

# ReMIX: Resilience for ML Ensembles using XAI at Inference against Faulty Training Data

**Abraham Chan,**

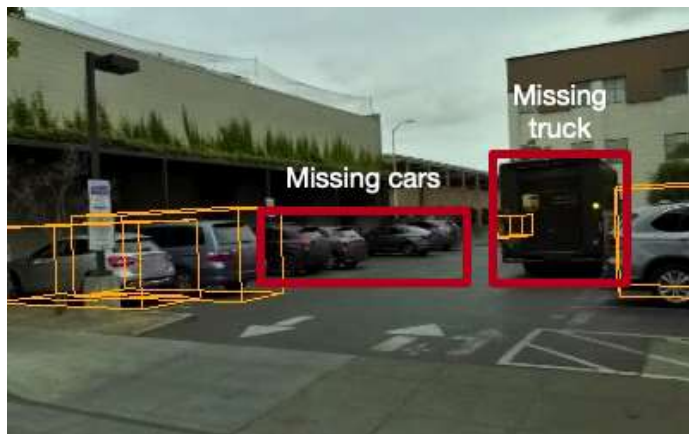
Arpan Gujarati, Karthik Pattabiraman, Sathish Gopalakrishnan



THE UNIVERSITY  
OF BRITISH COLUMBIA

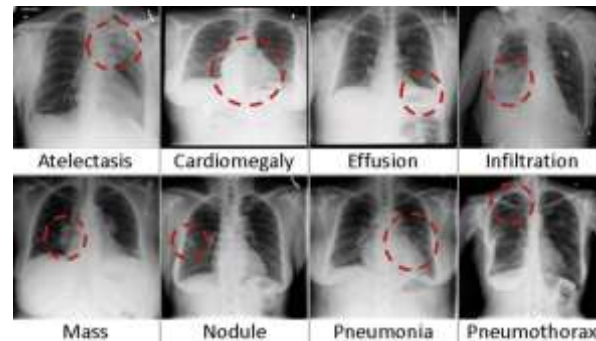
# Training Data Faults in Practice

70% of Lyft dataset missing, mislabelled [Kang et al, 2022]



**Autonomous Vehicles**

20% of ChestX-ray mislabelled [Tang et al, 2021]



**Healthcare**

# Training Data Faults

Image

Label



~~Speed 80km/h~~  
Stop sign



~~Stop sign~~  
Speed 80km/h



No entry

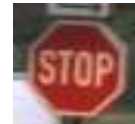
**Mislabeling**

Image

Label



Speed 80km/h



Stop sign



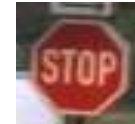
**Removal**

Image

Label



Speed 80km/h



Stop sign



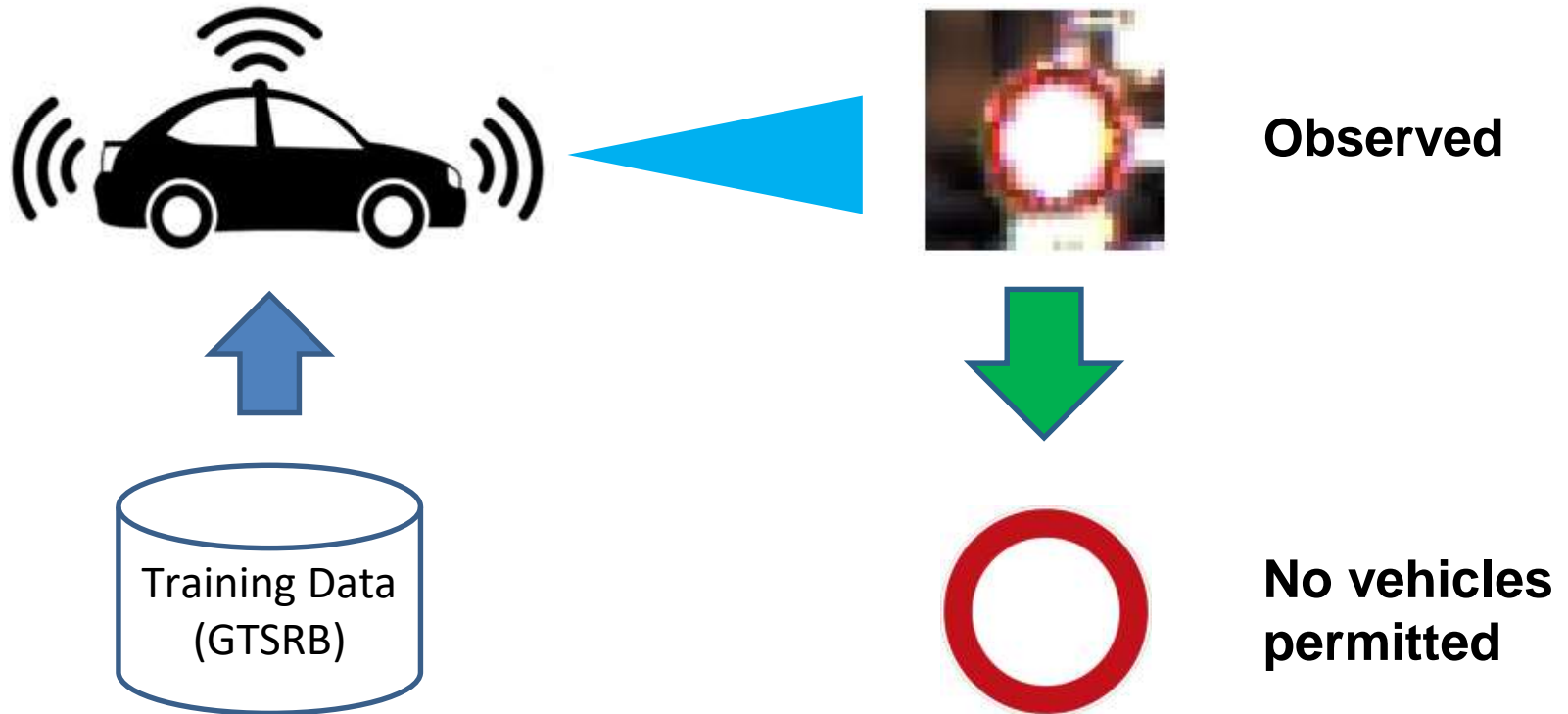
No entry



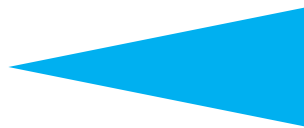
No entry

**Repetition**

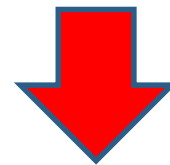
# Autonomous Vehicle Example



# Random Mislabelling

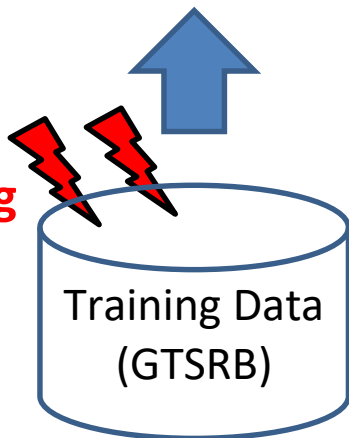


**Observed**



**100 km/hr  
speed limit**

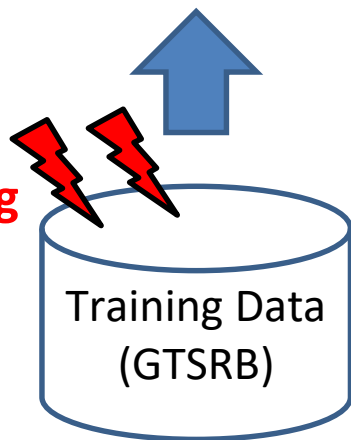
**30%  
Random  
Mislabelling**



# Resilience against Faulty Training Data



**30%  
Random  
Mislabelling**



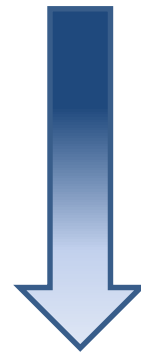
**Resilience**

# How to mitigate training data faults with minimal human effort?



1. Label Correction
2. Knowledge Distillation
3. Robust Loss
4. Label Smoothing
5. Ensembles

More Practitioner Effort



Less Practitioner Effort

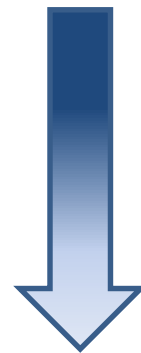
**Our Prior Work:** The Fault in Our Data Stars: Studying Mitigation Techniques against Faulty Training Data in ML Applications [Chan, DSN'22]

# How to mitigate training data faults with minimal human effort?



1. Label Correction
2. Knowledge Distillation
3. Robust Loss
4. Label Smoothing
5. **Ensembles**

More Practitioner Effort



Less Practitioner Effort

**Our Prior Work:** The Fault in Our Data Stars: Studying Mitigation Techniques against Faulty Training Data in ML Applications [**Chan, DSN'22**]



# How to mitigate training data faults with minimal human effort?



1. Label Correction

2. Knowledge Distillation

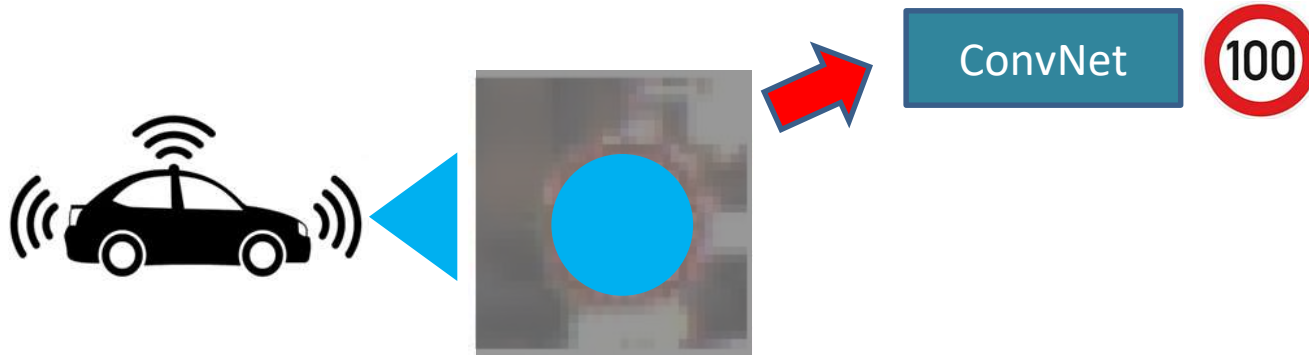
**Our Solution: Building Resilient Ensembles**

4. Label Smoothing

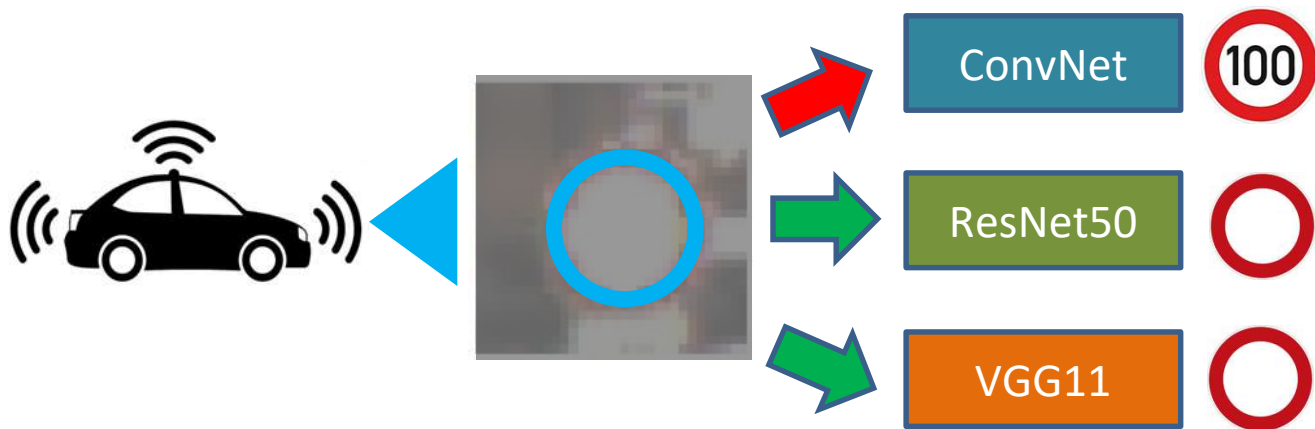
5. Ensembles

**Our Prior Work:** The Fault in Our Data Stars: Studying Mitigation Techniques against Faulty Training Data in ML Applications [Chan, DSN'22]

