgAng~~ | ~~0.05~~ | ~~0.241~~ |
| | Discontinuity | Mov | 0.9 | < 0.002 |
| | Rotation | Or | 0.43 | < 0.002 |

# AUTOMATED TEST GENERATION

# EVOLUTIONARY ALGORITHMS

# EVOLUTIONARY ALGORITHMS

# EVOLUTIONARY ALGORITHMS

# ILLUMINATION SEARCH

# DEEPHYPERION

# DEEPHYPERION

# DEEPHYPERION

# DEEPHYPERION

# DEEPHYPERION

# DEEPHYPERION

Initialize map → Generate Initial Population → Populate Map

Populate Map → Random Selection

Random Selection → Mutation → Evaluation → Update Map → Done?

Done? — Yes → Report

Done? — No → Random Selection

# MODEL-BASED INPUT PERTURBATIONS - DIGITS

BITMAP  SVG MODEL  SVG MODEL  BITMAP



**start_point** = (9.0, 20.85)
**BezierSegment**(
**c1** = (9.0, 20.22),
**c2** = (10.22, 17.30),
**end_point** = (11.70, 14.38)
)

**start_point** = (9.0, 20.85)
**BezierSegment**(
**c1** = (9.0, 20.22),
**c2** = (**8.10**, 17.30),
**end_point** = (11.70, 14.38)
)

# MODEL-BASED INPUT PERTURBATIONS - ROADS



ROAD

CATMULL-ROM MODEL

CATMULL-ROM MODEL

ROAD

# MISBEHAVIOUR DETECTION

# MISBEHAVIOUR DETECTION

MNIST

BEAMNG



# UP TO 10X
# MORE MISBEHAVIOURS

# EXPLORATION EFFECTIVENESS?

# EXPLORATION EFFECTIVENESS



MNIST

BEAMNG

# UP TO 8X
# MAP EXPLORATION

# MISBEHAVIOUR CHARACTERISATION

# MISBEHAVIOUR CHARACTERISATION

**Misbehaviour Probability Map (Run 1)**

**Misbehaviour Probability Map (Run 2)**

**Misbehaviour Probability Map (Run 3)**

**Final Misbehaviour Probability Map**

37

# MISBEHAVIOUR CHARACTERISATION

MNIST

BEAMNG

# USAGE: TEST ADEQUACY CRITERION

**Test Generator 1**

**VS**

**Test Generator 2**

Feature 2 / Feature 1 (axis labels)

# EXTENSION: DEEPHYPERION-CS

# EXTENSION: DEEPHYPERION-CS



$$\text{CONTRIBUTION SCORE (X)} = \frac{\text{\# TIMES A MUTANT OF X IS PLACED IN THE MAP}}{\text{\# TIMES X IS SELECTED}}$$

Initialize map → Generate Initial Population → Populate Map → Contribution-guided Selection → Mutation → Evaluation → Update Map → Done? → Yes → Report; No → Contribution-guided Selection

**Efficient and Effective Feature Space Exploration for Testing Deep Learning Systems**

TAHEREH ZOHDINASAB*, Università della Svizzera Italiana, Switzerland
VINCENZO RICCIO, Università della Svizzera Italiana, Switzerland
ALESSIO GAMBI, University of Passau, Germany
PAOLO TONELLA, Università della Svizzera Italiana, Switzerland

