

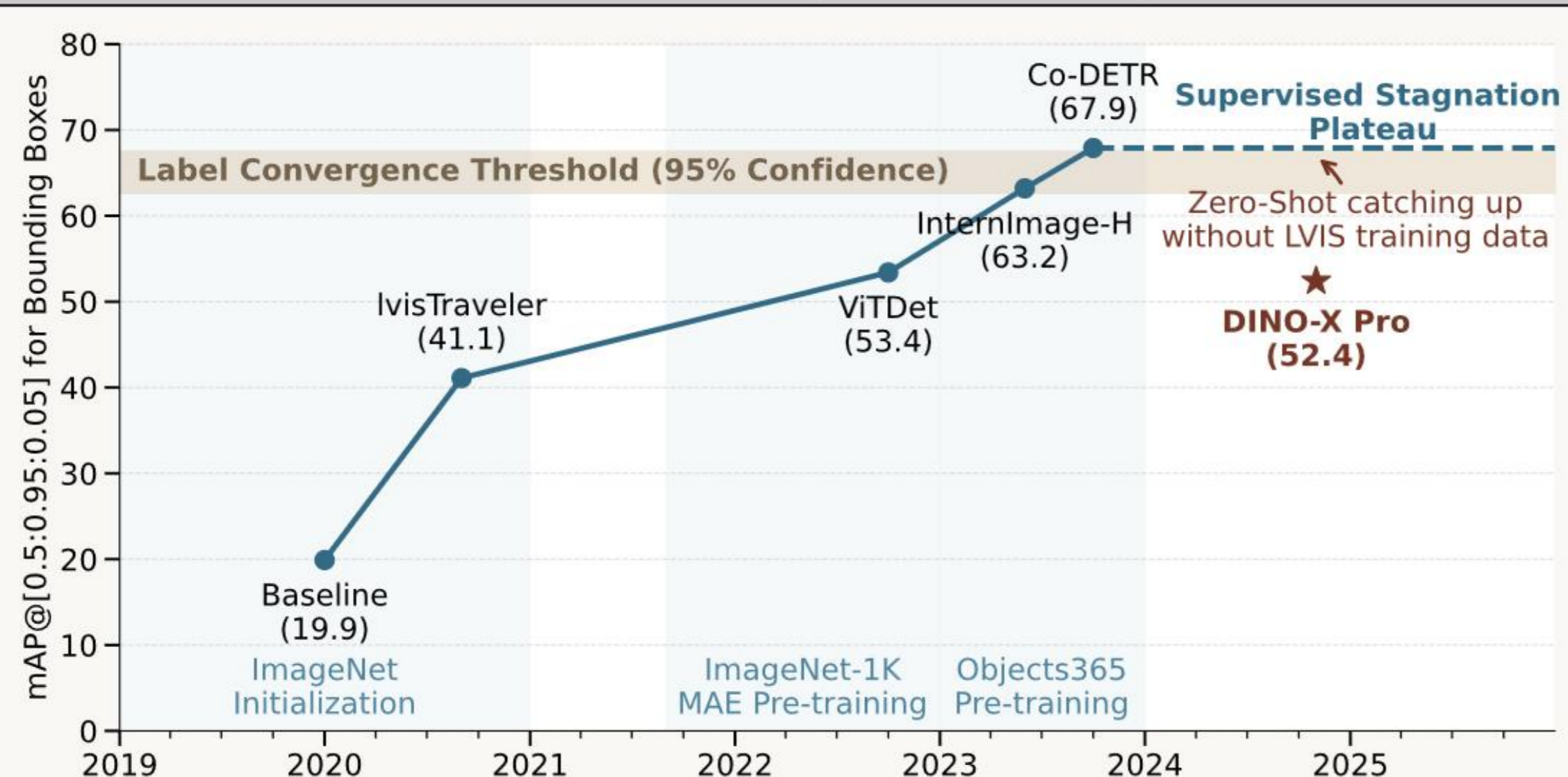
# KaLOS finds Consensus: A Meta-Algorithm for Evaluating Inter-Annotator Agreement in Complex Vision Tasks

DAVID TSCHIRSCHWITZ and VOLKER RODEHORST

BAUHAUS-UNIVERSITÄT WEIMAR, GERMANY

## I. Motivation and Problem Description

### 1. Motivation: Overcome Stagnation in Object Detection



History of top-performing models on LVIS v1.0 val, showing performance improvements stagnating at the label convergence threshold.