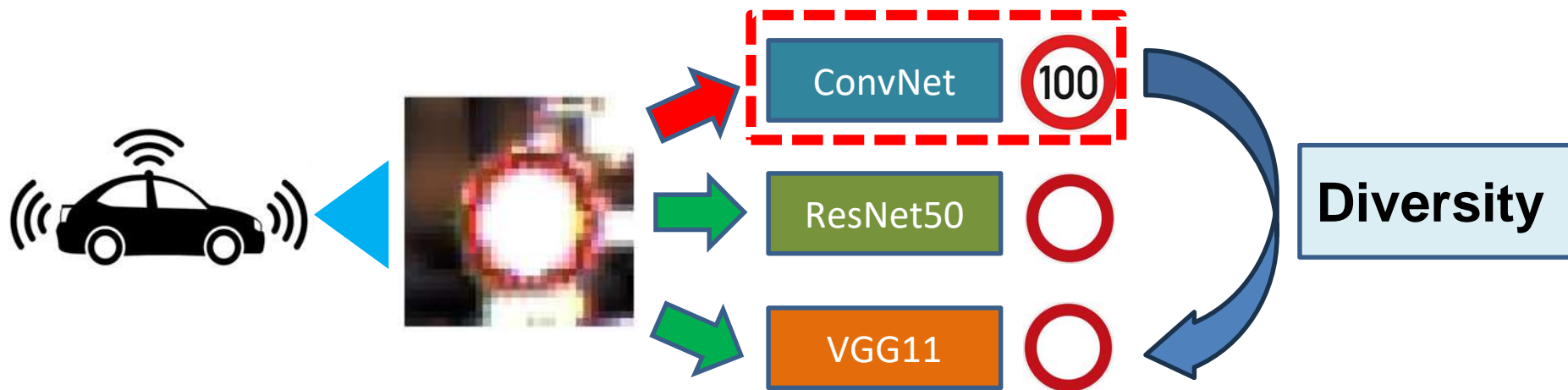
# Resilient Ensembles



# Resilient Ensembles

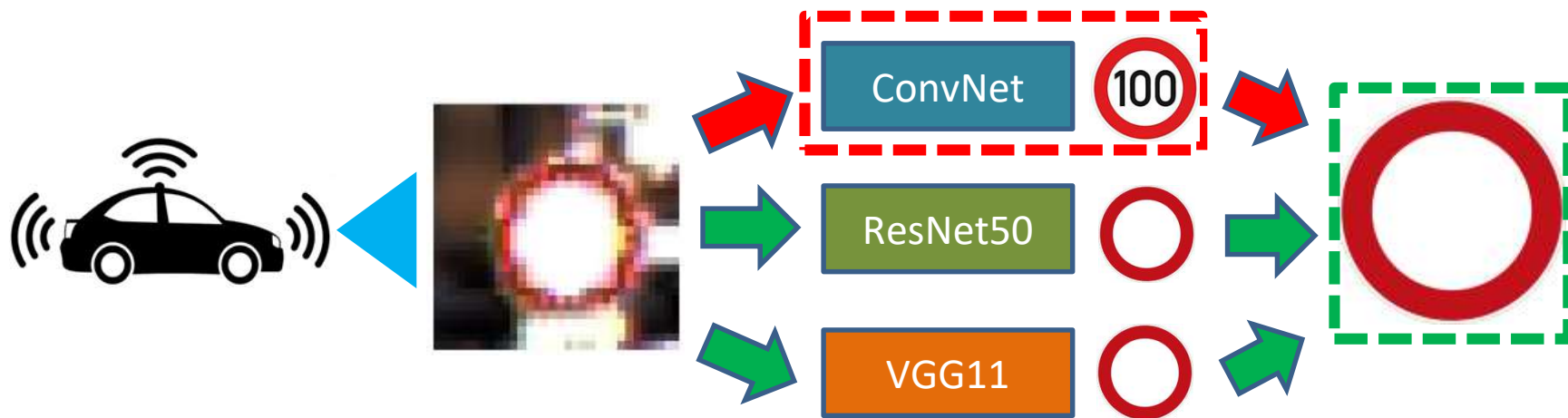


# Resilient Ensembles



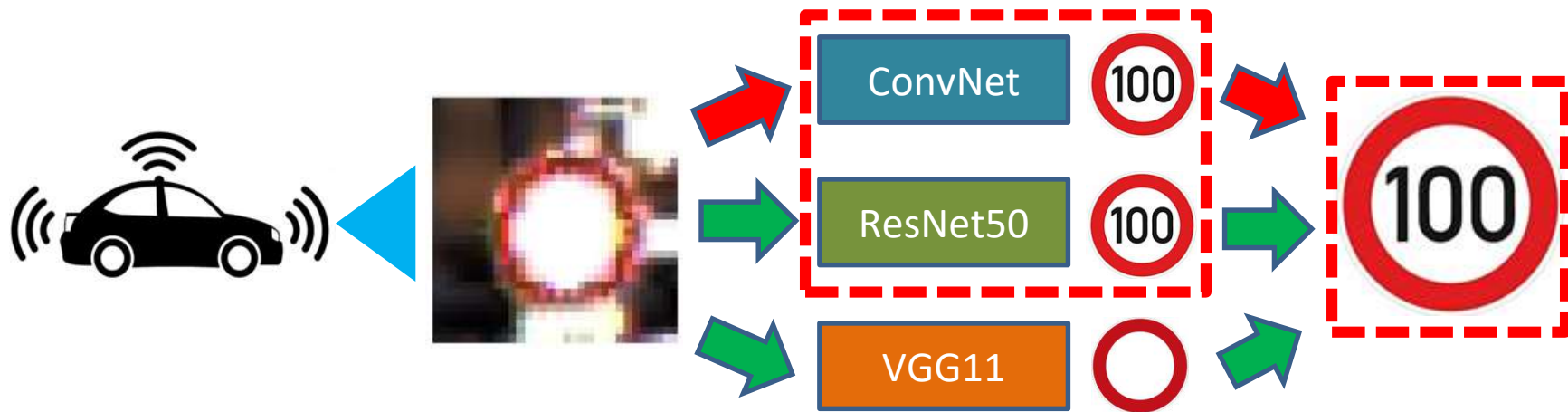
**Our Prior Work:** Understanding the Resilience of Neural Network Ensembles against Faulty Training Data [Chan, QRS'21]

# Resilient Ensembles



**Our Prior Work:** Understanding the Resilience of Neural Network Ensembles against Faulty Training Data [Chan, QRS'21]

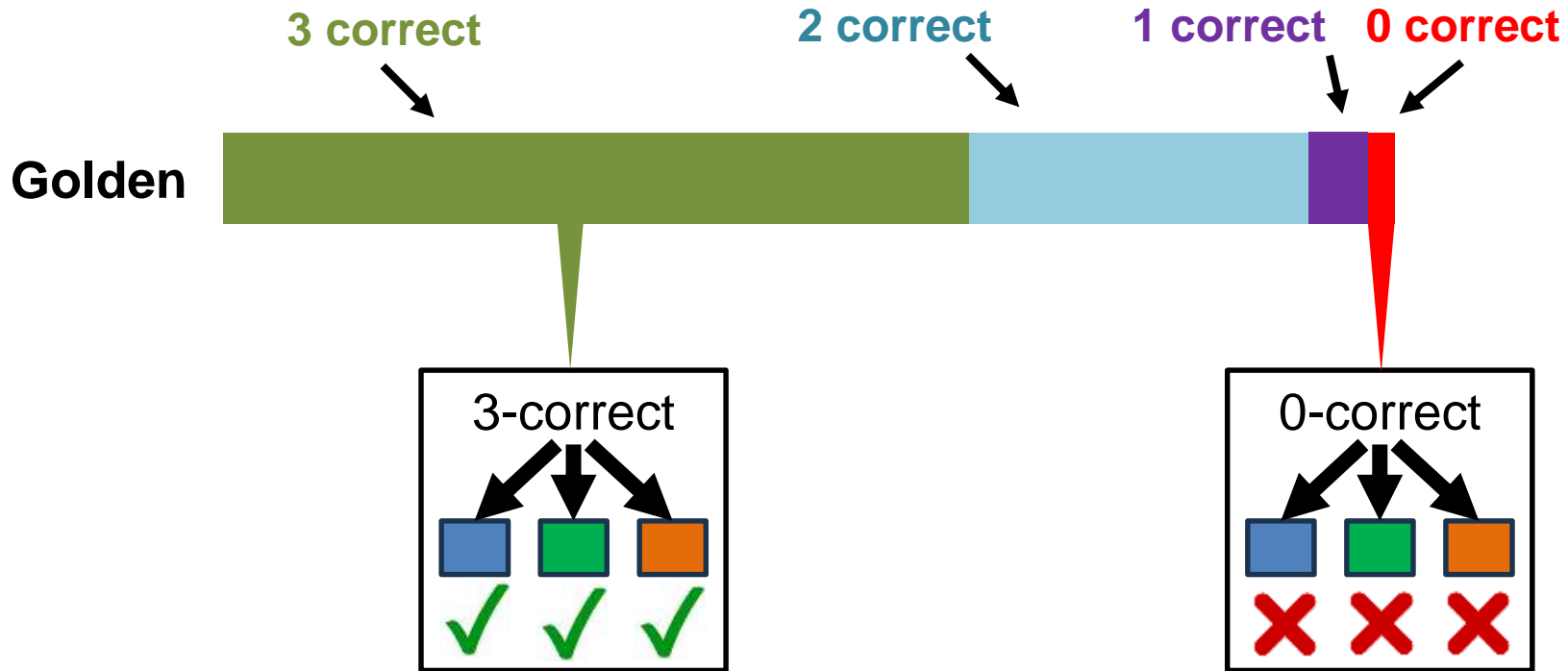
# When Ensembles Misclassify?



