# Effective Methods to Estimate Machine Learning Performance
## Yujie Wang, Shreya Shankar, Aditya Parameswaran
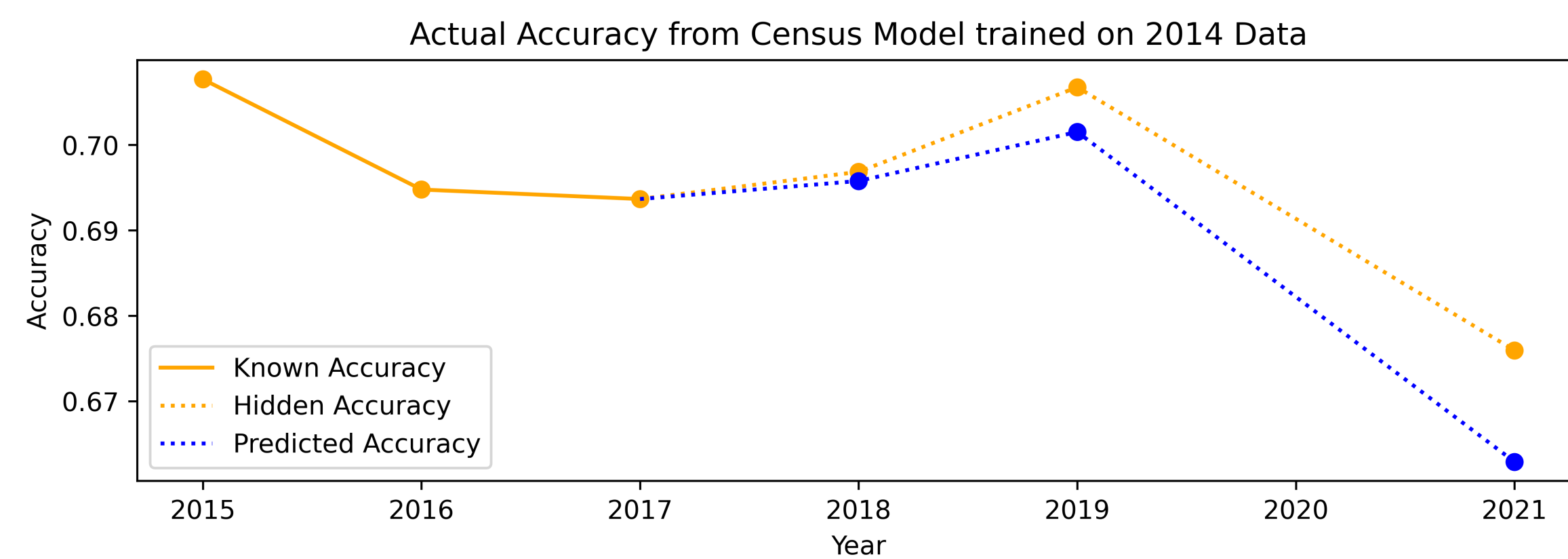
EPIC DATA lab — UC Berkeley

## Motivation

- ML model accuracy, precision, and recall typically change during production due to data shift
- ML Practitioners often have limited access to ground truth labels, preventing true performance tracking
- Performance estimation techniques approximate performance metrics, but there is an estimation gap due to reducible and irreducible error.

