

# A Calibration-Aware Framework for Geometric Support Evaluation of Synthetic Tabular Data

Lee Carlin Yuval Benjamini  The Hebrew University of Jerusalem

## Why Current Tabular Support Evaluation is Unreliable

- Evaluating synthetic generators is an ongoing challenge, especially for tabular mixed-type datasets.
- Characterized by high-cardinality features, non-linear dependencies, and lack of canonical representations.
- Practitioners have adopted *Geometric Support Metrics* (Precision/Recall) from Image/NLP domains.
- However, they assume a reliable embedding space. In tabular data, used as "black-box" with unknown validity.
- Sparse categorical structures and conditional constraints distort support estimation, leading to unreliable diagnostics.

## Support Metrics: Formal Setup

Date:  $\mathcal{D}_{\text{real}} = \{x_i\}_{i=1}^N, \mathcal{D}_{\text{syn}} = \{y_j\}_{j=1}^M \subset \mathcal{X}$ .

Embedder:  $\Phi: \mathcal{X} \rightarrow \mathbb{R}^p$  maps  $\mathcal{Z}_{\text{real}} = \{\Phi(x_i)\}_{i=1}^N \sim P, \mathcal{Z}_{\text{syn}} = \{\Phi(y_j)\}_{j=1}^M \sim Q$

The support  $S_P(\alpha)$  denotes a high-density region of  $P$ , for a probability level  $\alpha$ :

$$S_P(\alpha) = \{z \in \mathbb{R}^p : f(z) \geq \tau_\alpha\}.$$

**Precision (Fidelity):**

$$\Pr(P, Q) = 1 - \Pr_{z \sim Q}[z \in \widehat{S}_P(\alpha)]$$

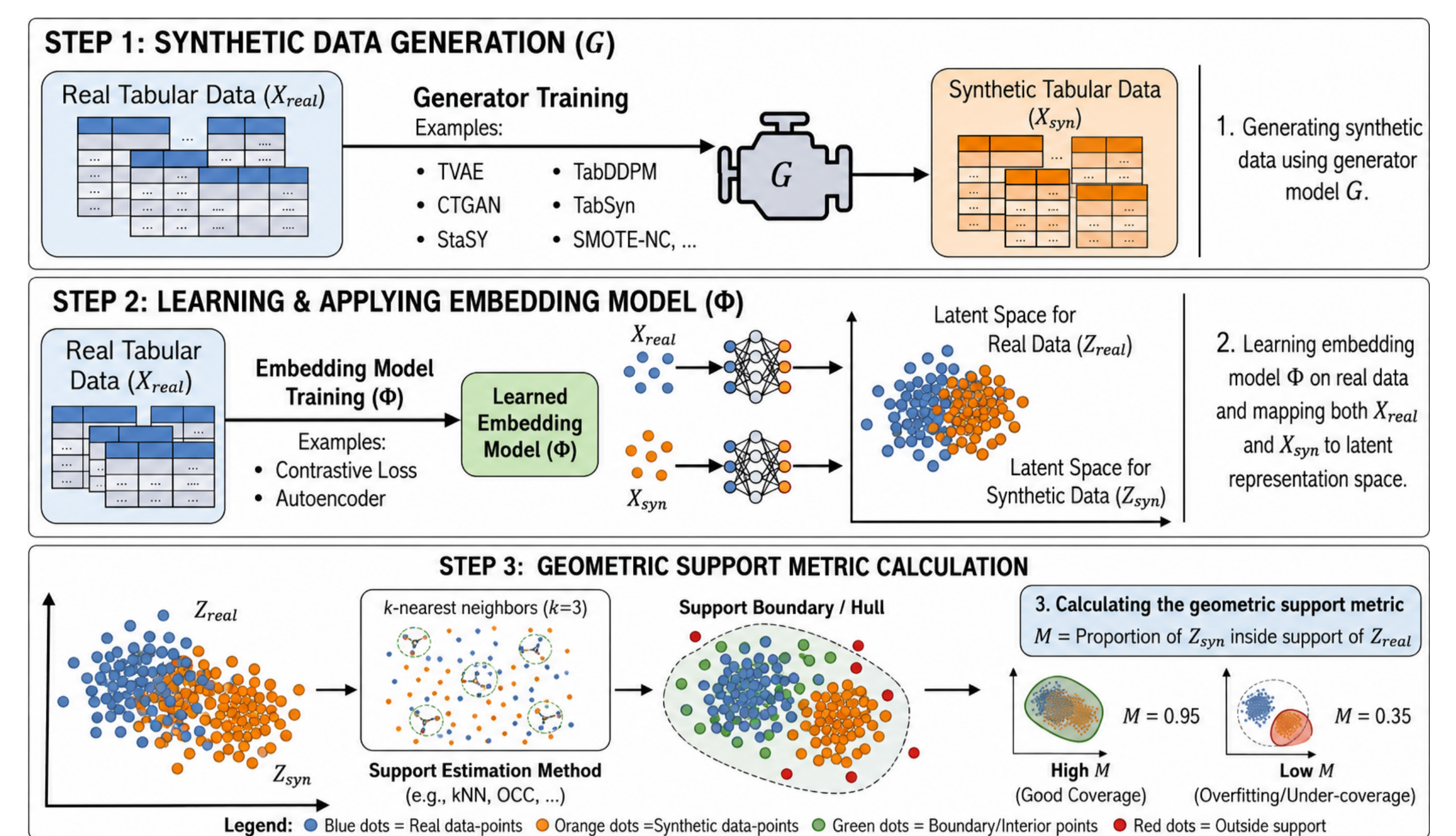
Fraction of synthetic samples *inside* real support.