Assessing the quality of Deep Learning (DL) systems is crucial, as they are increasingly adopted in safety-critical domains. Researchers have proposed several input generation techniques for DL systems. While such techniques can expose failures, they do not explain which features of the test inputs influenced the system's (mis-) behaviour. DeepHyperion was the first test generator to overcome this limitation by exploring the DL systems' feature space at large. In this paper, we propose DeepHyperion-CS, a test generator for DL systems which enhances DeepHyperion by promoting the inputs that contributed more to feature space exploration during the previous search iterations. We performed an empirical study involving two different test subjects (i.e., a digit classifier and a lane-keeping system for self-driving cars). Our results proved that the contribution-based guidance implemented within DeepHyperion-CS outperforms 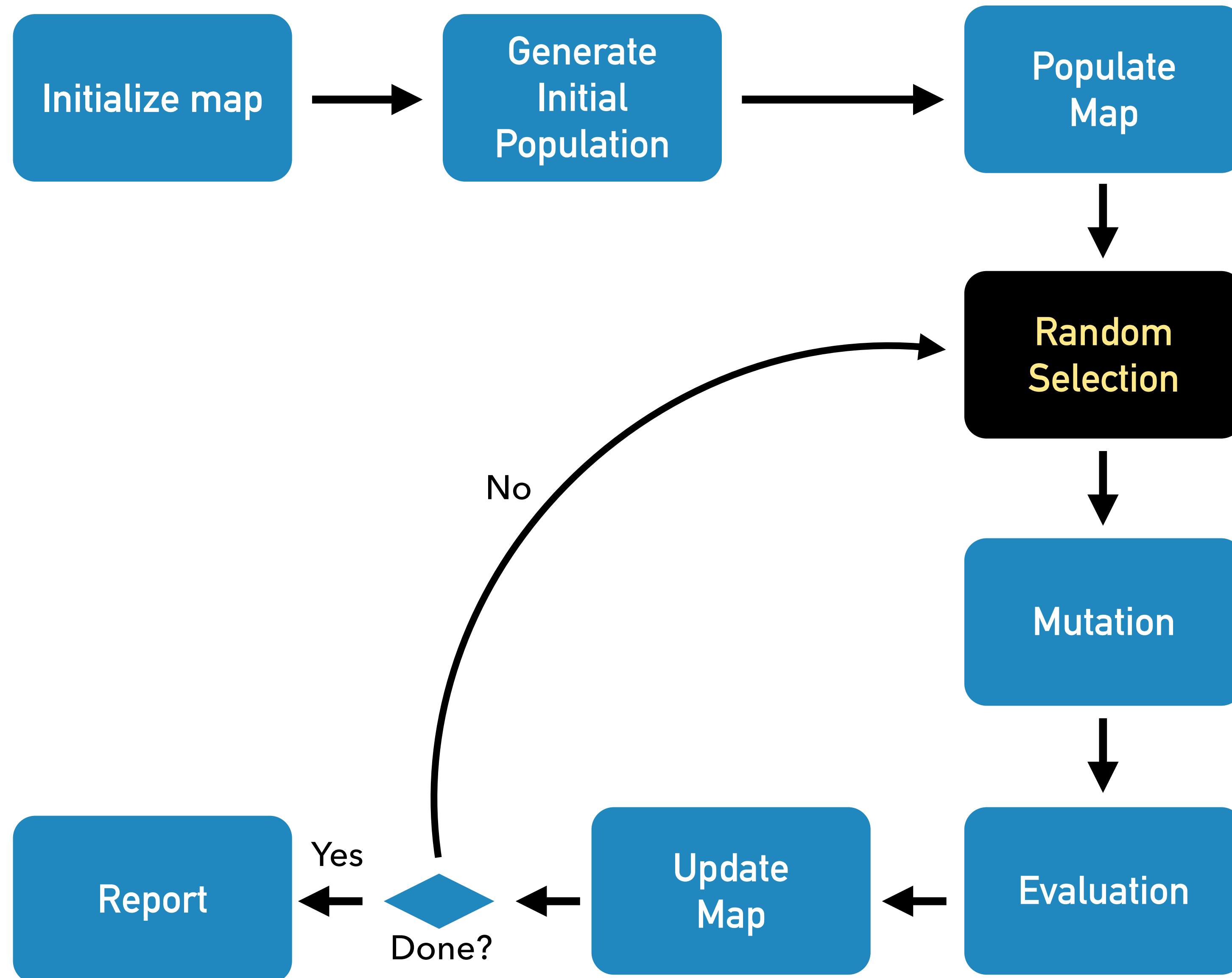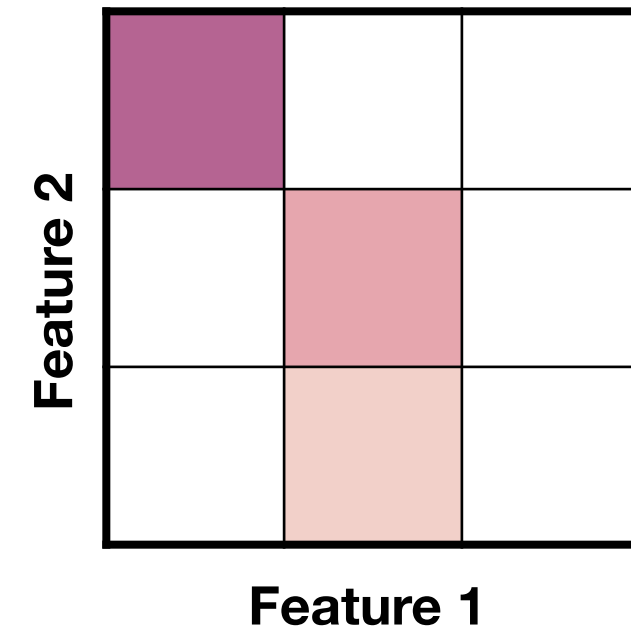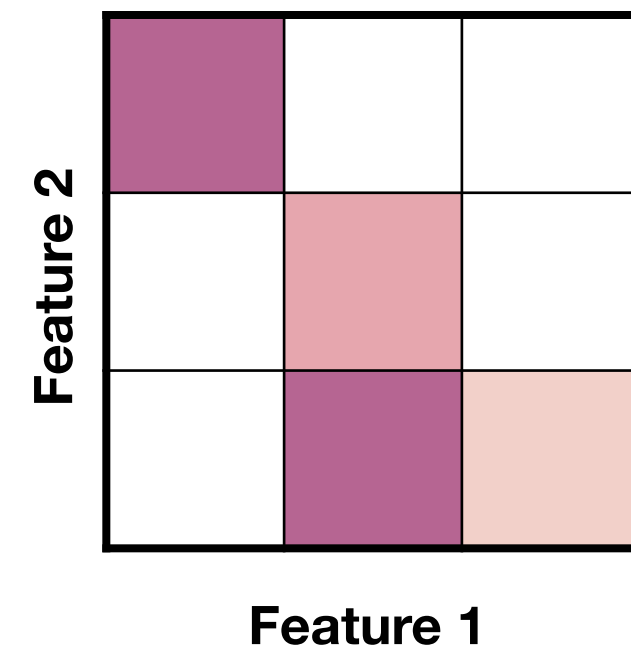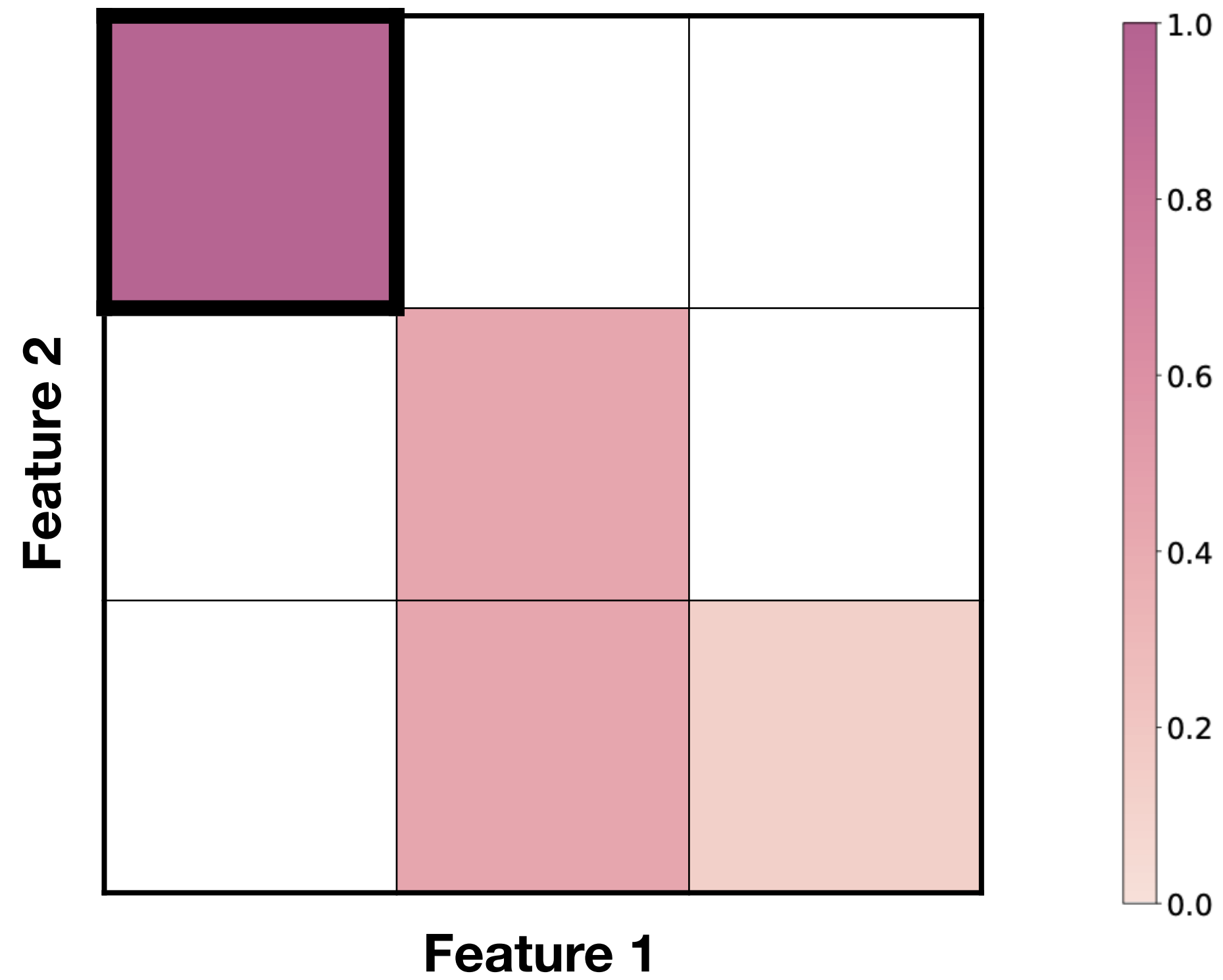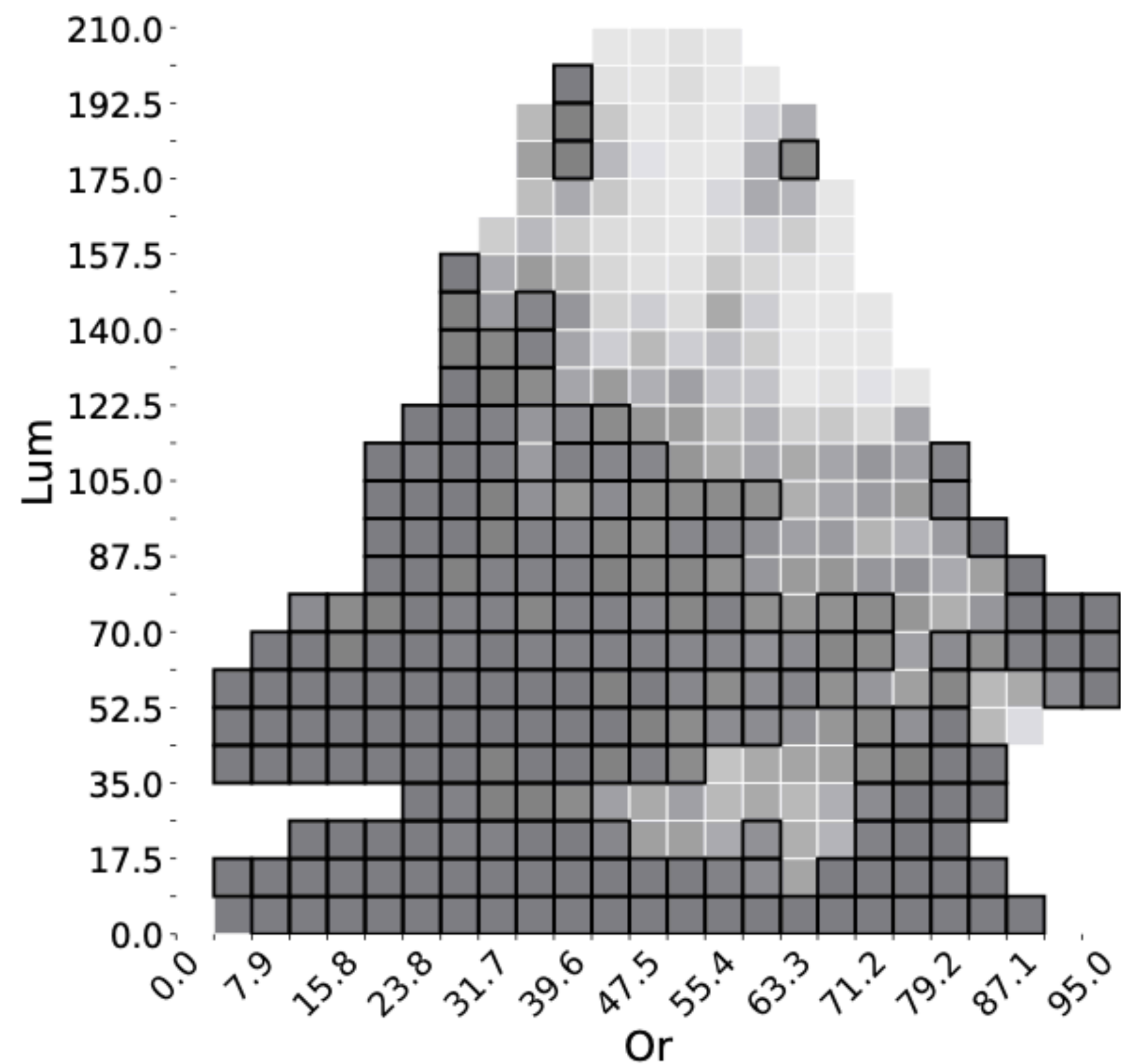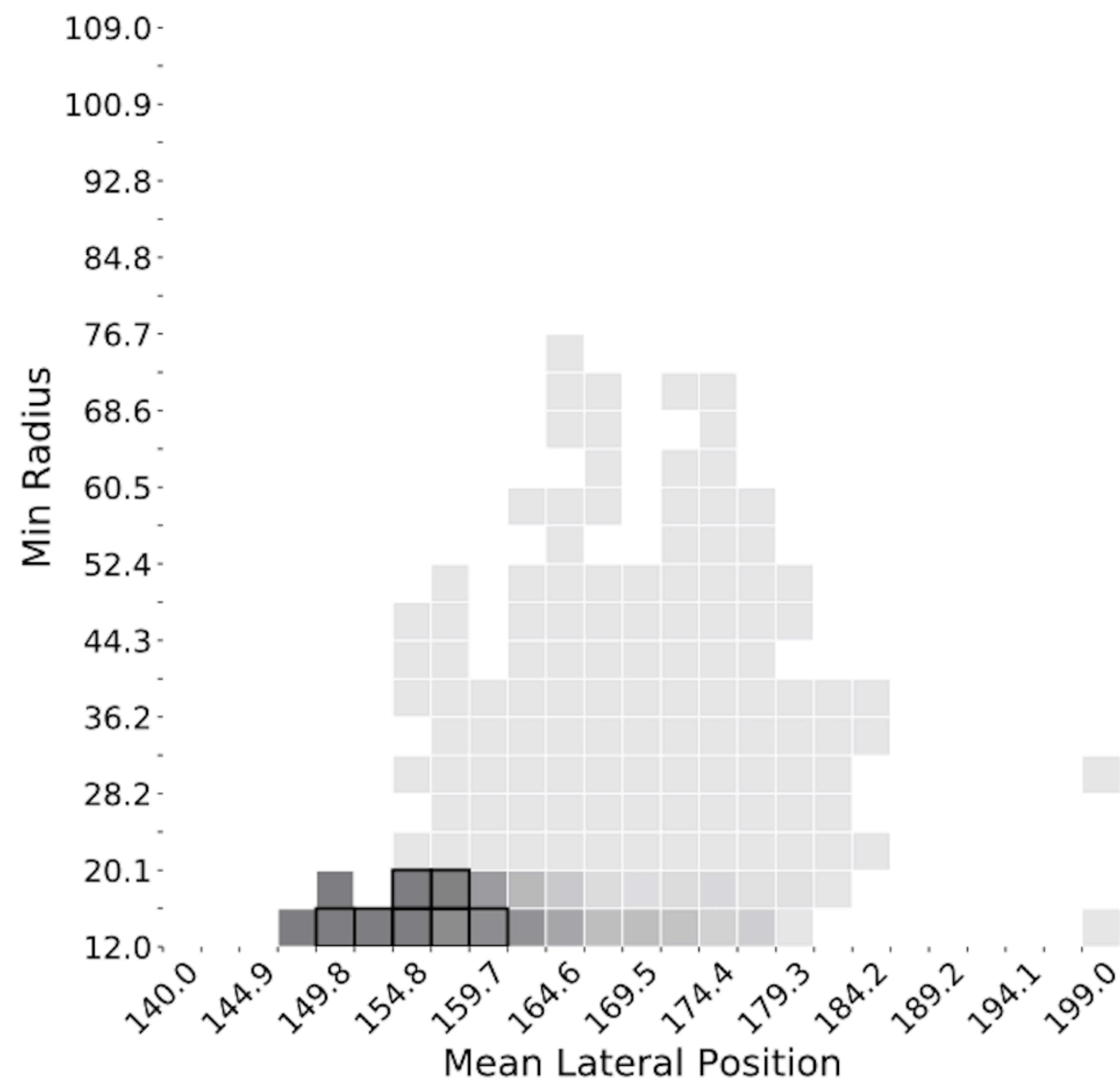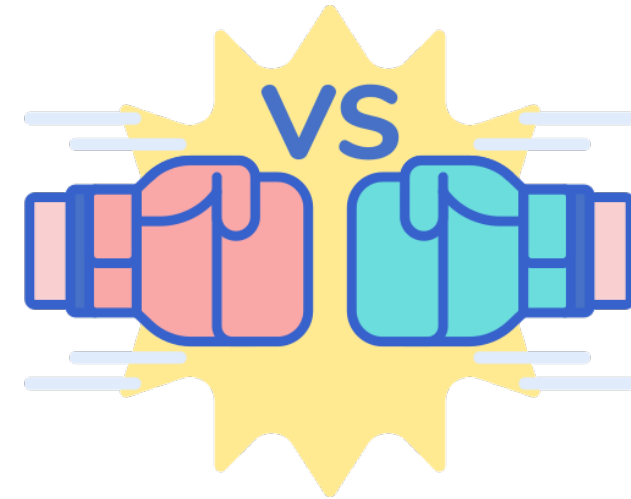state-of-the-art tools and significantly improves the efficiency and the effectiveness of DeepHyperion. DeepHyperion-CS exposed significantly more misbehaviours for 5 out of 6 feature combinations and was up to 65% more efficient than DeepHyperion in finding misbehaviour-inducing inputs and exploring the feature space. DeepHyperion-CS was useful for expanding the datasets used to train the DL systems, populating up to 200% more feature map cells than the original training set.