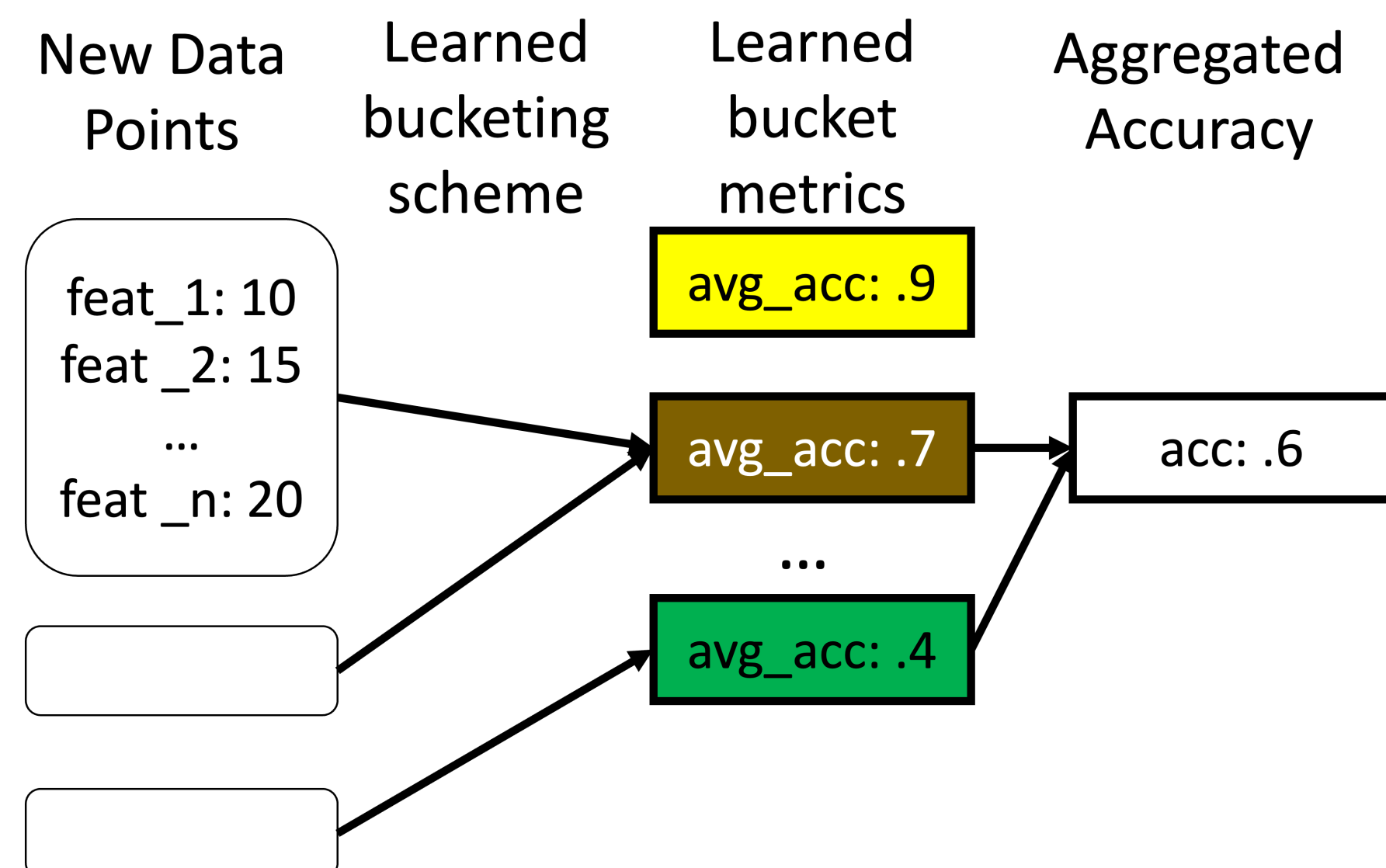**Can we improve the current tools used to estimate ML performance?**