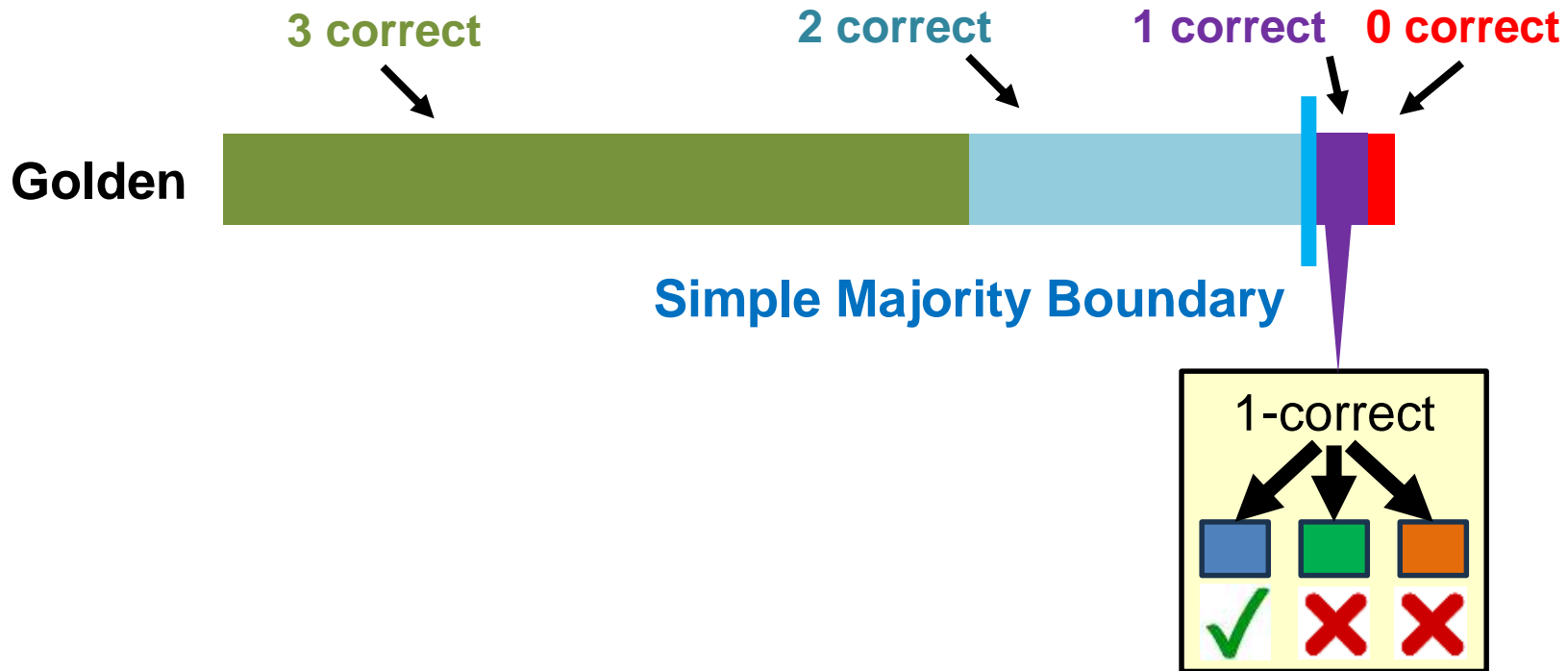
**This paper's contribution:**

**ReMIX** reduces ensemble misclassifications

# # Correct Models in Ensemble



# # Correct Models in Ensemble





# # Correct Models in Ensemble



⚡⚡  
30% Mislabelling



# # Correct Models in Ensemble

Golden

1 correct

4x increase under  
training data faults

30% Mislabelling

# # Correct Models in Ensemble

Golden



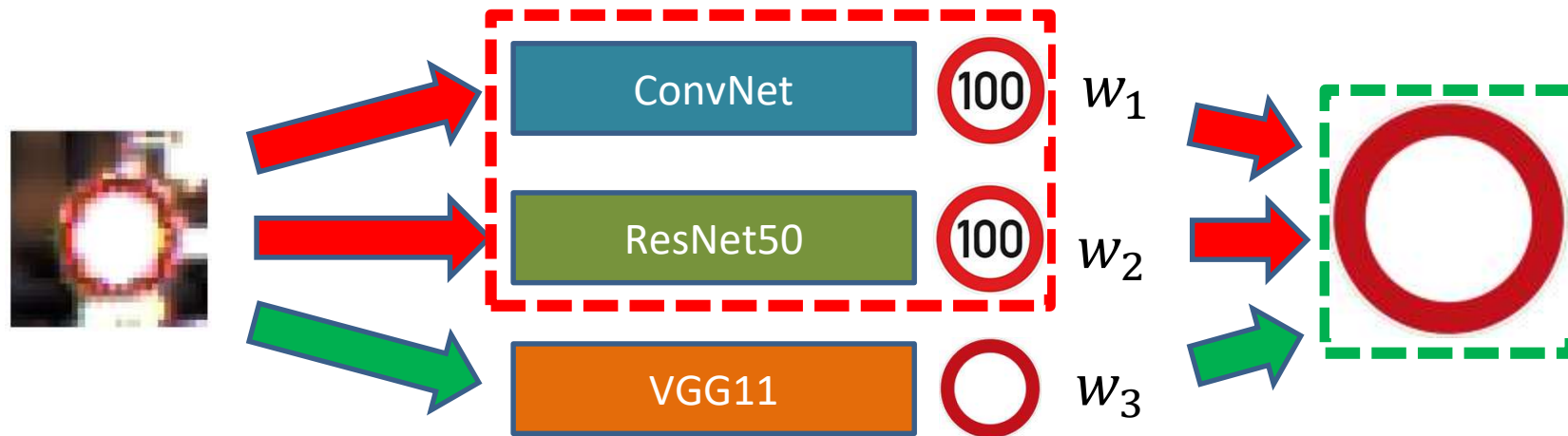
1 correct

**Observation:** Simple majority voting is vulnerable under faulty training data

30% Mislabelling



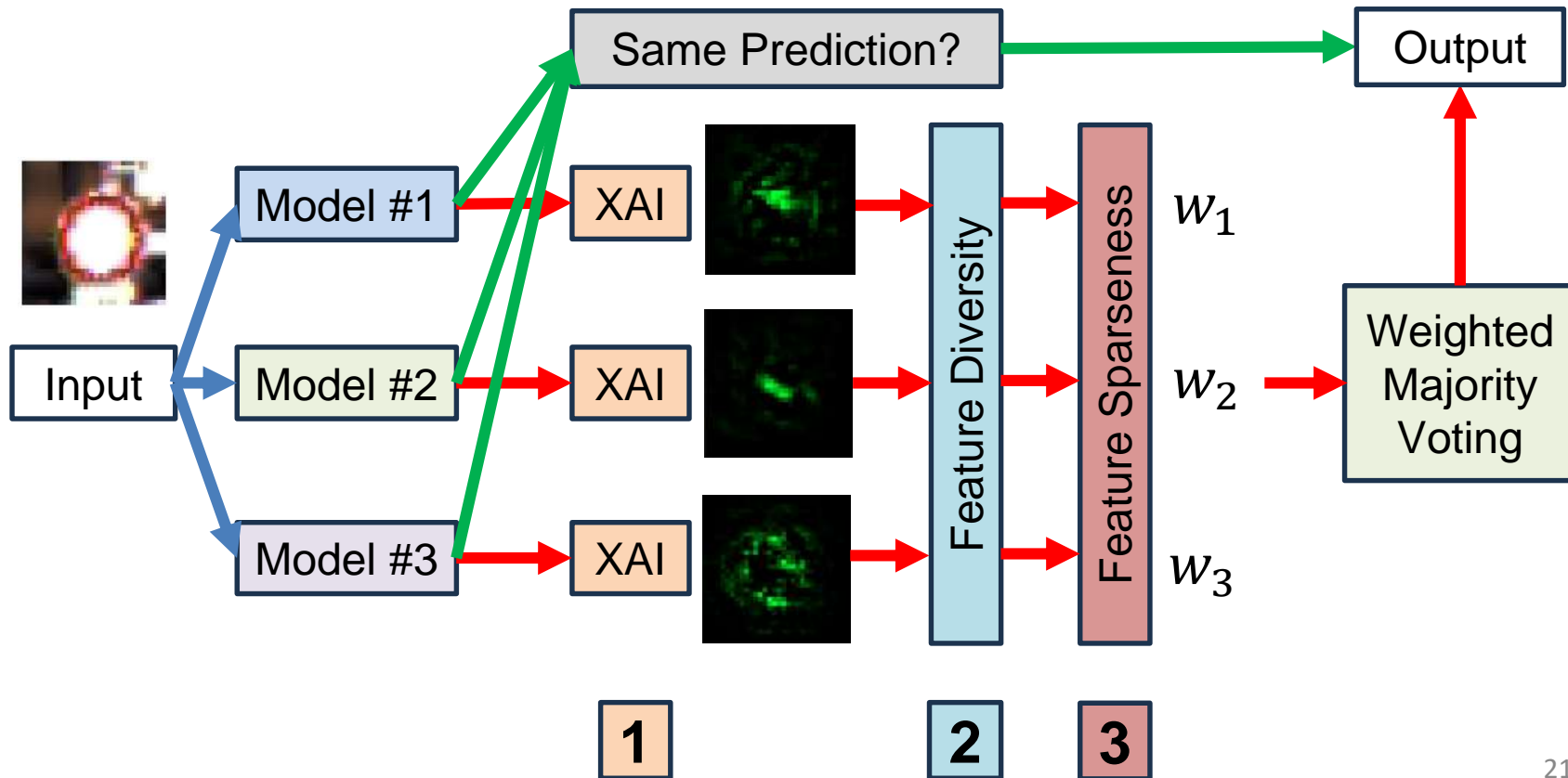
# Dynamically Weighted Ensembles



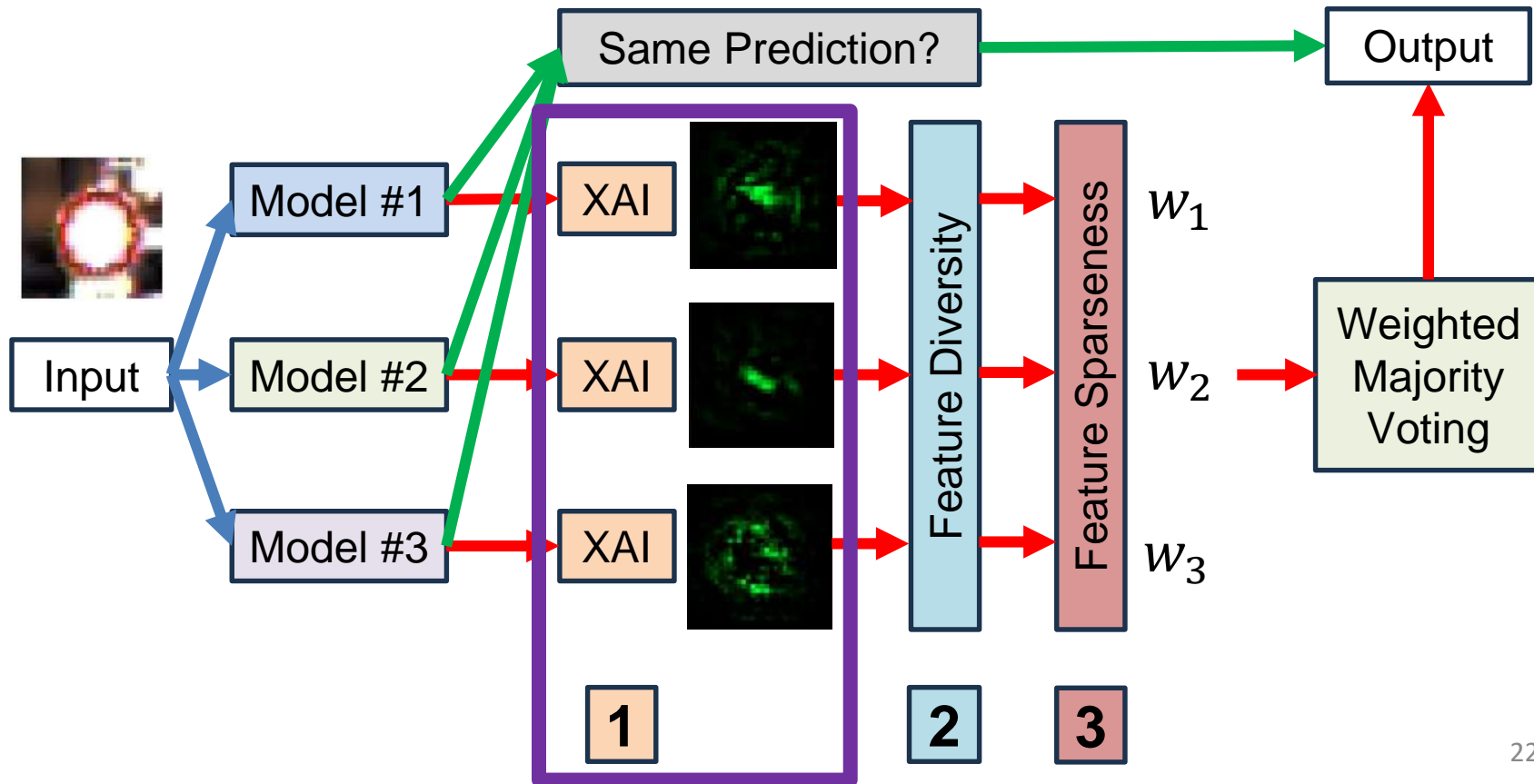
How to determine weights?  
**ReMIX** uses Feature Space  
Diversity!