### 2. Problem: Contradicting Annotations cause a Bottleneck

Label noise creates a **Supervised Stagnation Plateau**, capping model progress. Improving data quality is now more critical than scaling architectures.

#### Our Solution: KaLOS

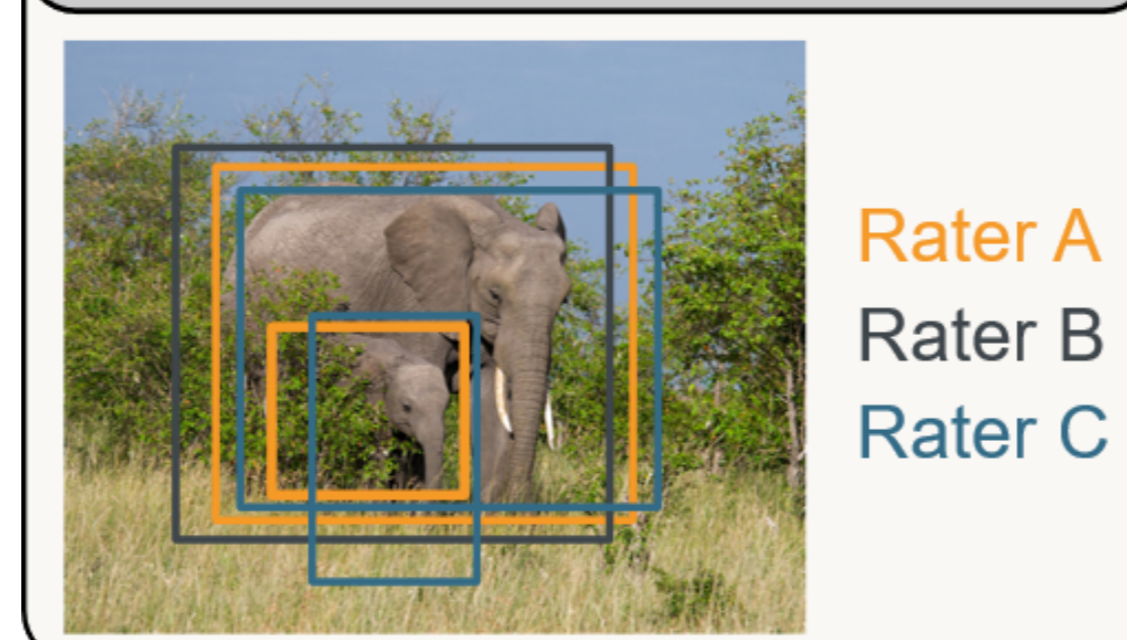
We introduce **KaLOS** to unify Inter-Annotator Agreement (IAA) across complex tasks. It identifies systematic errors in human, synthetic and auto-annotated data.

**LVIS Example:** Disputed 'Train' topology (Lead car vs. entire consist). How do we quantify consensus for such conceptual ambiguity?