Actual Accuracy from Census Model trained on 2014 Data
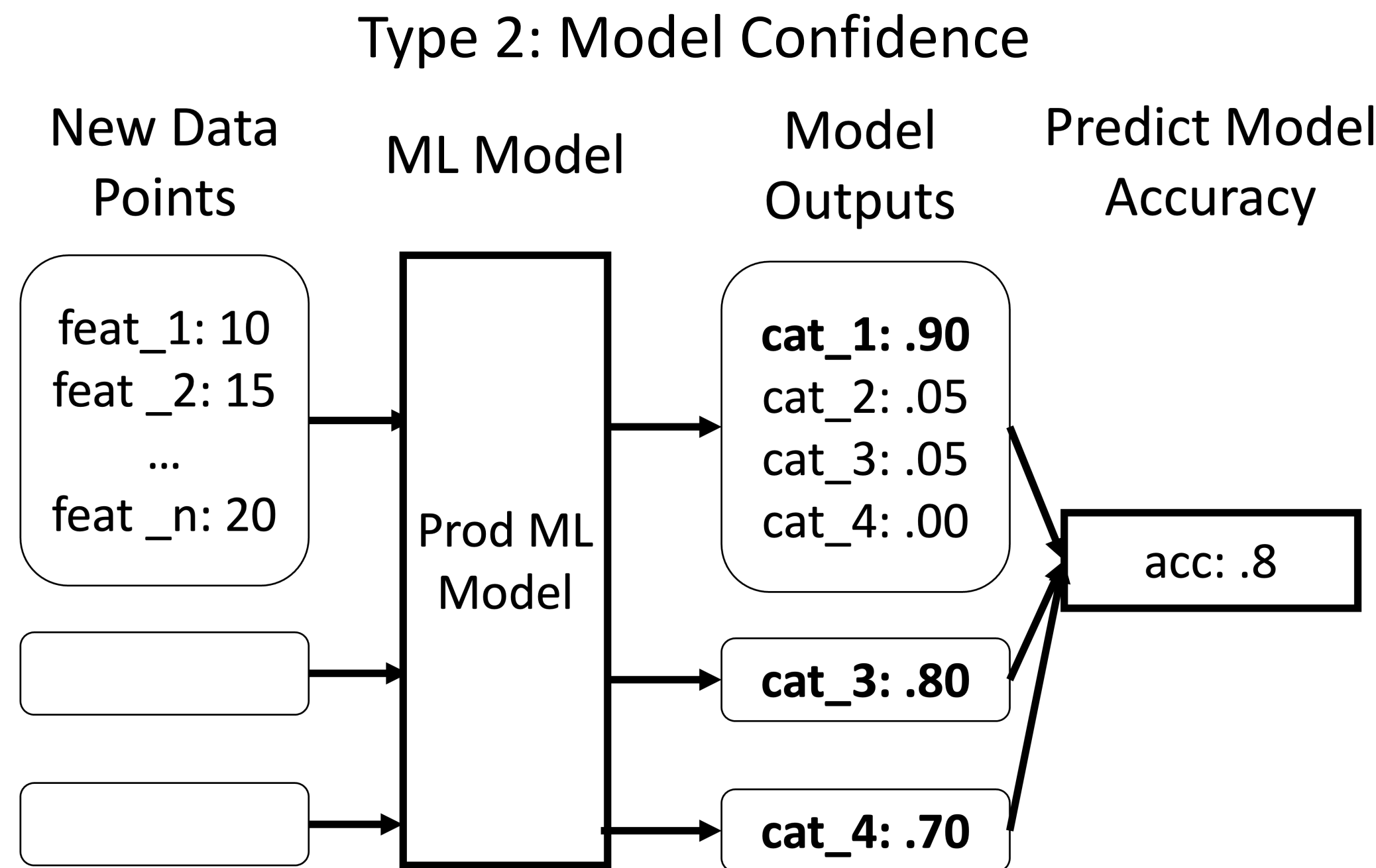
## Introduction

Performance estimation for ML models broadly follows 2 types of techniques.

### Type 1: Importance-Weighted

- Separate data into buckets by features, e.g., feat_1 < 10, feat_1 >= 10.
- Learn accuracy of validation dataset for each bucket



New Data Points → Learned bucketing scheme → Learned bucket metrics → Aggregated Accuracy

avg_acc: .9 / avg_acc: .7 / avg_acc: .4 → acc: .6

Problem: buckets based on individual features scale poorly with high-dimensional inputs (see Evaluation)