$$w_1, w_2 < w_3$$

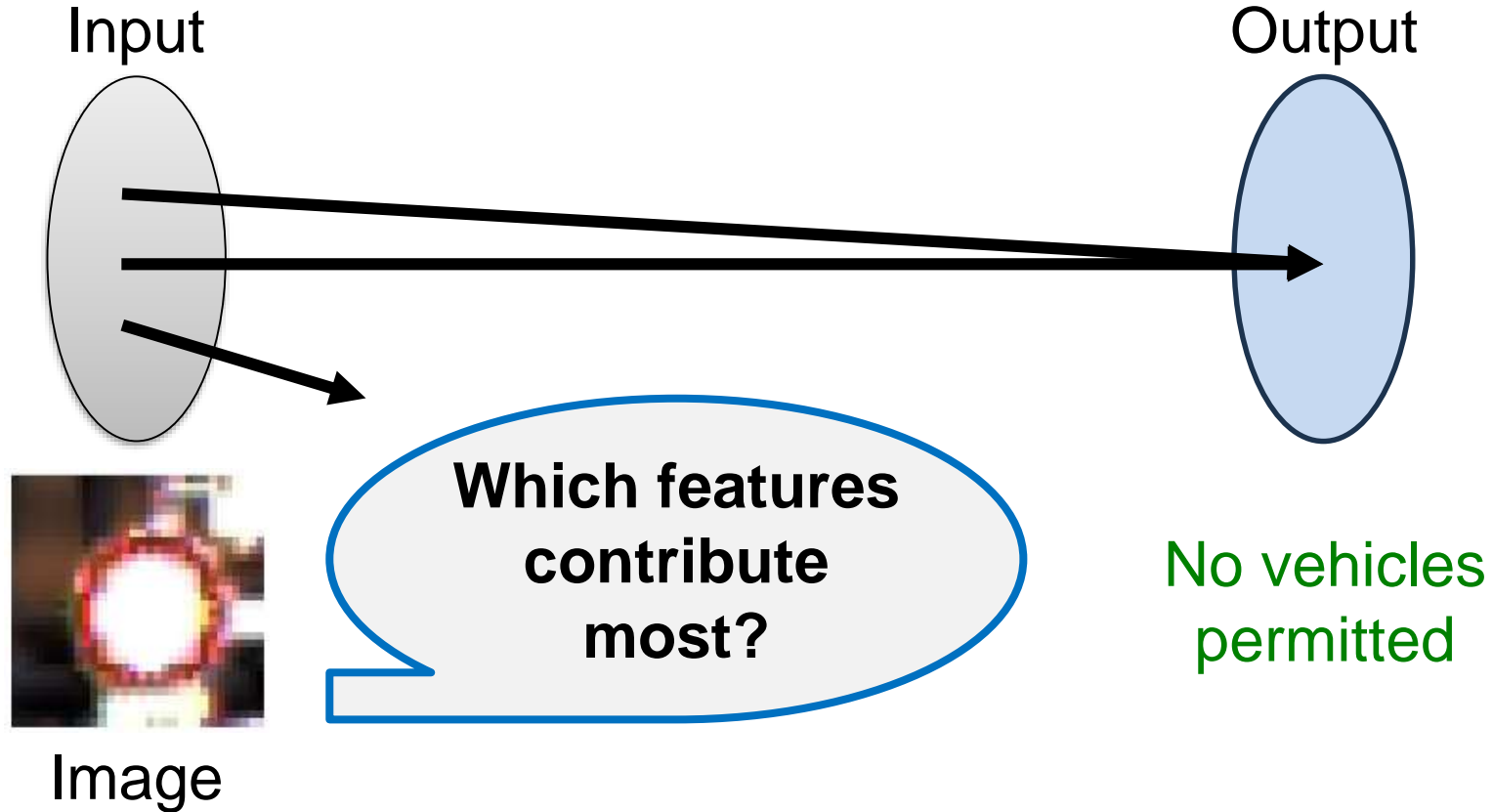
# ReMIX Workflow



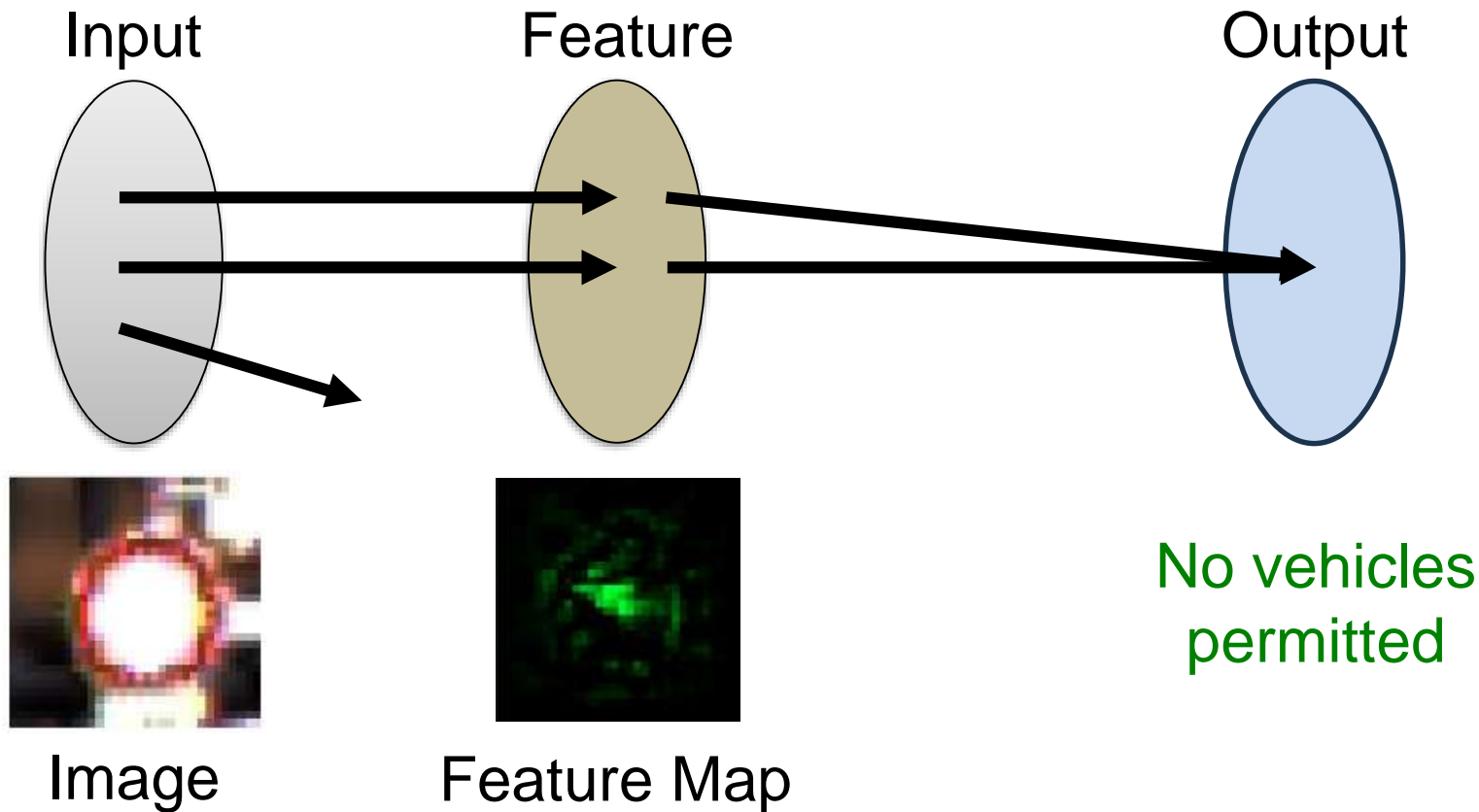
# Step 1 – Explainable AI (XAI)



# Input - Output Space



# Feature Space





# Post-Hoc Local Explainable AI (XAI)

Post-Hoc XAI



ML Model

No vehicles  
permitted

# Post-Hoc Local Explainable AI (XAI)

1. Smooth Gradients (SG)
2. Integrated Gradients (IG)
3. SHAP
4. LIME
5. Counterfactual Explanations

Faithful?

Robust?

Efficient?



ML Model

No vehicles  
permitted

# Post-Hoc Local Explainable AI (XAI)



1. **Smooth Gradients (SG)**
2. Integrated Gradients (IG)
3. SHAP
4. LIME
5. Counterfactual Explanations

Faithful?

Robust?

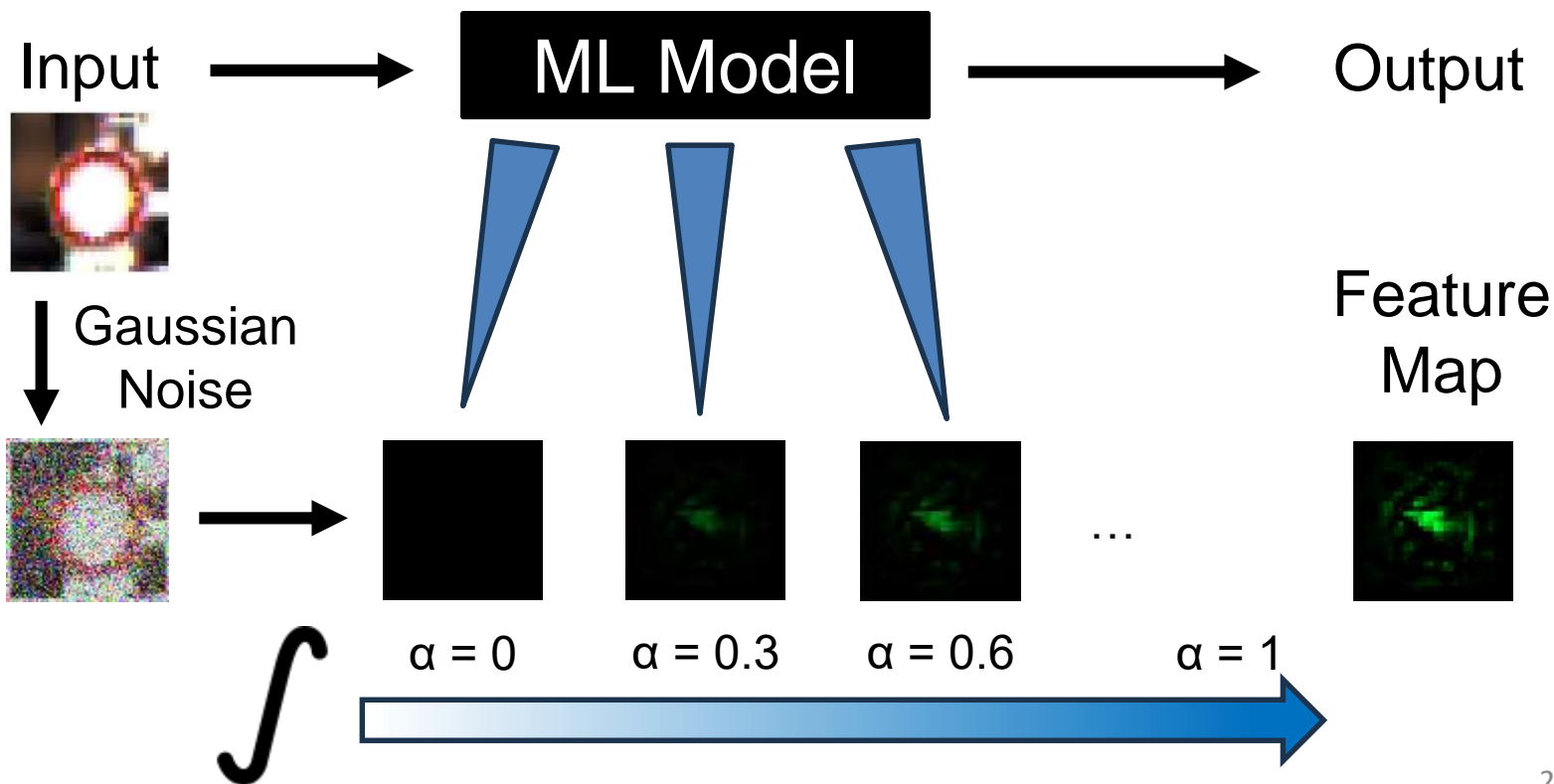
Efficient?



