

# Interpretable Generative Stress Testing

## A Framework to Audit Decision Boundaries through Synthetic Data Generation



Me!

Inês Gomes<sup>1,2,4</sup>, Luís F. Teixeira<sup>1,3</sup>, Jan N. van Rijn<sup>4</sup> and Thomas Bäck<sup>4</sup>

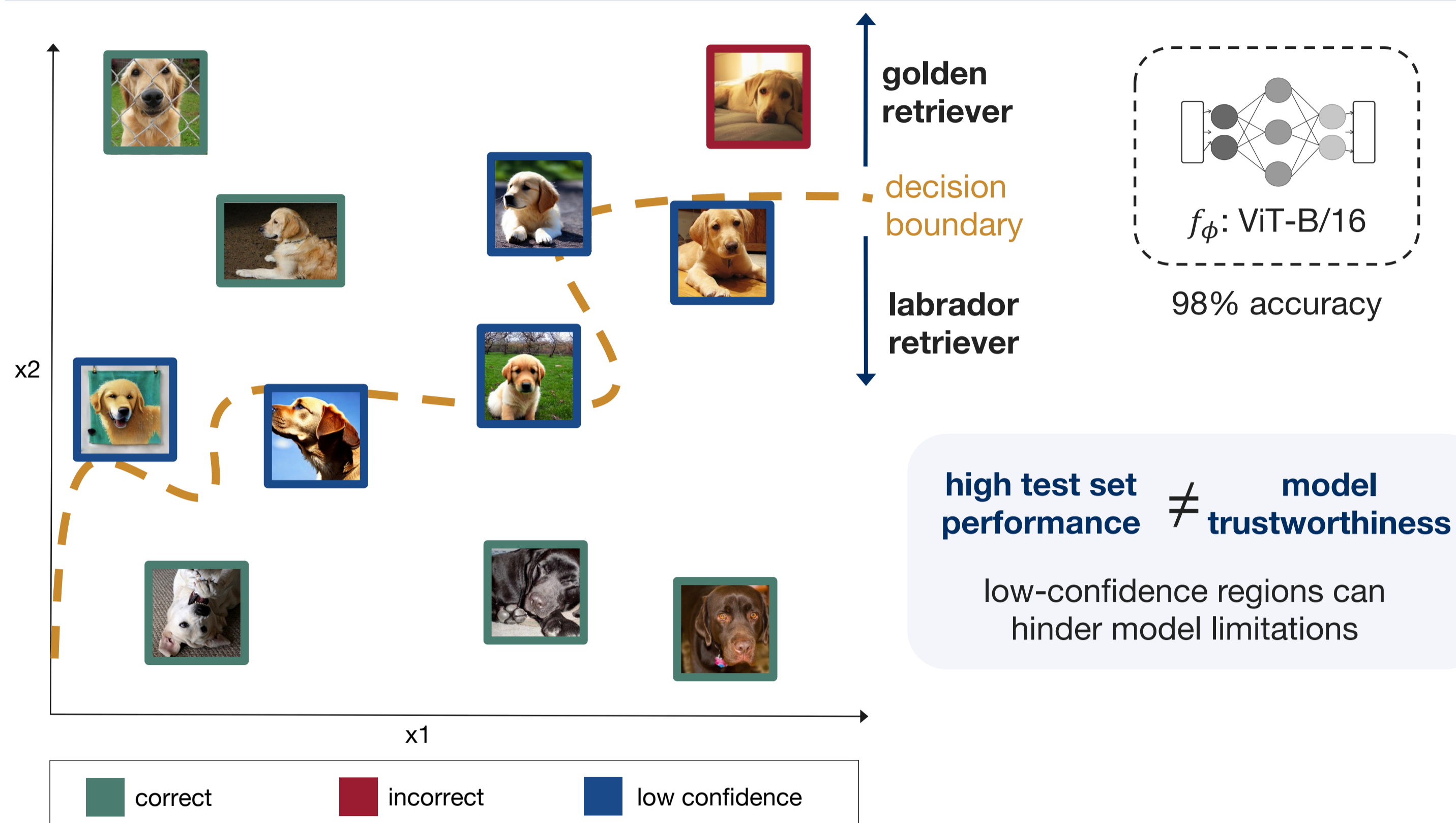
<sup>1</sup>Faculty of Engineering, University of Porto, Portugal

<sup>2</sup>Artificial Intelligence and Computer Science Laboratory, Portugal

<sup>3</sup>INESC TEC, Portugal

<sup>4</sup>Leiden Institute of Advanced Computer Science (LIACS), Leiden University, The Netherlands