## Type 2: Model Confidence



New Data Points → ML Model → Model Outputs → Predict Model Accuracy

feat_1: 10, feat_2: 15, ... feat_n: 20 → Prod ML Model → cat_1: .90, cat_2: .05, cat_3: .05, cat_4: .00 → acc: .8
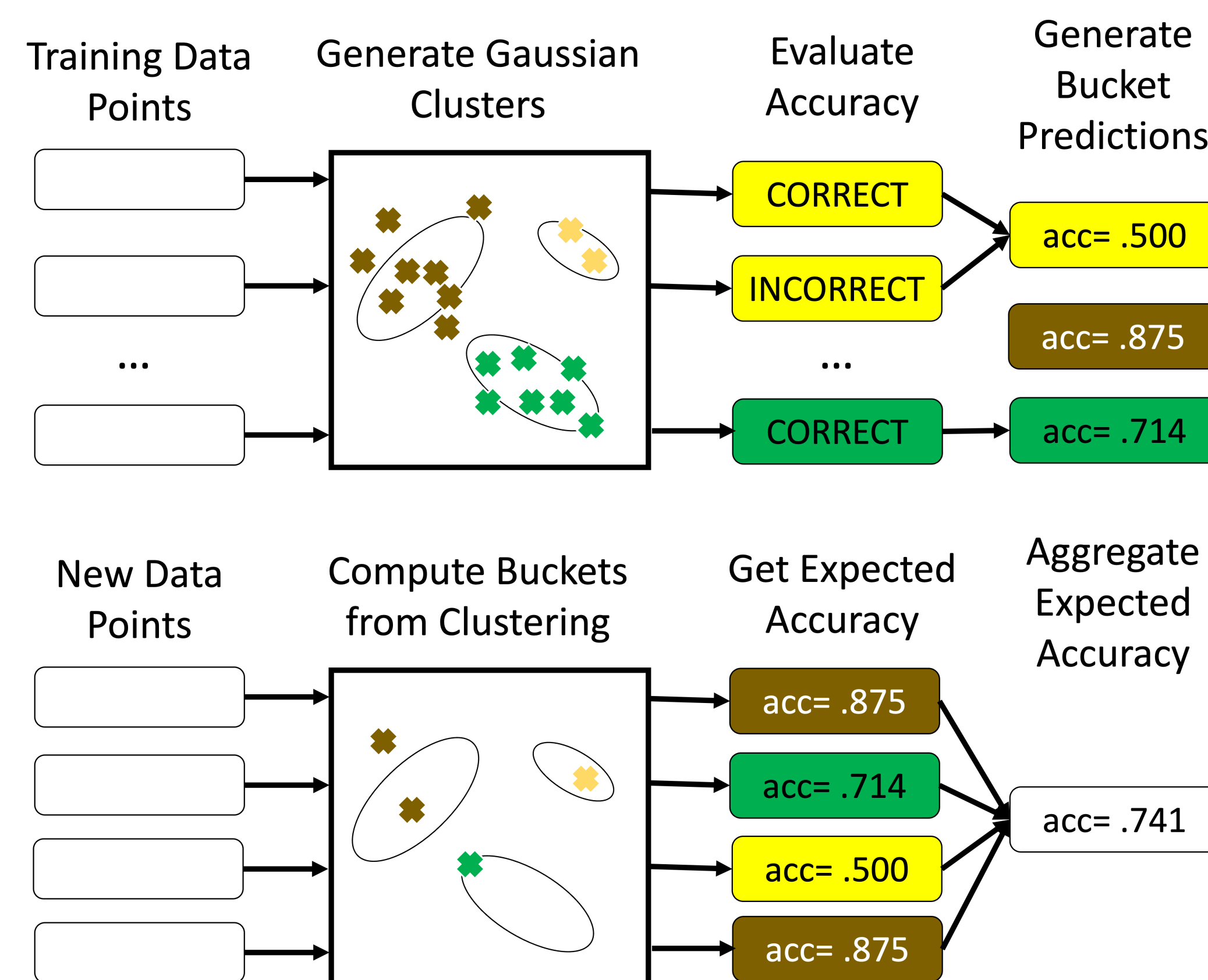
cat_3: .80