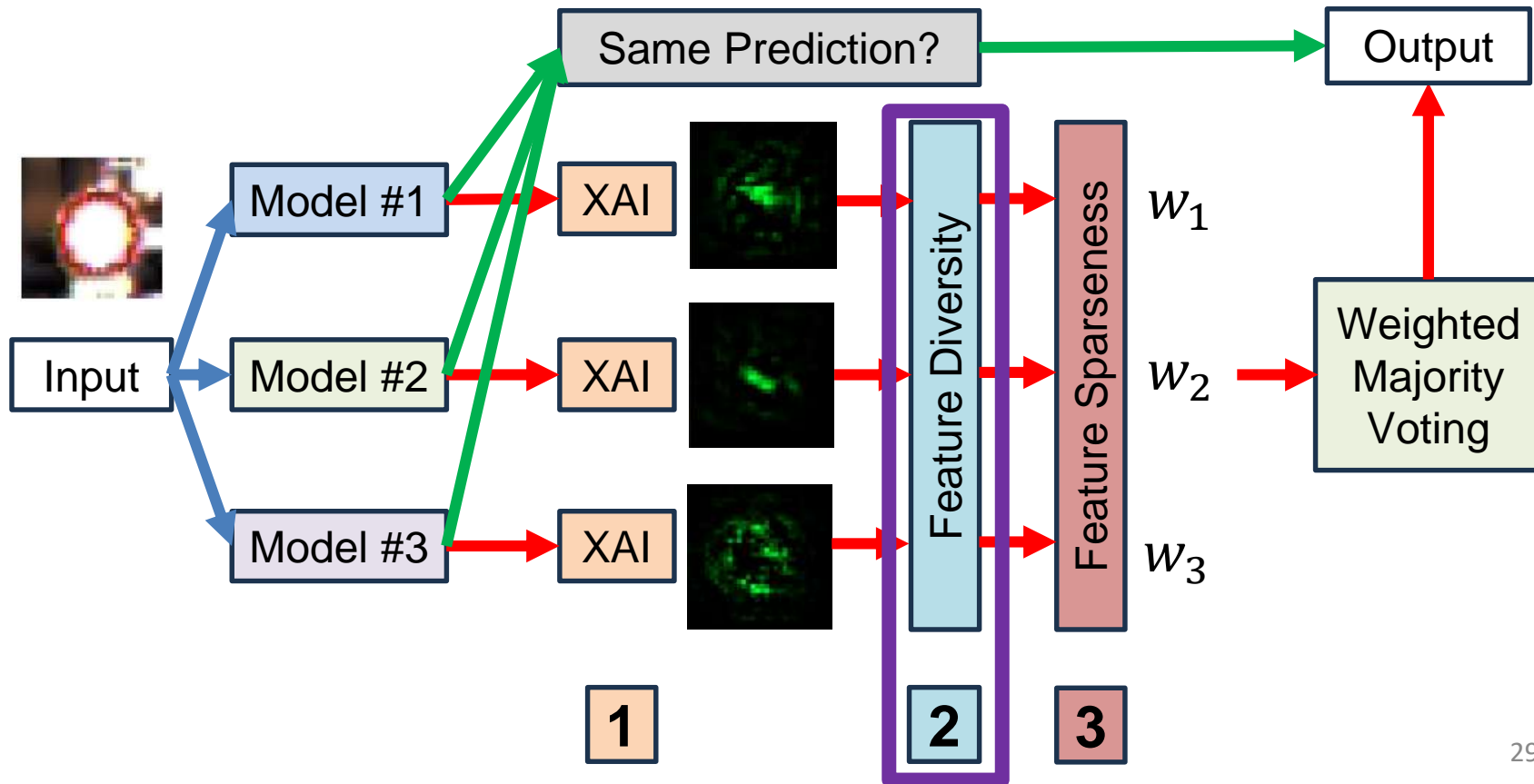
ML Model

No vehicles  
permitted

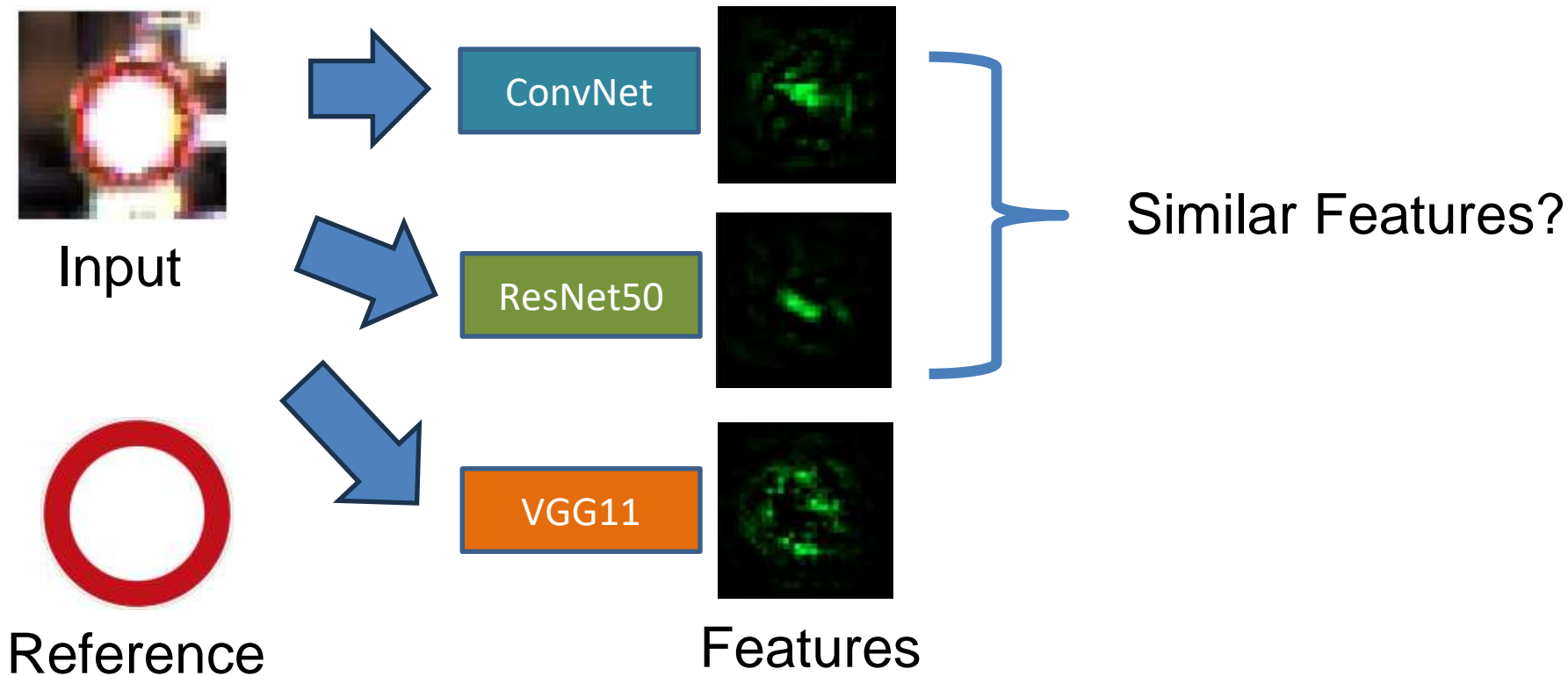
# Smooth Gradients (SG)