CCS Concepts: • **Software and its engineering → Software testing and debugging**.

Additional Key Words and Phrases: software testing, deep learning, search based software engineering, self-driving cars

**ACM Reference Format:**
Tahereh Zohdinasab, Vincenzo Riccio, Alessio Gambi, and Paolo Tonella. 2021. Efficient and Effective Feature Space Exploration for Testing Deep Learning Systems. 1, 1 (August 2021), 38 pages. https://doi.org/10.1145/1122445.1122456

**1 INTRODUCTION**

Using Deep Learning (DL) has become widespread for modern software systems that must process complex inputs and timely solve challenging tasks. For example, image classifiers [31, 62] can analyse images to diagnose diseases, while intelligent driving agents use sensor information (e.g., from cameras and LiDARs) to drive vehicles [10]. Since DL systems are applied also in safety-critical domains, ensuring their dependability is *literally* vital.

Unlike traditional software, DL systems' behaviour is not explicitly coded, being instead indirectly learned from training examples [49]. This fundamental difference of DL systems from traditional software has profound implications on how their quality is assessed.

Authors' addresses: Tahereh Zohdinasab, tahereh.zohdinasab@usi.ch, Università della Svizzera Italiana, Lugano, Switzerland, 6900; Vincenzo Riccio, Università della Svizzera Italiana, Lugano, Switzerland, vincenzo.riccio@usi.ch; Alessio Gambi, University of Passau, Passau, Germany, alessio.gambi@uni-passau.de; Paolo Tonella, Università della Svizzera Italiana, Lugano, Switzerland, paolo.tonella@usi.ch.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2018 Association for Computing Machinery.
Manuscript submitted to ACM

Manuscript submitted to ACM    1

**ACM TRANSACTIONS ON SOFTWARE ENGINEERING AND METHODOLOGY (TOSEM) 2022**

41

42

44

45