### MOTIVATION

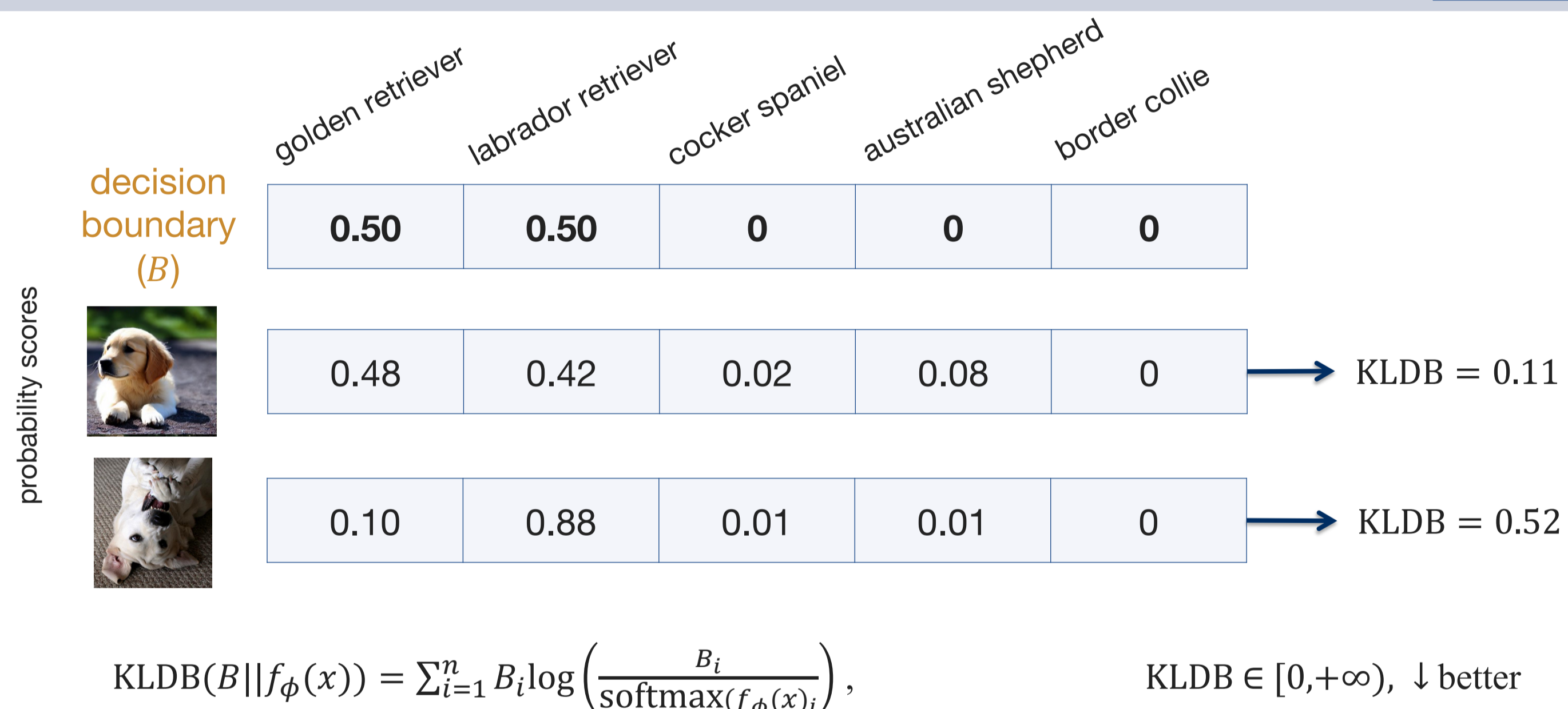


### GOAL

Develop a framework to expose model limitations associated with low-confidence regions in image classifiers