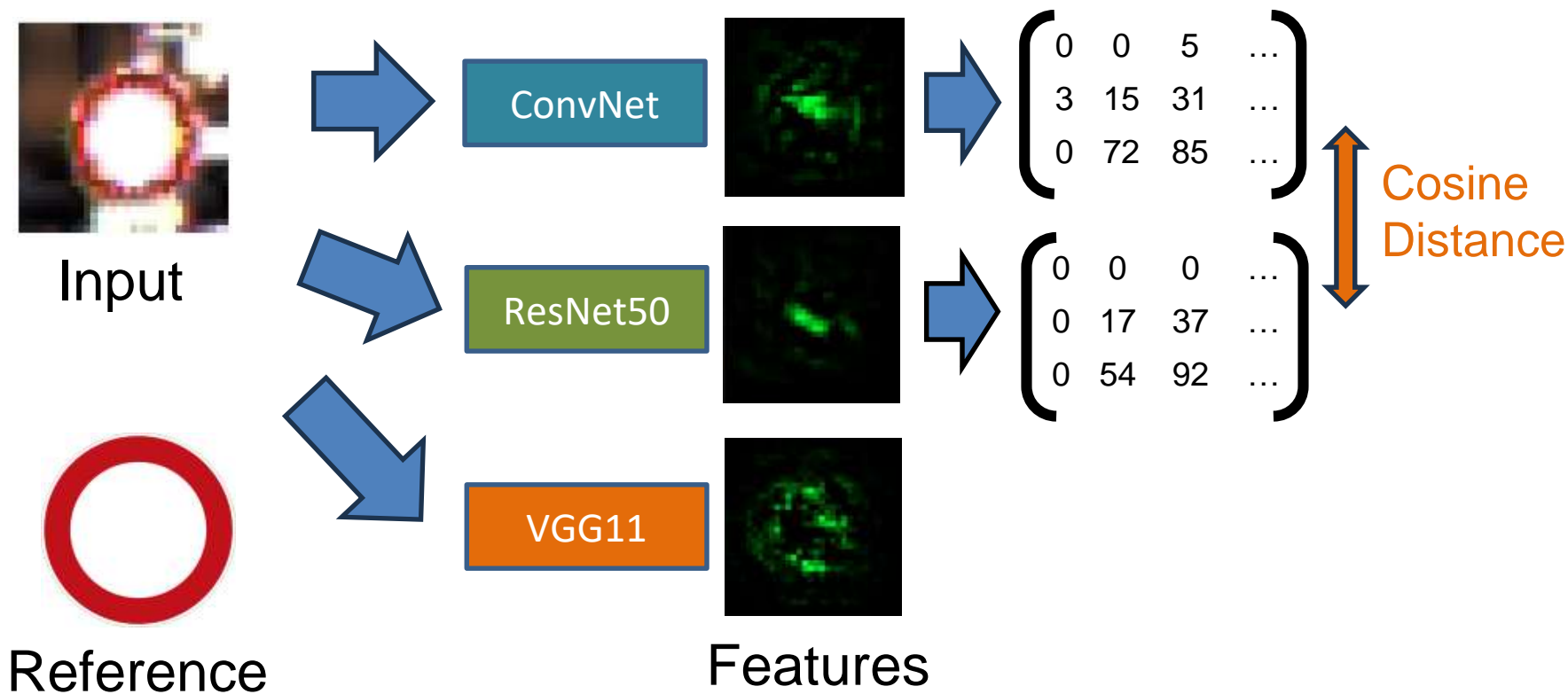
# Step 2 – Feature Diversity



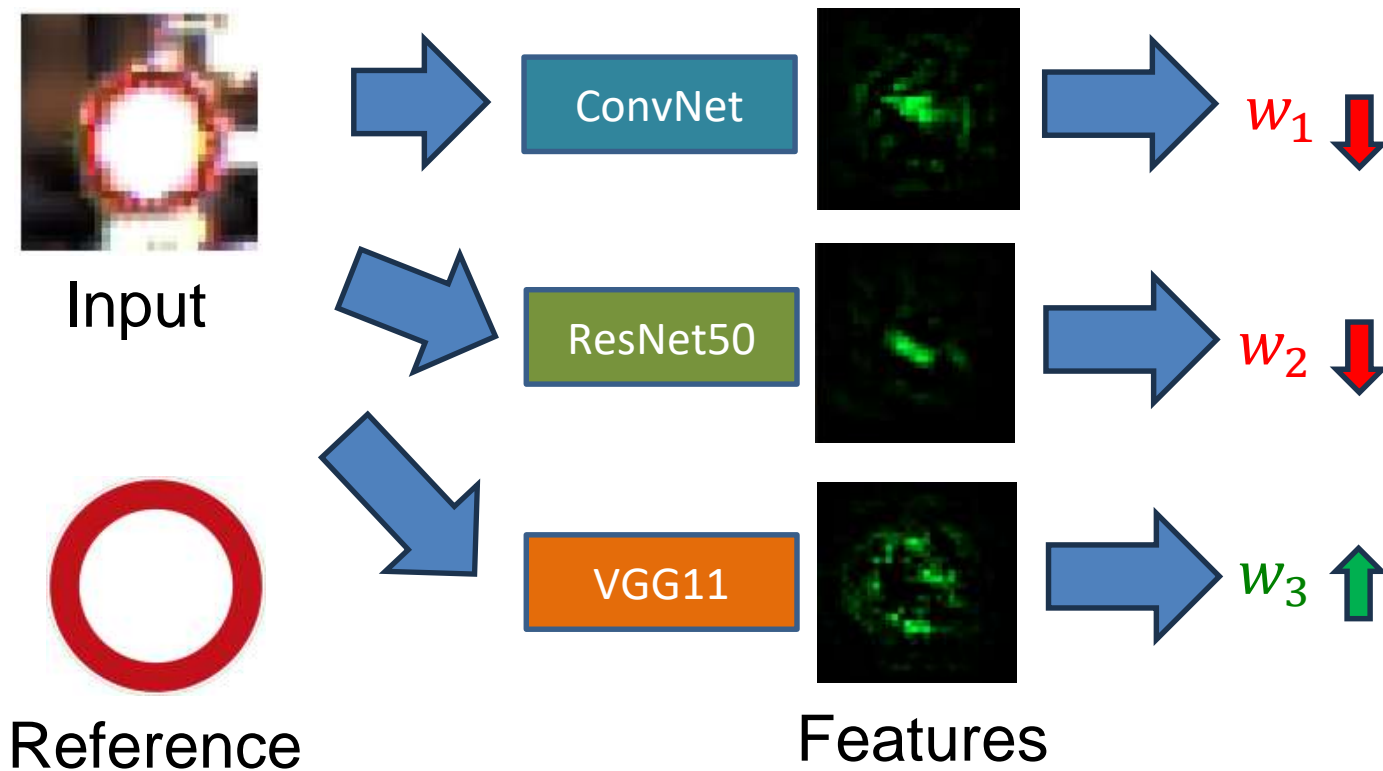
# Feature Diversity using SG



# Feature-Space Diversity

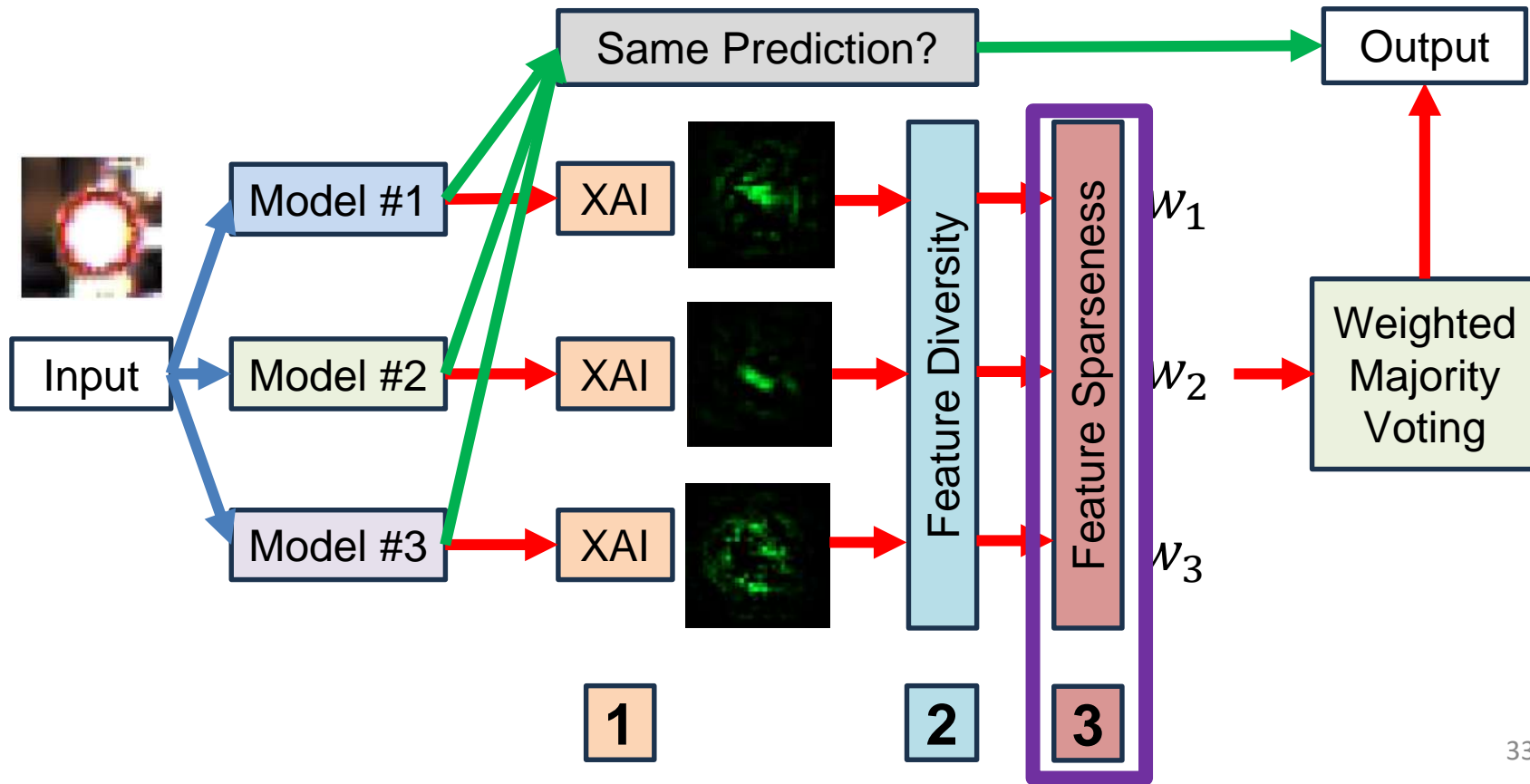


# Dynamic Weights using Feature Diversity





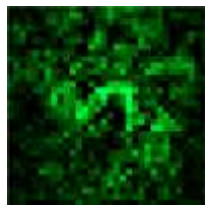
# Step 3 – Feature Sparseness



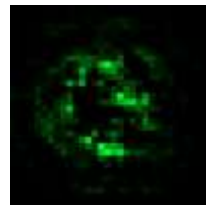
# Feature Sparseness



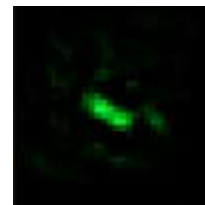
Input



Low



Moderate



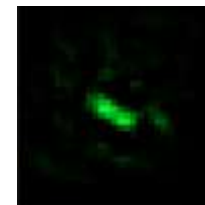
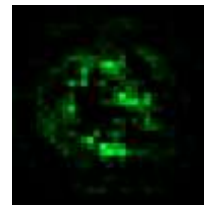
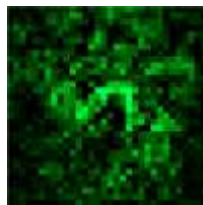
High

**? Useful Diversity ?**

# Feature Sparseness



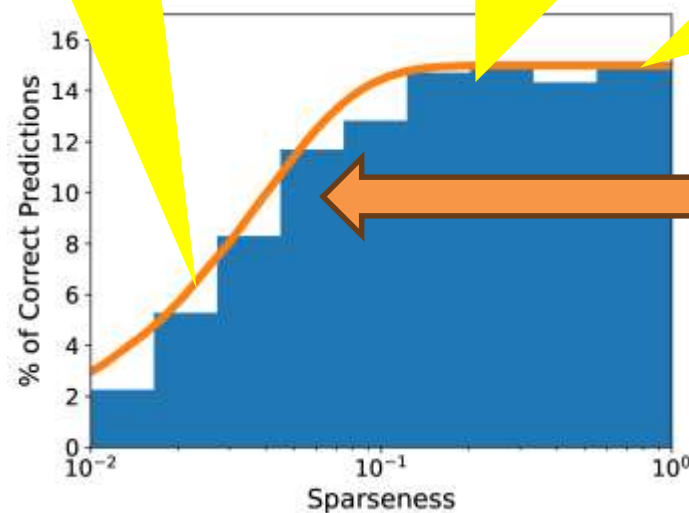
Input



Low

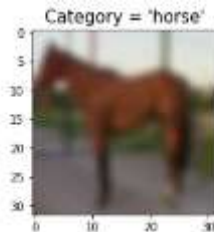
Moderate

High



Tangent  
Activation

# Evaluation Datasets



**CIFAR-10**  
Object Classification



**GTSRB**  
Self-Driving Cars



**Pneumonia**  
Medical Diagnosis

Safety-Critical Applications

# Neural Networks

ML Model Name	Depth (# of Layers)
ConvNet	Shallow
DeconvNet	Shallow
MobileNet	Deep
ResNet18	Deep
ResNet50	Deep
VGG11	Deep
VGG16	Deep

# Resilience Metrics

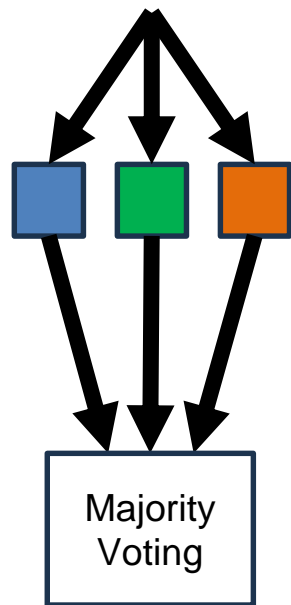
- **Balanced Accuracy**
  - Compatible with imbalanced datasets
- **F1 score**
  - Focus on false positives/negatives than true negatives (e.g. **Pneumonia** [focus case] vs **Benign**)

