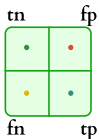


## (A) Inputs (Audits)

### Counts audits



Partitions:  
sex, race, sex x race  
Datasets/time windows  
vary → Classification  
rates

### Metric-only audits

$\Delta$  (DP, EO)

AIR (log ratio)

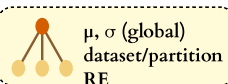
EOdds

CI / SE

Transforms applied  
scaled-logit( $\Delta$ ) for difference  
→ Measurement model

## (B) Bayesian Meta-Evaluator

### Rate-level hierarchy (Partial pooling)



$\mu, \sigma$  (global)  
dataset/partition  
RE

$p, \theta^{\text{TPR}}, \theta^{\text{FPR}}$

### Measurement model (metric-only audits)



Normal( $\mu, \sigma^2$ )  
on transforms

Scaled-logit for  
 $\Delta$  log for AIR  
ratio

### Lattice coherence (coarse ↔ fine subgroups)



Coarse (sex)  
Fine (sex×race)

Soft constraints  
weights  $W_{\cdot|g}$

## (C) Decision-Ready Outputs

### Interval-valued disparities

DP [-0.08, 0.12]

EO [-0.03, 0.15]

AIR [0.76, 0.94]

### Policy-risk $r(\delta)$

Decision thresholds (e.g.,  $\delta = 0.8$ )

Deploy

Defer

Mitigate

$P(\Delta < \delta \mid \text{data})$

### Heterogeneity ( $I^2$ )

Variance decomposition + LOAO



$I^2 = 0.64$

(moderate hetero.)

### Vol: rank next audits

Audit type A

$\Delta W / \text{cost (EVSI)}$  Audit type B