# Trust but Verify: Archaeologist

Sahil Bhatia

# Based on: Sample and Seek

1. **AQP** trades off **accuracy** for **speed** in data analysis
   a. Generates approximate answers to queries using sampling or statistical techniques
2. Issues with existing AQP:
   a. Precision metrics
      i. Confidence Intervals : focus on estimating parameters within each group or subset but may not reveal the overall distribution errors
   b. Unbounded Errors
      i. CI Width $\propto$ std(Sm) / $\sqrt{m}$.

# Sample and Seek

**Support AQP with a user-specified error bound**

# Sample and Seek

1. Distributional Precision
   a. Measures overall precision over all the groups
2. **L2 distance** between normalized distributions of the **approximate answer** and the **exact one**.
3. System guarantees to produce approximations bounded by user-specified error ε.

# Sample and Seek

1. Measure-based Sampling
   a. Rows with higher measure attribute values are more likely to be included in the sample.
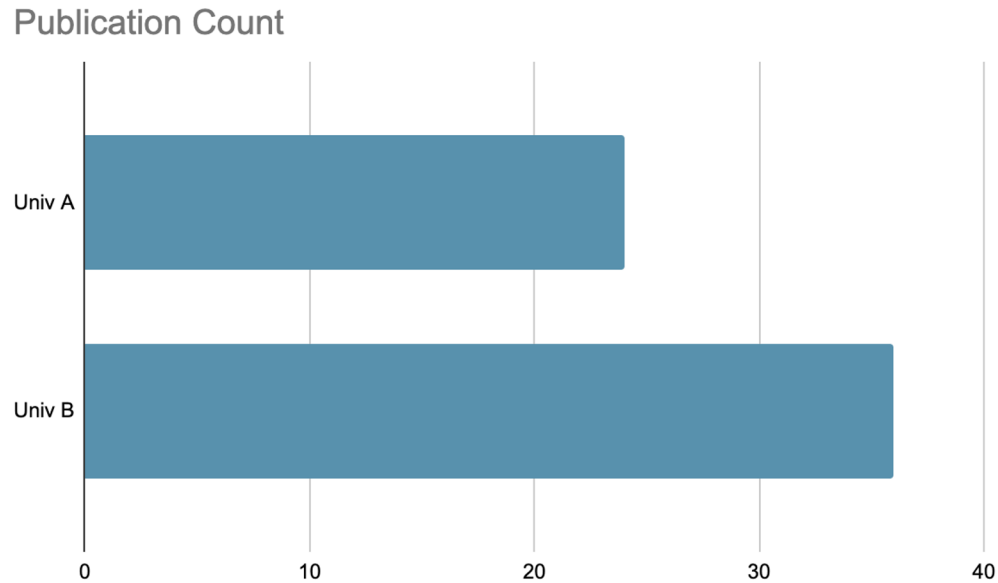   b. A($25), B($20), C($10)
2. Indexes for selective predicates
   a. measure -augmented inverted index - maps values to their respective rews
   b. low- frequency group index - identify rows that belong to low-frequency groups and store them sequentially on disk

# Influenced: ProReveal: Progressive Visual Analytics With Safeguards

1. Progressive Visual Analysis allows access to partial result in middle of computations
   a. Infeasible to compute precise results
2. Not all systems guarantees how **long** they have to wait for **exact results**
3. Users can take a decision using partial results
   a. Account for worst case

# ProReveal

1. Represent intermediate knowledge as as guards (logical formulas).

**Publication Count**



pubCount(Univ A) < pubCount(Univ B)

# ProReveal

1. Represent intermediate knowledge as as guards (logical formulas).
2. System continuously gives feedback on validity
   a. If invalidated systems notifies user

# ProReveal

1. Moritz et al.'s research on optimistic visualization [9] is one of the studies that motivated this research.
2. Differences:
   a. Precise results is not obtainable
   b. Continuous feedback on validity of guard
   c. Intermediate knowledge can be represented structurally

# ProReveal

$$\langle \mathit{PVA\text{-}Guard} \rangle := \langle \mathit{variable} \rangle \; \langle \mathit{operator} \rangle \; \langle \mathit{operand} \rangle$$

**where**

$$\langle \mathit{operand} \rangle := \; \mathit{empty} \; | \; \langle \mathit{variable} \rangle \; | \; \langle \mathit{constant} \rangle.$$

1. Variables can be single value (cell of heatmap) or even a distribution of values

# ProReveal