1 = Most Resilient

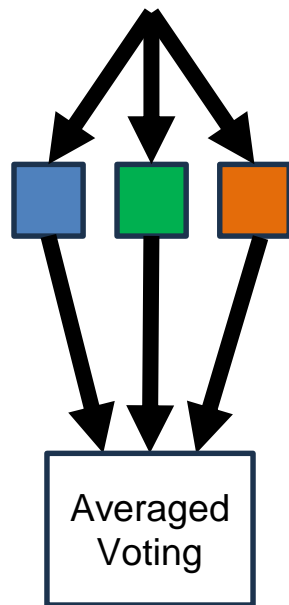


0 = Least Resilient

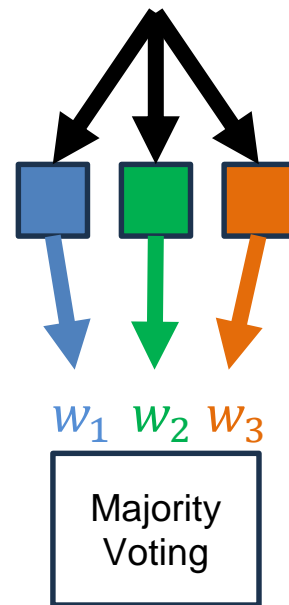
# Ensembling Baselines



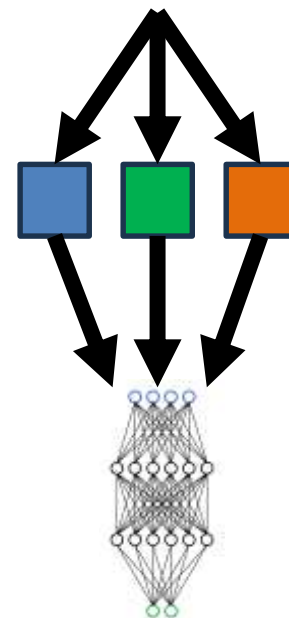
UMaJ



UAvg



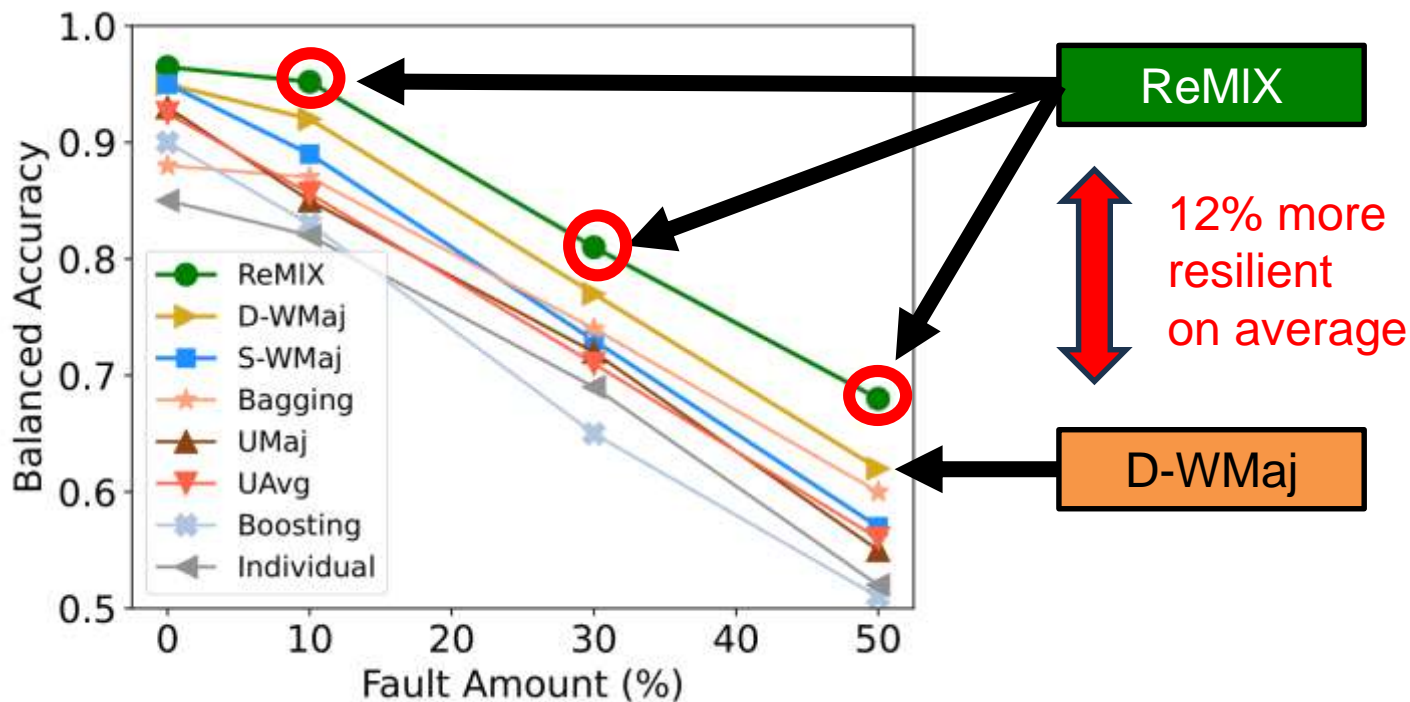
S-WMaJ



D-WMaJ  
(Stacking)

# RQ1: Resilience of ReMIX vs Baselines

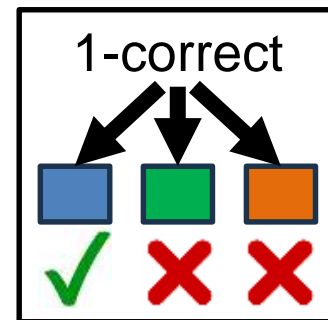
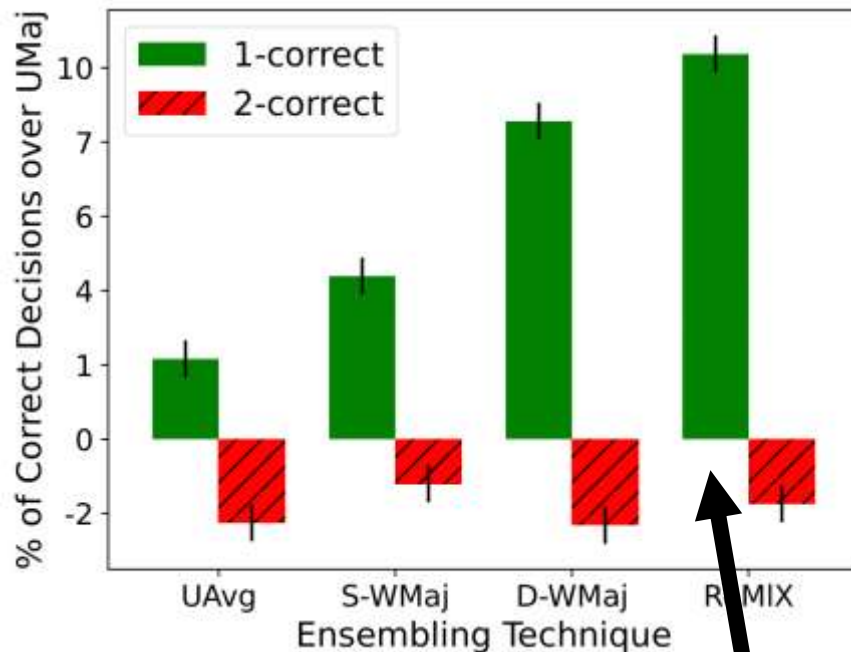
Dataset:  
GTSRB





# RQ1: Resilience of ReMIX vs Baselines

Dataset:  
GTSRB

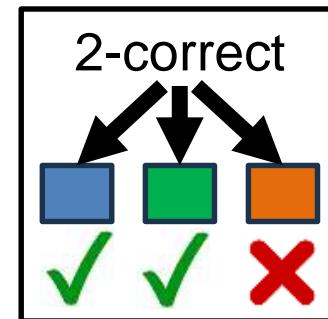
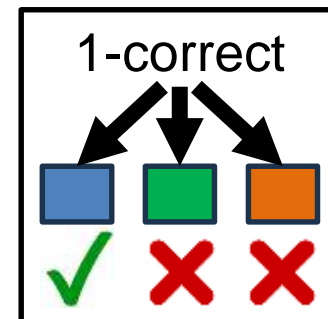
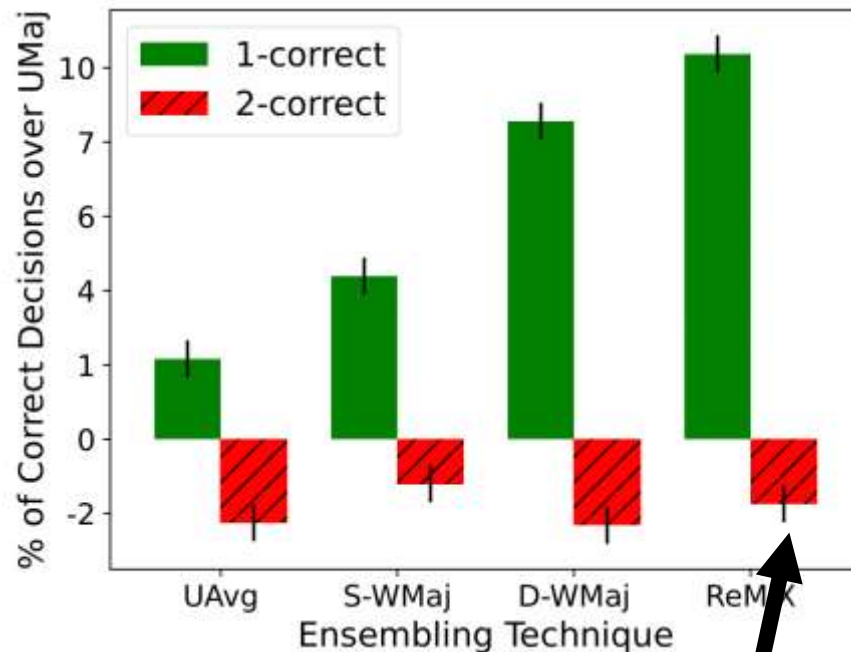