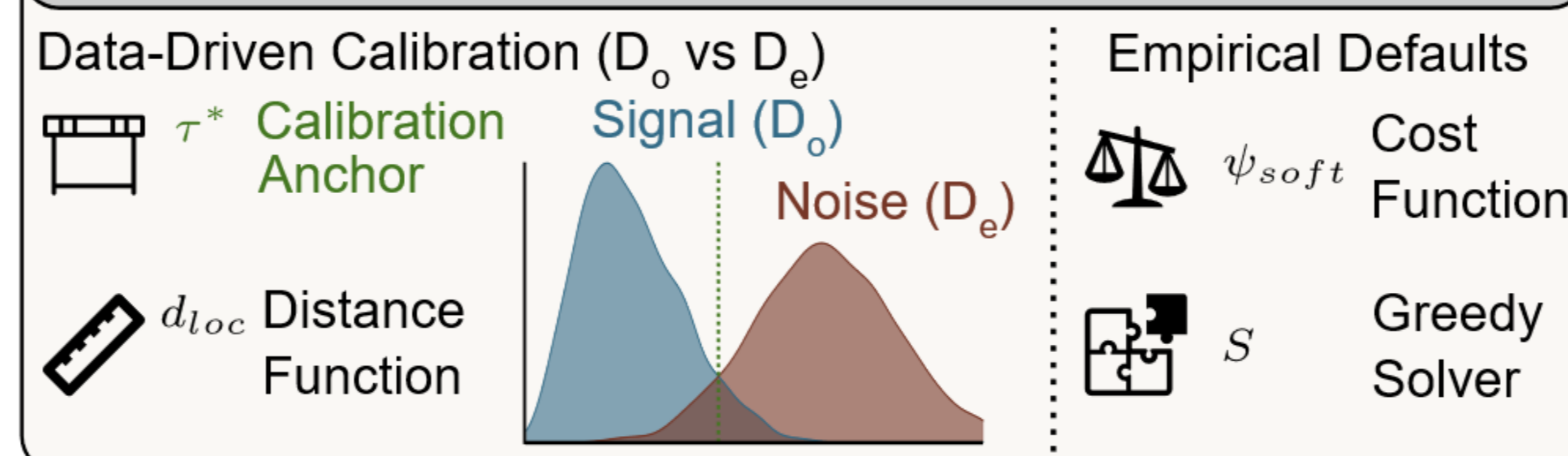


## II. Method: Evaluate Dataset Quality

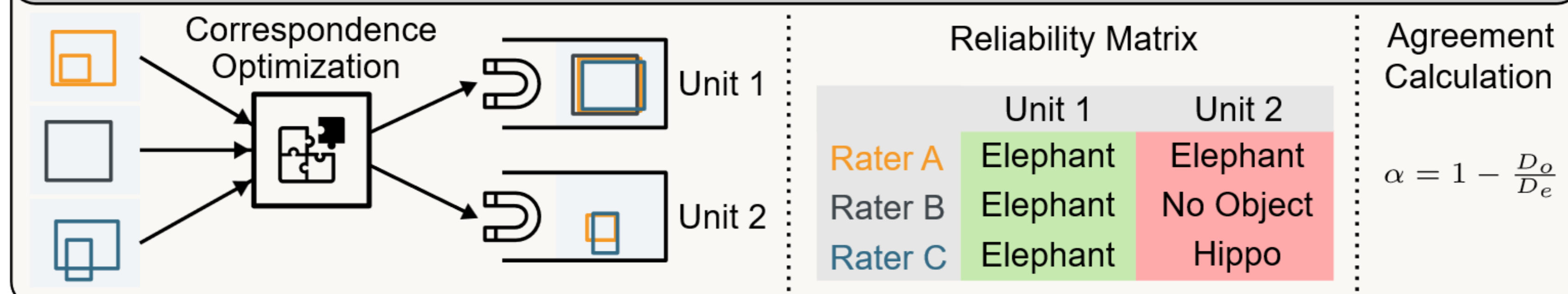
### 1. Input: Active Disagreement