cat_4: .70

Problem: model confidence performs erratically when out of distribution
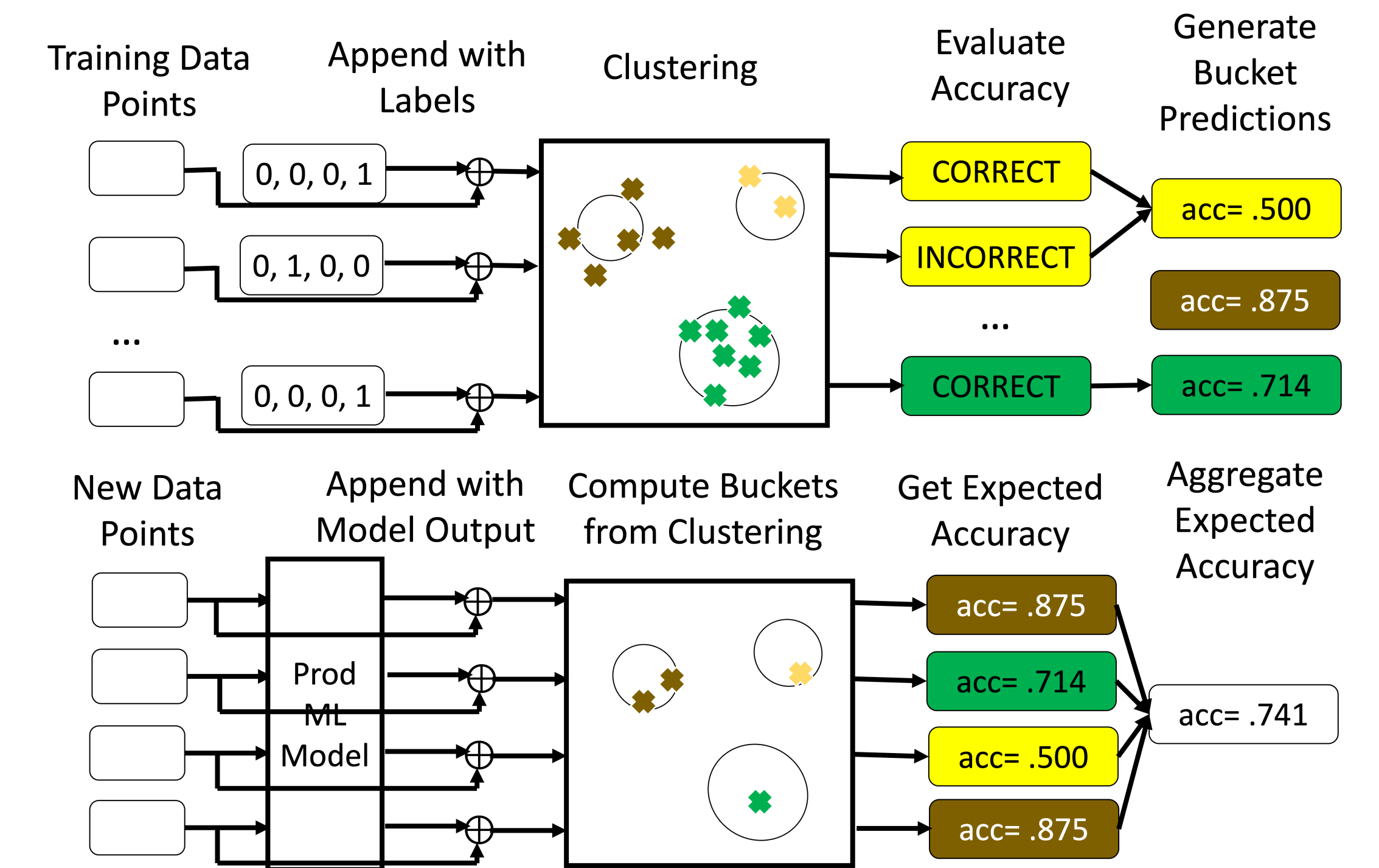
## Methods Under Research