**Recall (Diversity):**

$$\text{Re}(P, Q) = 1 - \Pr_{z \sim P}[z \in \widehat{S}_Q(\alpha)]$$

Fraction of real samples *covered* by synthetic support.

## Current Geometric Support Framework

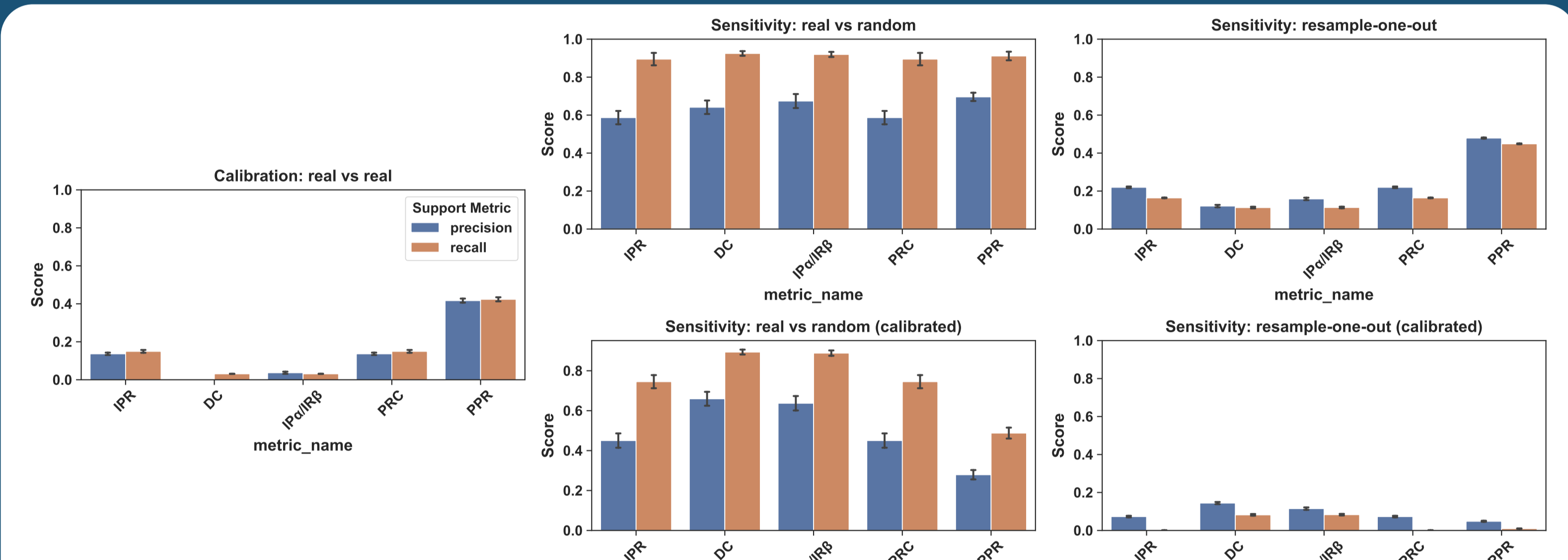


Geometric Support Metric ( $M$ ) workflow for synthetic tabular data. The process involves generating synthetic data ( $X_{\text{syn}}$ ) from real data ( $X_{\text{real}}$ ) using generator  $G$  (Step 1), learning an embedding model  $\Phi$  on real data and mapping both  $X_{\text{real}}$  and  $X_{\text{syn}}$  to latent representations  $Z_{\text{real}}$  and  $Z_{\text{syn}}$  (Step 2), and finally calculating ( $M$ ) as the proportion of  $Z_{\text{syn}}$  points falling within the estimated geometric support boundary (hull) of  $Z_{\text{real}}$  (Step 3).