ReMIX maximizes  
predictions of 1-  
correct cases

ReMIX

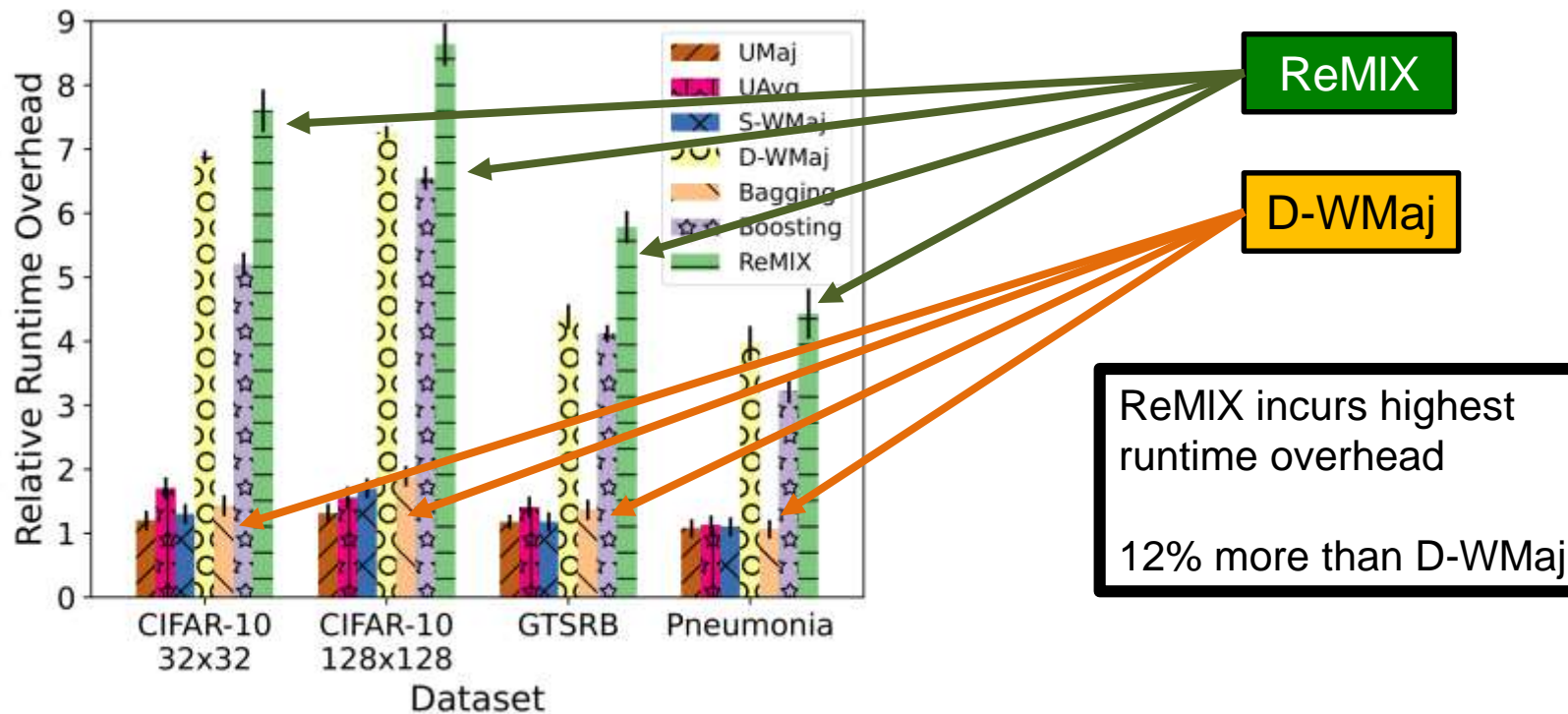
# RQ1: Resilience of ReMIX vs Baselines

Dataset:  
GTSRB



ReMIX

# RQ2: Runtime Overhead



# RQ2: Runtime Overhead

## Application



AVs - GTSRB

**ReMIX  
Average**

0.07s

**Industry  
Maximum**

**0.83s**

**Associated  
Risk**

Safe braking



VR Telesurgery - Pneumonia

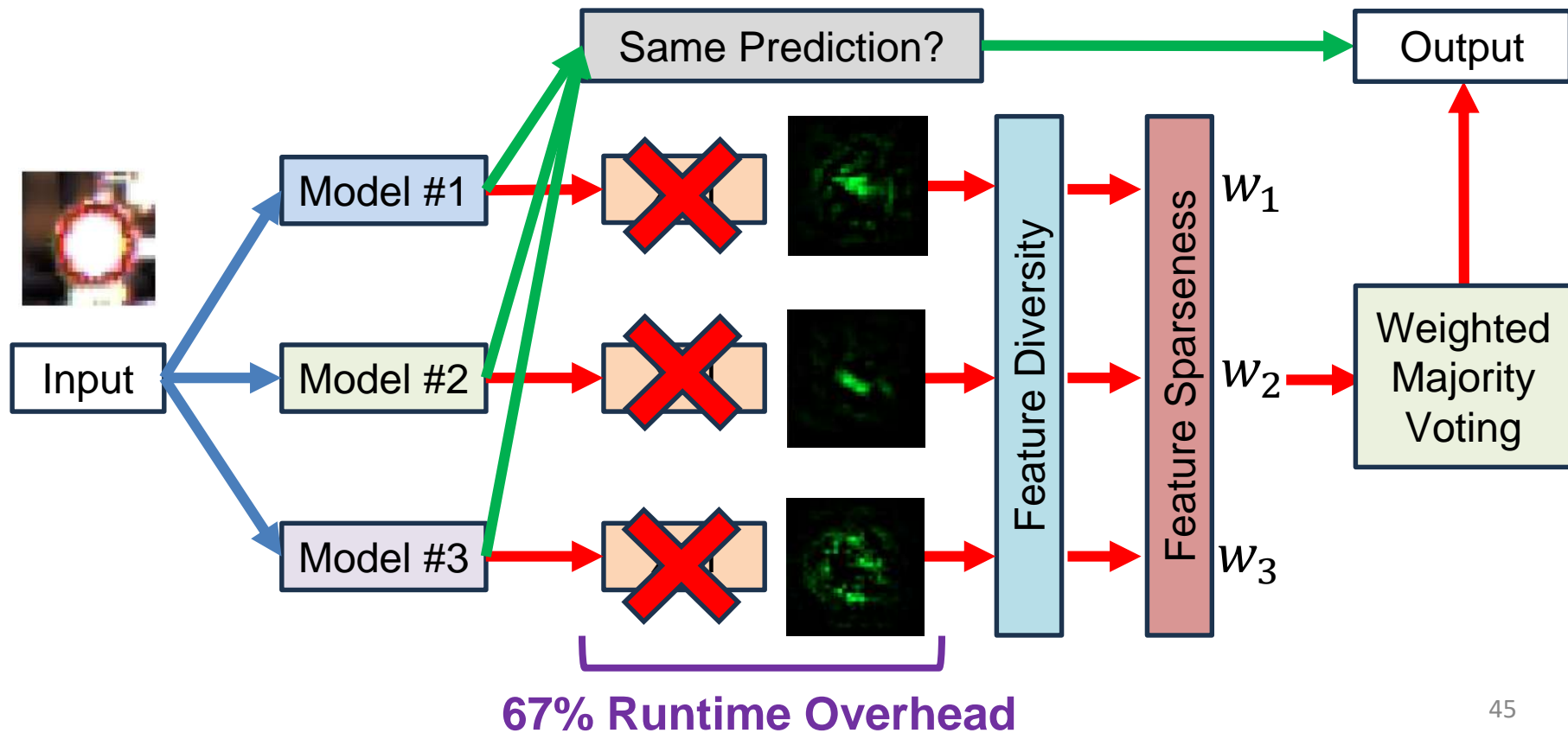
0.31s

**0.50s**

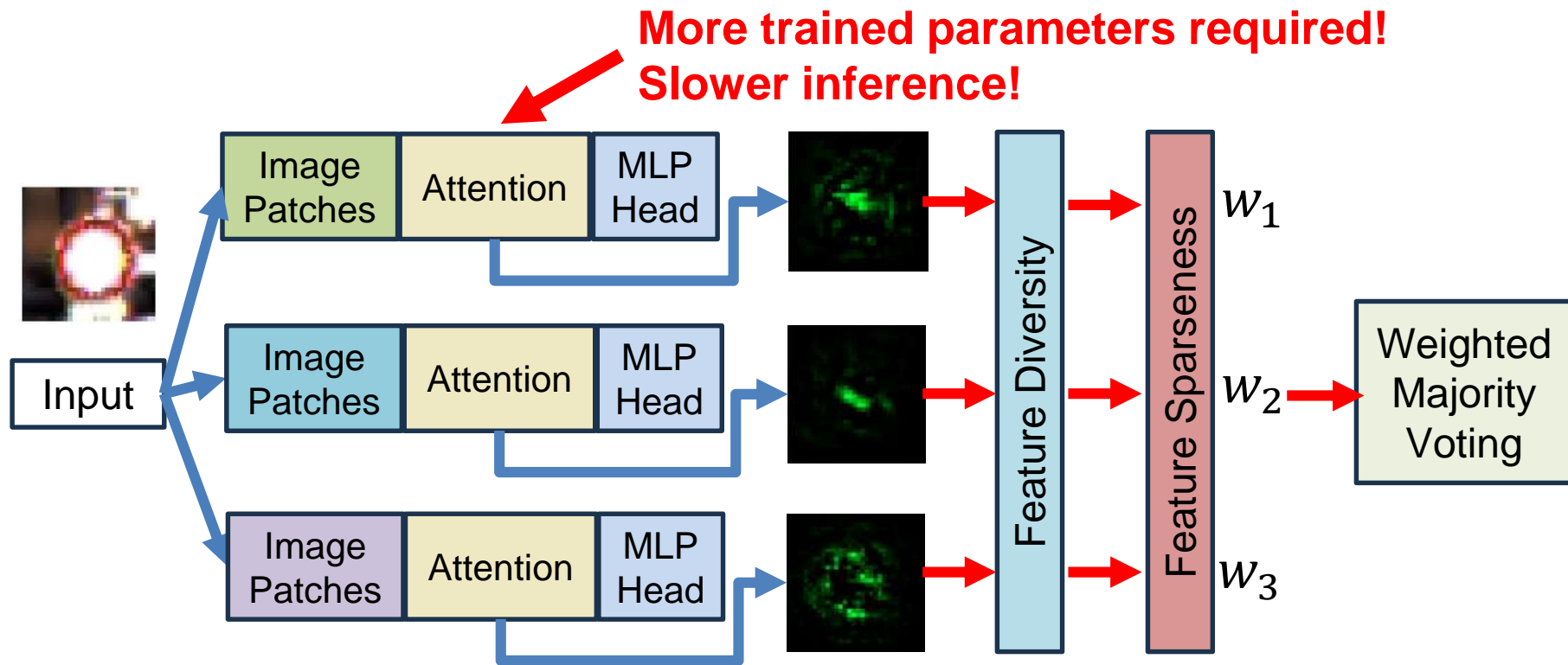
VR sickness



# Optimization using Ante-Hoc XAI



# Optimization using Ante-Hoc XAI

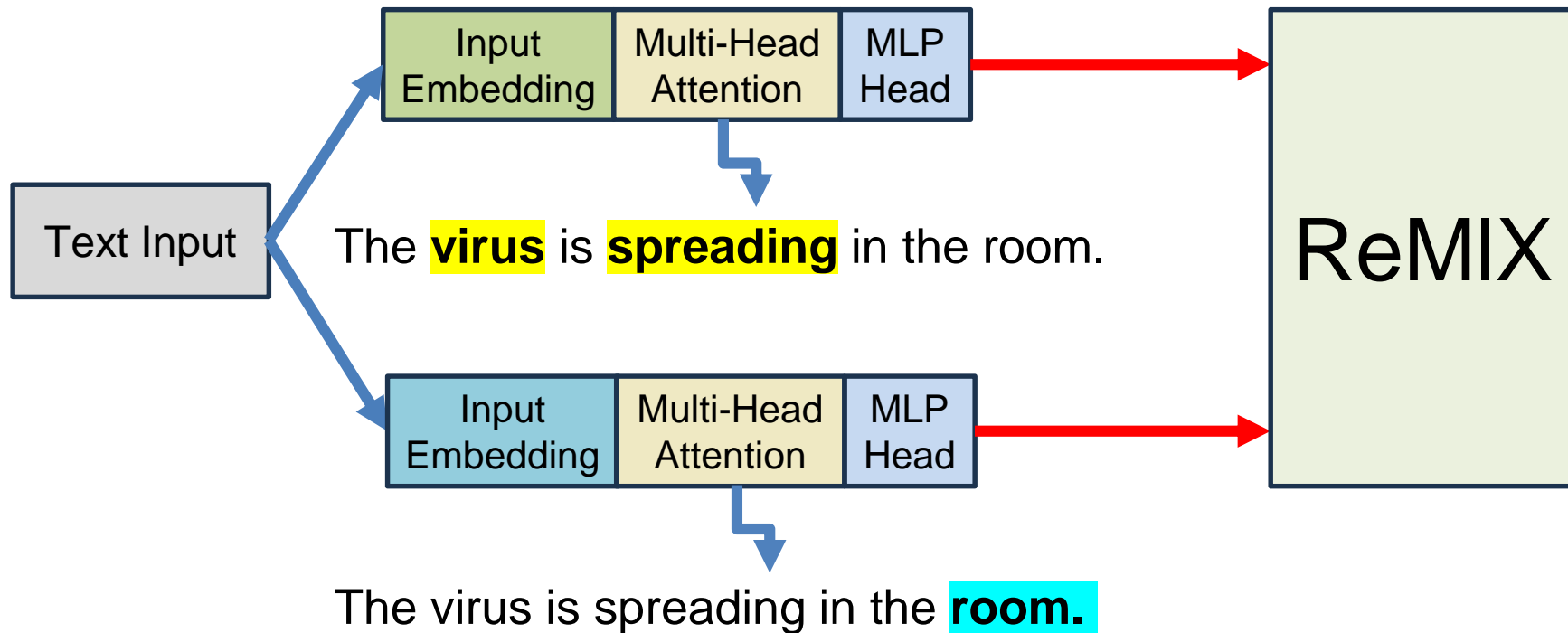


# Future – ReMIX for Sentiment Analysis

Text Input

The virus is spreading in the room.

# Future – ReMIX for Sentiment Analysis





# Summary

1. **Problem:** Reducing misclassifications by ensembles
2. **Approach:** (ReMIX) Resilience of ML Ensembles using XAI
3. **Results:** ReMIX improves resilience by 12% but with 15% overhead over D-WMaj / Stacking (best baseline)

**Email:** [abrahamc@ece.ubc.ca](mailto:abrahamc@ece.ubc.ca)

**Website:** <https://people.ece.ubc.ca/abrahamc/>



Paper



Code



Code  
Available



Code  
Reviewed



Code  
Reproducible