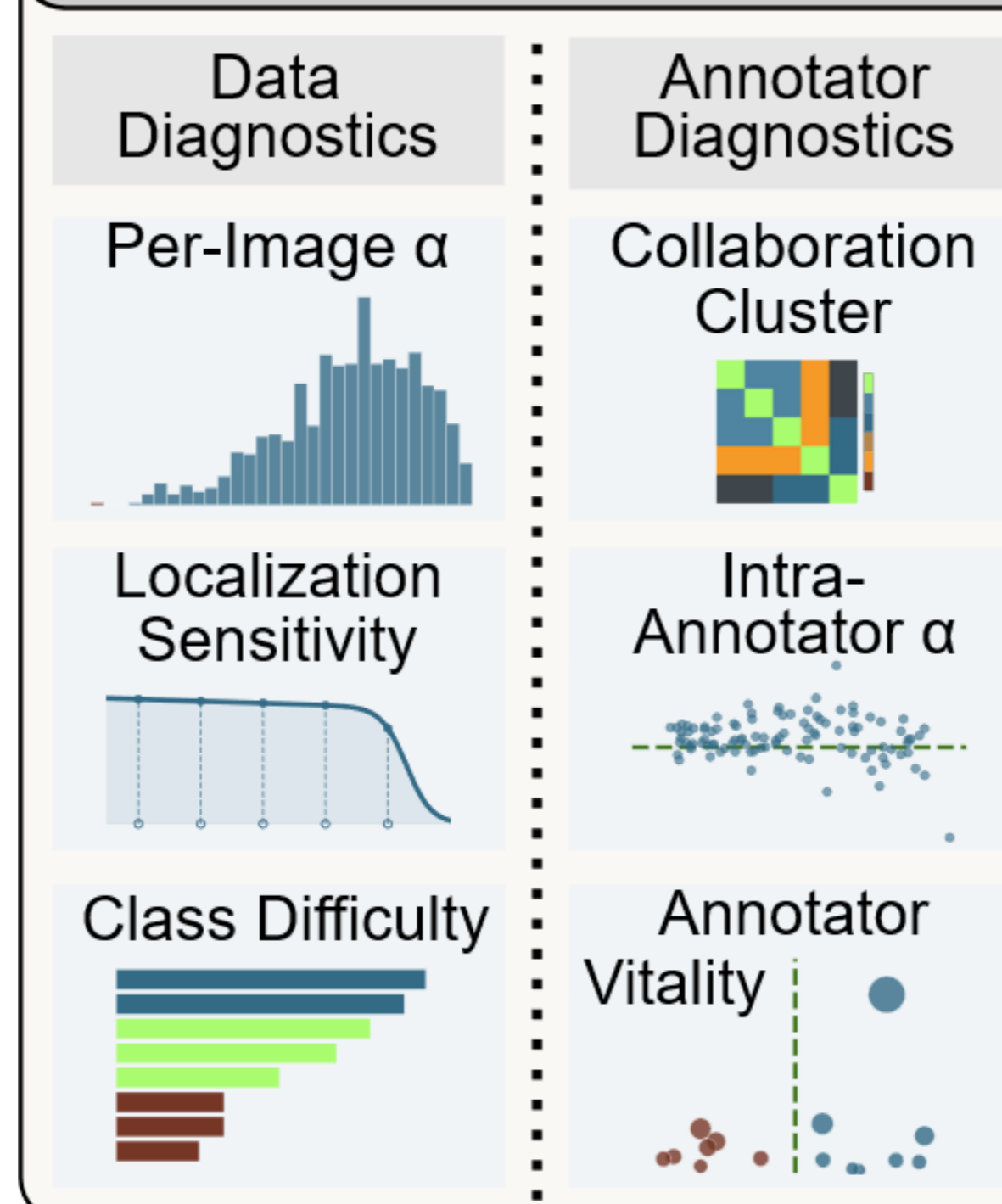
### 2. Principled Configuration



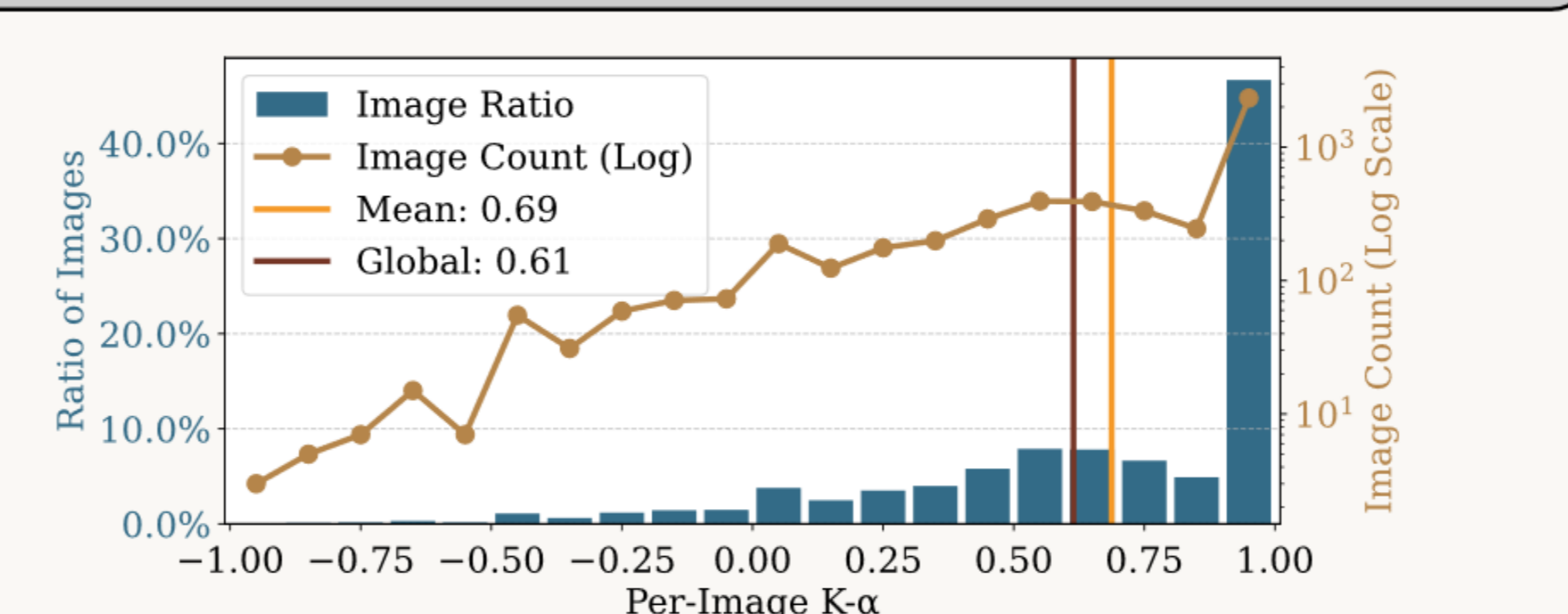
### 3. KaLOS Execution



### 4. Downstream Analysis