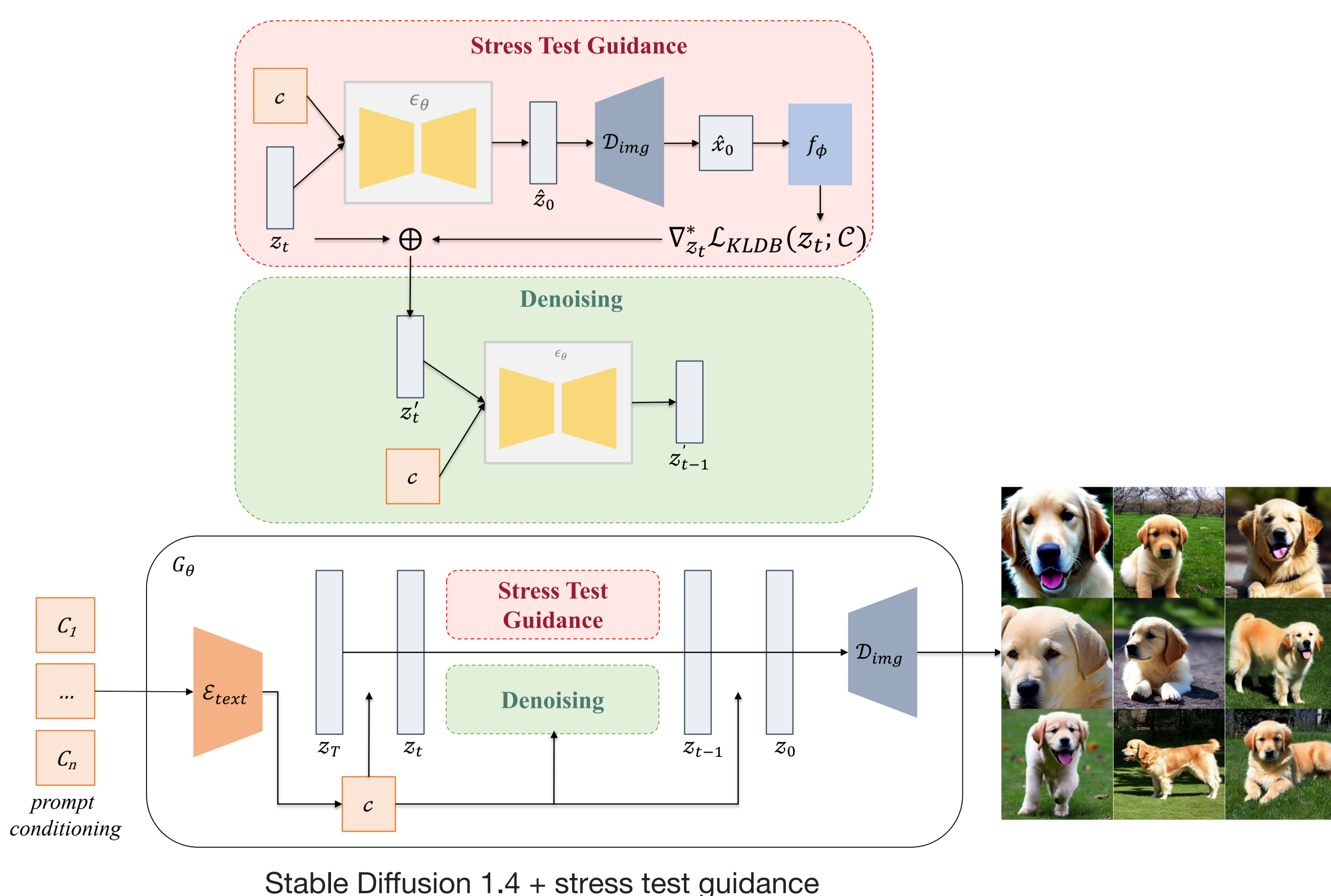
### 1.1. ARBITRARY DECISION BOUNDARIES

submitted under review

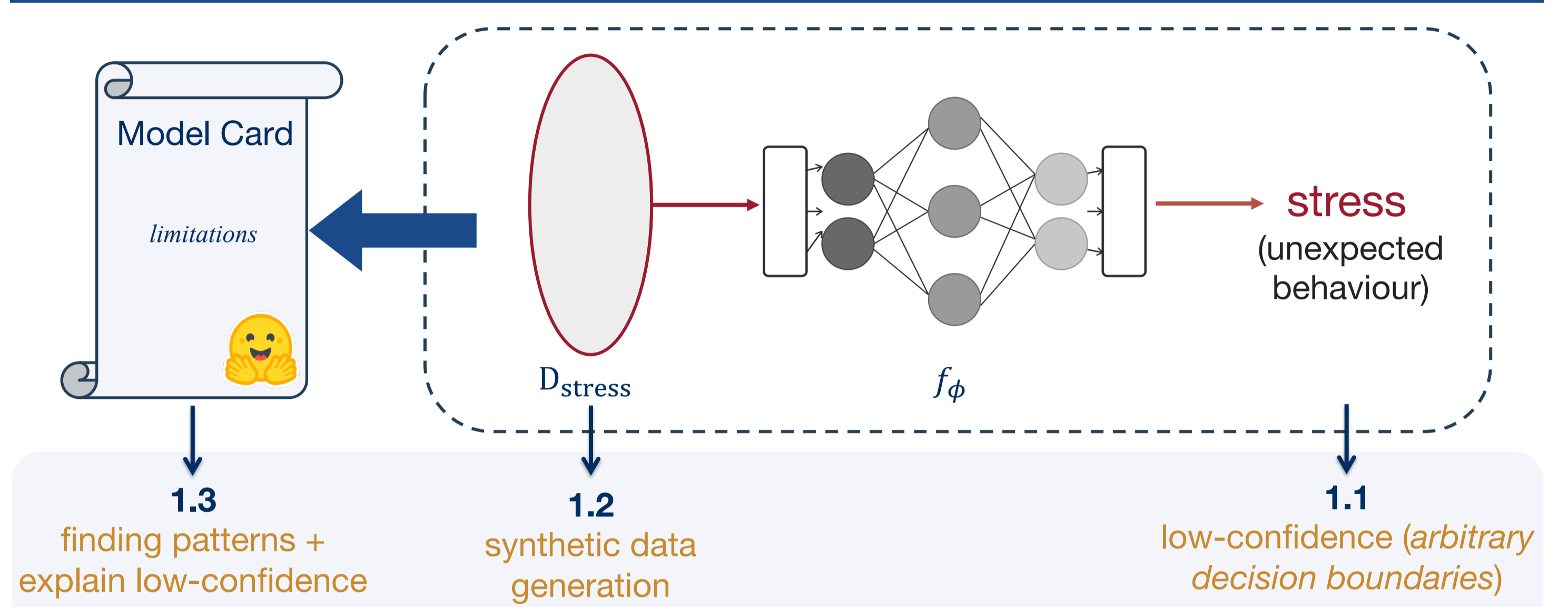


### 1.2. SYNTHETIC DATA GENERATION

submitted under review

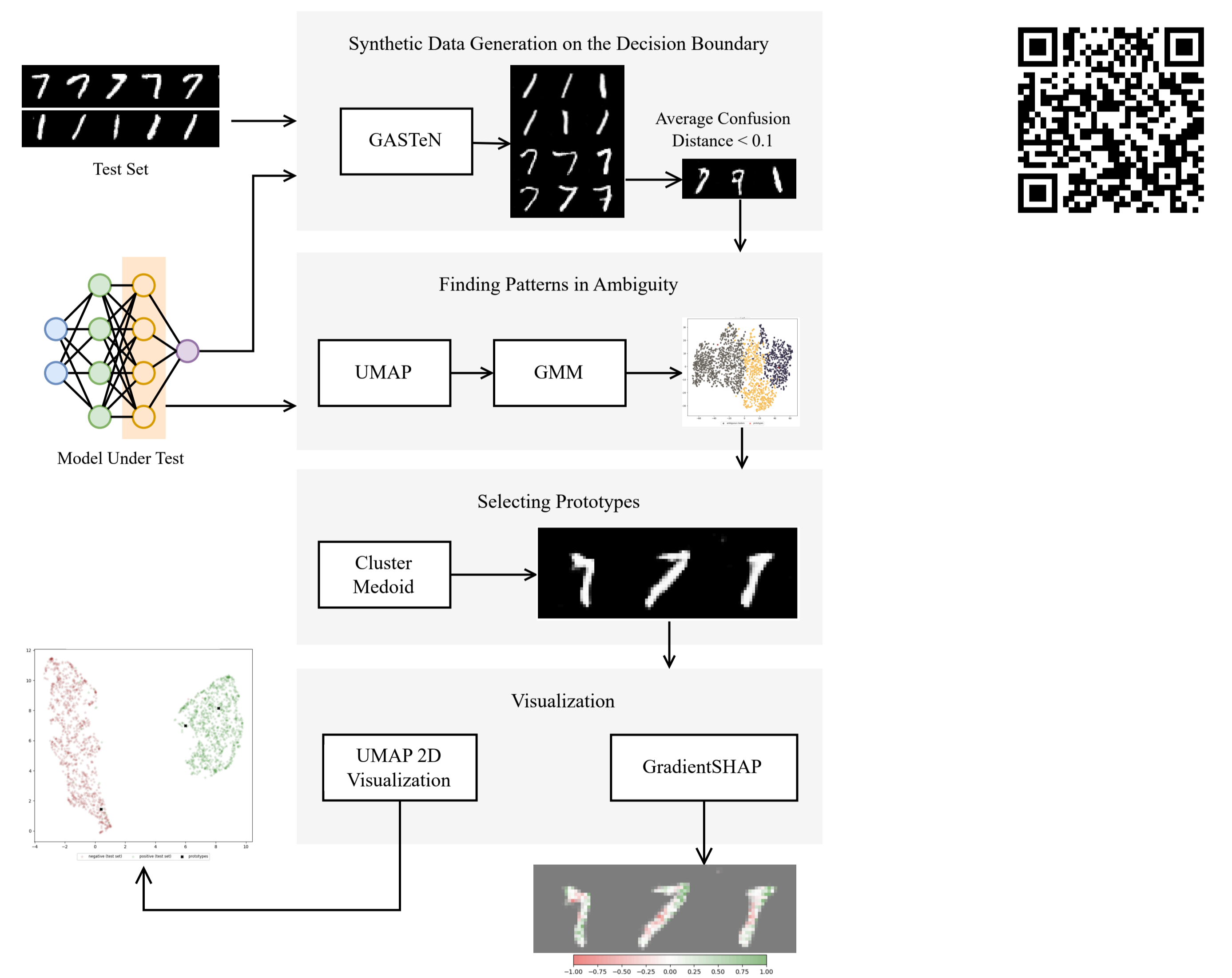


### STRESS TESTING FRAMEWORK



### 1.3. FINDING PATTERNS

published (CVPRW)