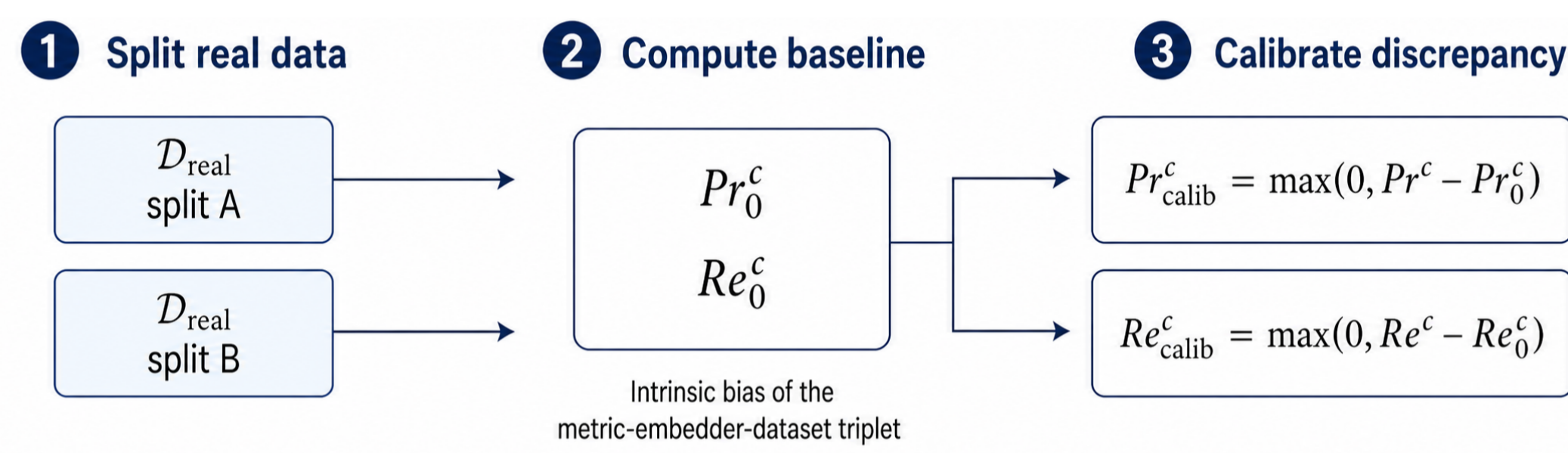
## Baseline Bias

- **Empirical Bias:** Support metrics often assign non-zero discrepancy to samples drawn from *identical* distributions.
- **Sensitivity:** Raw scores ( $Pr, Re$ ) are highly sensitive to:
  1. Embedding architecture (PCA vs. DSVD vs. Contrastive).
  2. Preprocessing (Scaling, Encoding).
  3. Hyperparameters (e.g.,  $k$ -nearest neighbors).
- **Inconsistency:** Generators are often ranked differently depending purely on the choice of the embedder, not the data quality.