### Novel Method 1: Gaussian Mixture Importance-Weighting

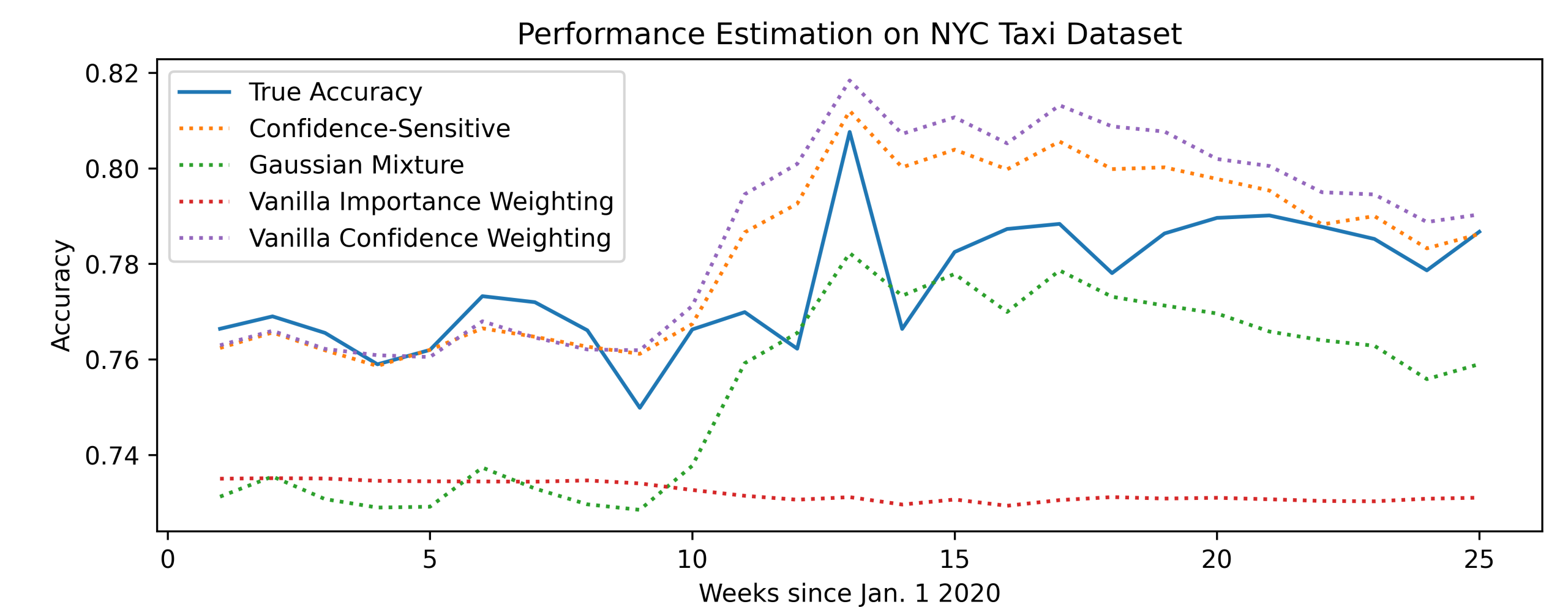A type of Importance-Weighted Performance Prediction (Type 1) whose bucketing scheme is across multiple features.



Training Data Points → Generate Gaussian Clusters → Evaluate Accuracy → Generate Bucket Predictions

CORRECT / INCORRECT / ... / CORRECT → acc: .500 / acc: .875 / acc: .714

New Data Points → Compute Buckets from Clustering → Get Expected Accuracy → Aggregate Expected Accuracy

acc= .875 / acc= .714 / acc= .500 / acc= .875 → acc= .741

Contribution: Allows importance-weighted technique (Type 1) to generalize to higher dimensions.

## Novel Method 2: Confidence-Sensitive Clustering



Training Data Points → Append with Labels → Clustering → Evaluate Accuracy → Generate Bucket Predictions

0, 0, 0, 1 / 0, 1, 0, 0 / ... / 0, 0, 0, 1 → CORRECT / INCORRECT / ... / CORRECT → acc= .500 / acc= .875 / acc= .714

New Data Points → Append with Model Output → Compute Buckets from Clustering → Get Expected Accuracy → Aggregate Expected Accuracy