### 5. Example: Per-Image alpha on LVIS

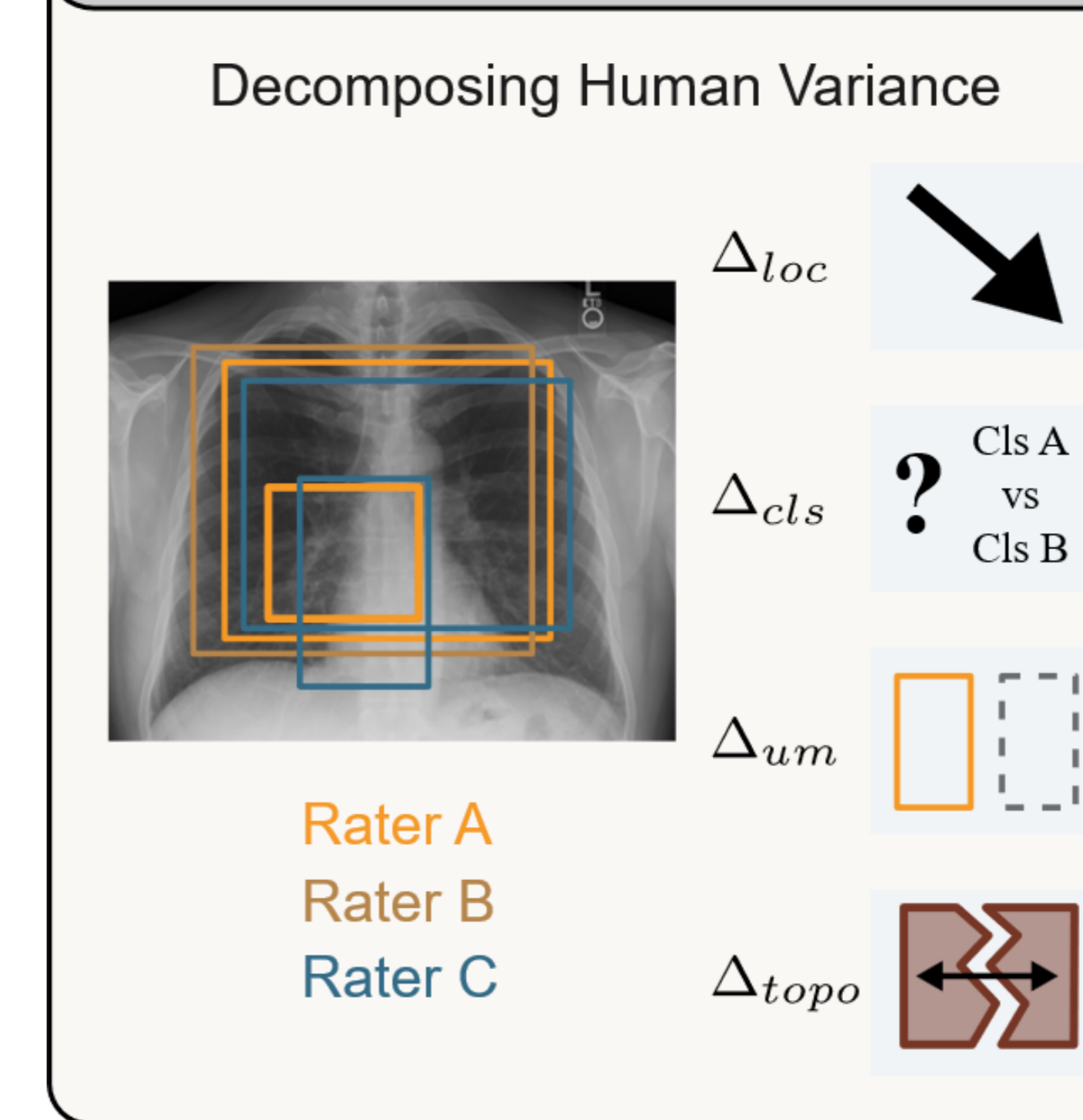


#### Beyond a Single Score (Diagnostic Pipeline):

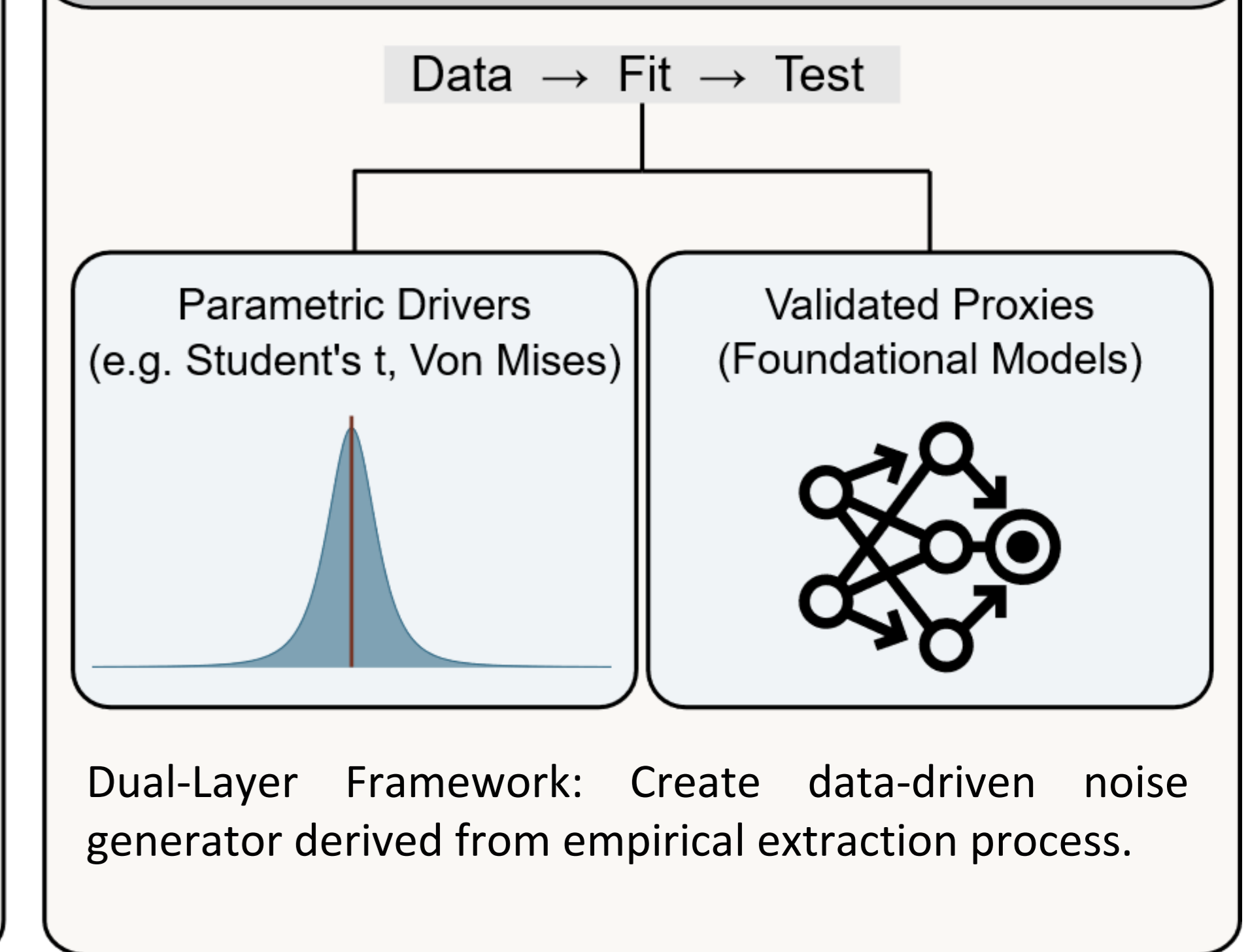
- Guideline Audit: Identify ambiguous classes.
- Rater Vitality: Measure individual consistency.
- Precision Limits: Quantify spatial sensitivity.