## Baseline Bias: Calibration-Aware Scores



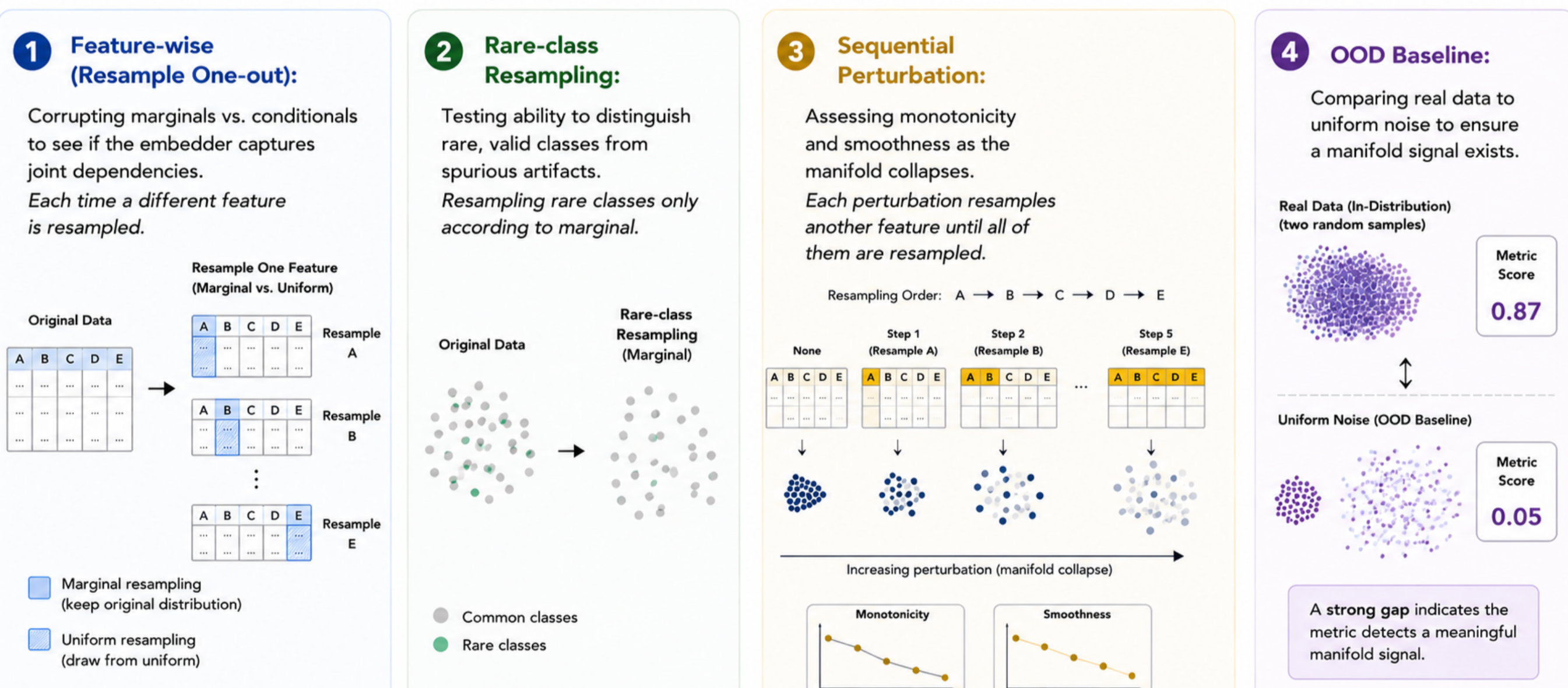
## Contribution 1: Baseline Calibration Protocol



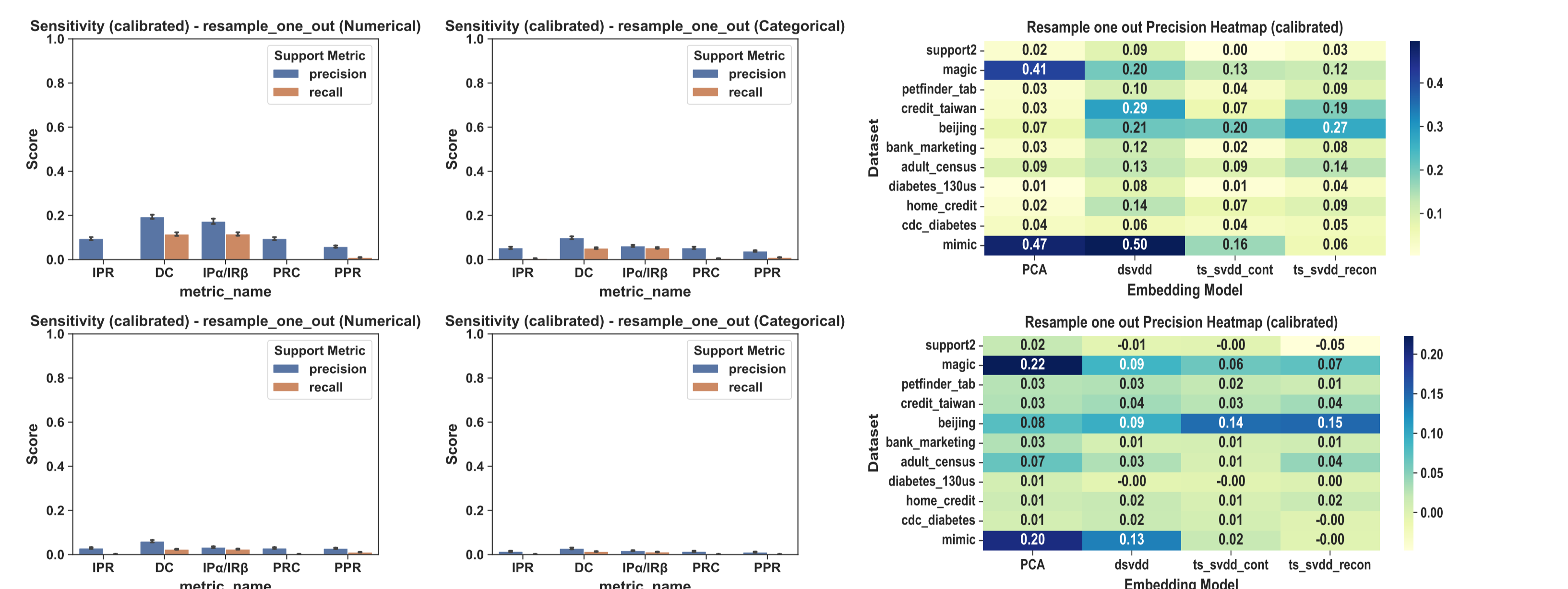
Raw geometric support scores can be dominated by embedder-induced bias  
Calibration materially improves interpretability and ranking stability