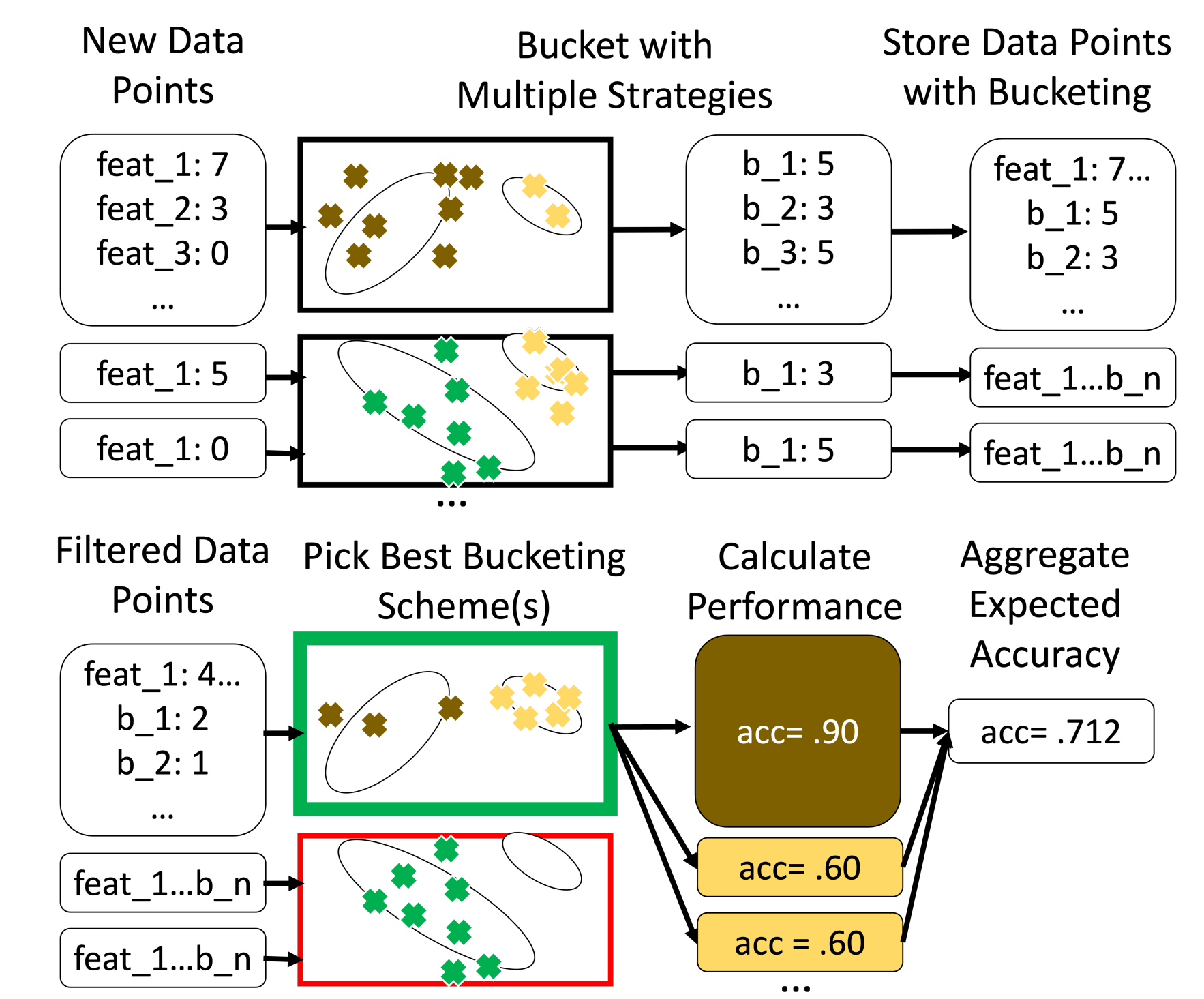Prod ML Model → acc= .875 / acc= .714 / acc= .500 / acc= .875 → acc= .741

Contribution: Merges Importance-Weighting (Type 1) and Model Confidence (Type 2), allowing for confidence-aware clusters

## Preliminary Evaluation on Challenging Dataset



Performance Estimation on NYC Taxi Dataset

Legend: True Accuracy, Confidence-Sensitive, Gaussian Mixture, Vanilla Importance Weighting, Vanilla Confidence Weighting

## Roadmap Vision: Dynamic Bucket Strategy



New Data Points → Bucket with Multiple Strategies → Store Data Points with Bucketing

feat_1: 7, feat_2: 3, feat_3: 0, ... → b_1: 5, b_2: 3, b_3: 5, ... → feat_1: 7..., b_1: 5, b_2: 3

feat_1: 5 → b_1: 3

feat_1: 0 → b_1: 5

Filtered Data Points → Pick Best Bucketing Scheme(s) → Calculate Performance → Aggregate Expected Accuracy

feat_1: 4..., b_1: 2, b_2: 1 → acc= .90 / acc= .60 / acc = .60 → acc= .712

- Ensemble buckets enable better importance-weighting coverage
- Storage of bucket information allows for fast accuracy estimates for arbitrary data, enabling faster lookup of accuracy drop cause