## III. Method Validation

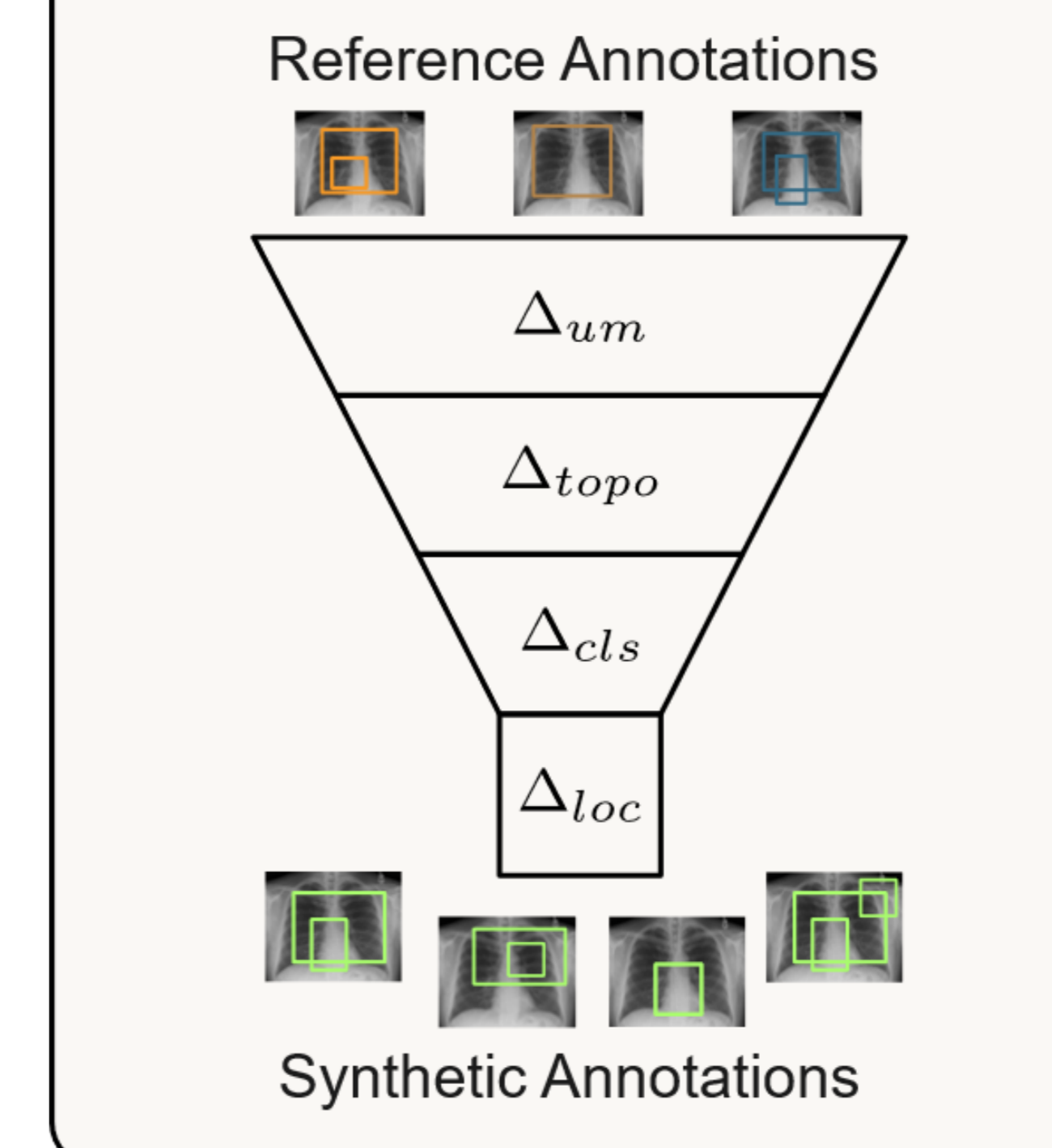
### 1. Error Extraction



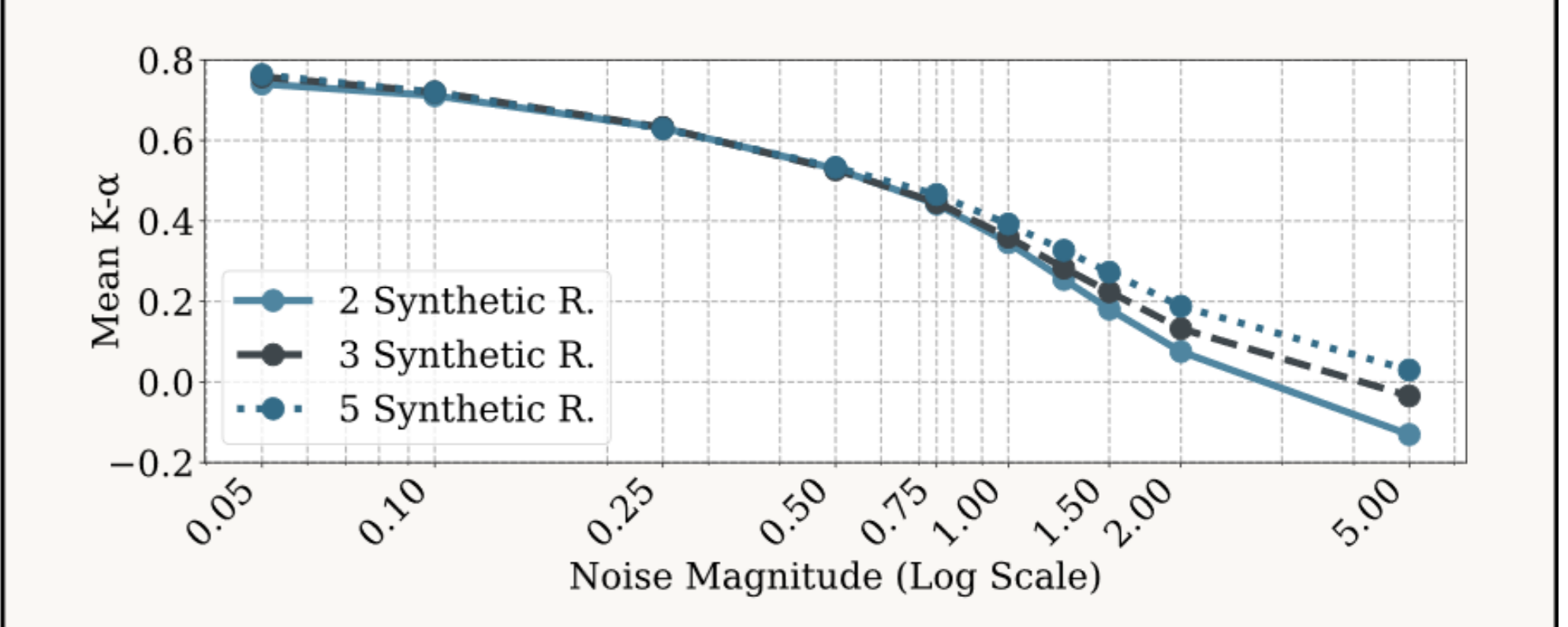
### 2. Noise Modelling



### 3. Generator Composition



### 4. Evidence Based on Synthetic Test-Suit



We validate the core properties of KaLOS (monotonic decrease, full-range utilization, and no sudden jumps) on a synthetic testbed with increasing noise magnitudes, based on an empirically derived noise generator.