### 1.3. EXPLAINING LOW-CONFIDENCE

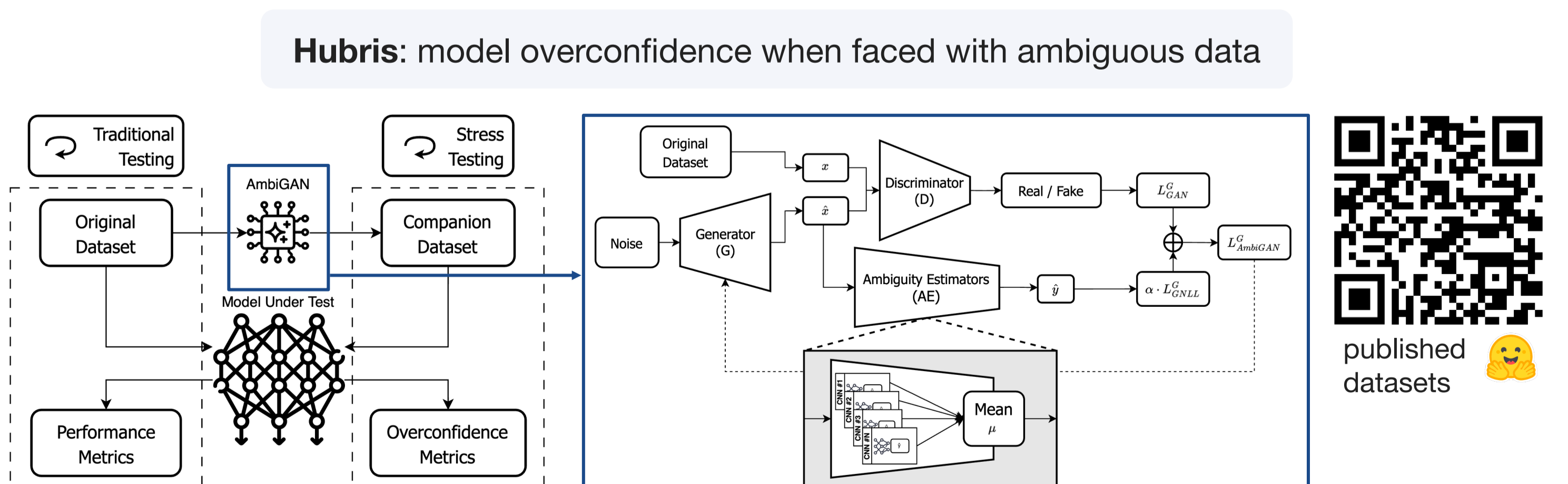
WIP (WACV)

- subgroup discovery to identify groups of images driving low-confidence OR concept-based XAI to identify visual features driving low-confidence
- ✓ Global textual explanations easily integrated in model cards
  - ✓ Pinpoints specific visual concepts (local and global explanations)
  - ⚠ Images may belong to multiple subgroups
  - ⚠ Low confidence may come from missing features
  - ⚠ Descriptions depend on VLM quality
  - ⚠ Mixed-class samples hard to detect

Open question: how to characterize the decision boundary? (OOD/hardness/ambiguity/noise/...)

### HUBRIS BENCHMARKING

submitted under review



### ACKNOWLEDGMENTS