## Contribution 2: Sensitivity Validation Tests

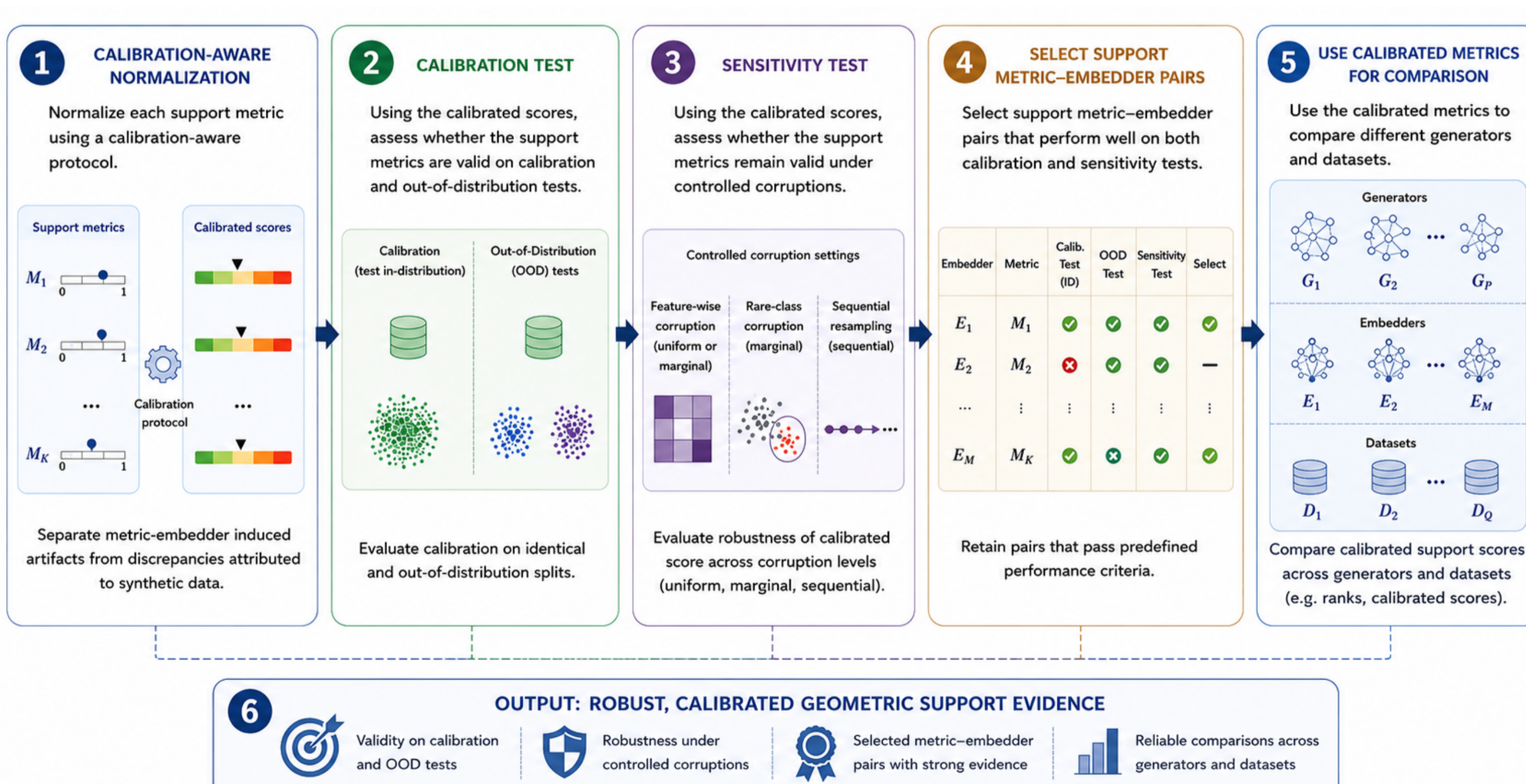
Before trust, we verify the calibrated metric's sensitivity to tabular-specific shifts:



## Sensitivity Validation Tests: Tabular-Specific Shifts

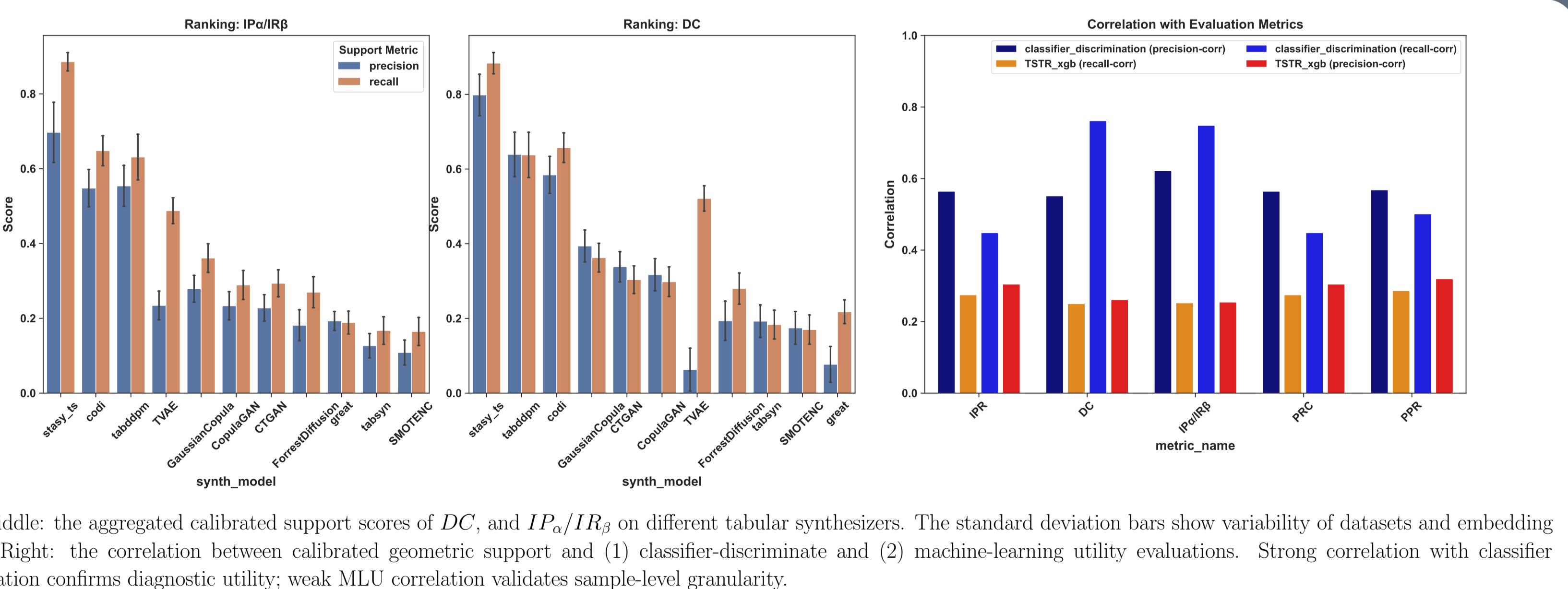


## Calibration-Aware Geometric Support Framework



12 Datasets, 5 Metrics, 4 Embedders, 11 Generators

## Calibration-Aware Generator Support Benchmark



## Takeaways

- Support metrics are **not** black-box tools; they are pipeline-dependent and require dataset-specific calibration.
- Our protocol separates **baseline representation bias** from actual **synthetic data discrepancies**.
- Benchmarking reveals current metrics favor majority modes and **struggle with tabular conditional dependencies**.