# Online Aggregation
# 10-Year Feedback

Joseph M. Hellerstein (Berkeley)
Peter J. Haas (IBM)
Helen J. Wang (Berkeley → MSR)

# A Demo is Worth a Thousand Words

# A Demo is Worth a Thousand Words

# The OA Back-Story

Jeff Naughton
UW Madison

Peter, Jeff, & S. Seshadri:
Fixed-precision estimation
For COUNT queries over joins

Joe and Jeff:
"Sloppy Databases"
(email exchange)

Helen comes
to Berkeley

SIGMOD '95:
Keynote + hallway chat

SIGMOD '97
intro paper

Ripple join (SIGMOD '99)
CONTROL project

# A Three-Way Synthesis

**Systems**

**Statistics**

Index Striding

Ripple Joins

Online Reordering

Approximate Answers
Confidence Intervals
- CLT approximate intervals
- Hoeffding conservative intervals

Over relational ops

OA

Continuous interaction

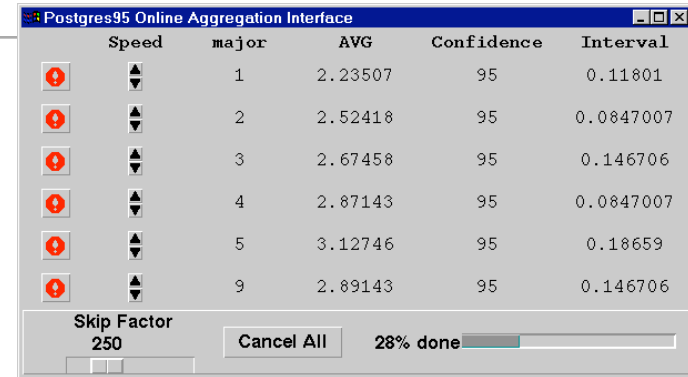Multiresolution visualization

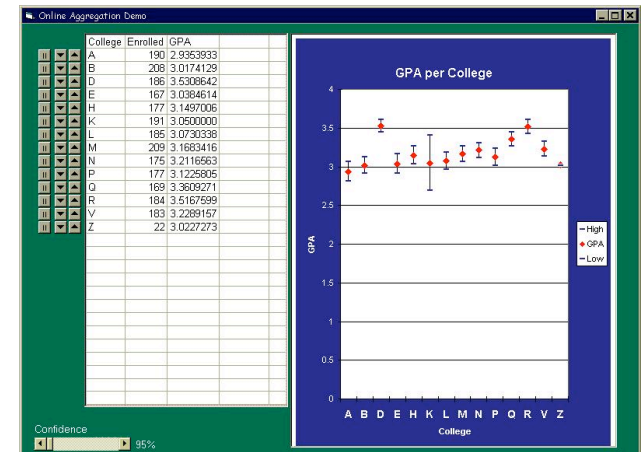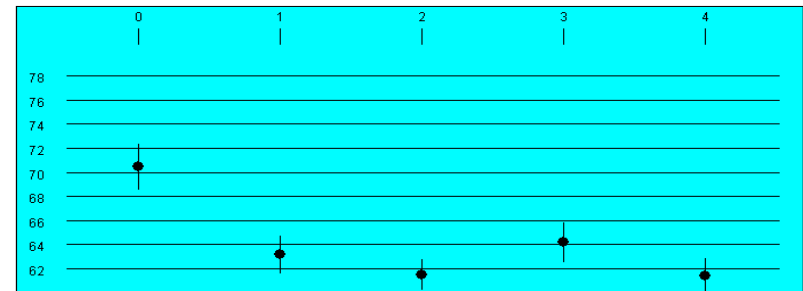Computation follows focus

**HCI**

# Implementations

- ## Postgres
  - Index stride, ripple join in engine
  - Simple Tcl/Tk interface

- ## DB2
  - JDBC client application (differential pacing)
  - Java swing interface

- ## Informix Universal Server prototype
  - Index stride, ripple join, online reordering
  - Integrated with Metacube OLAP tool
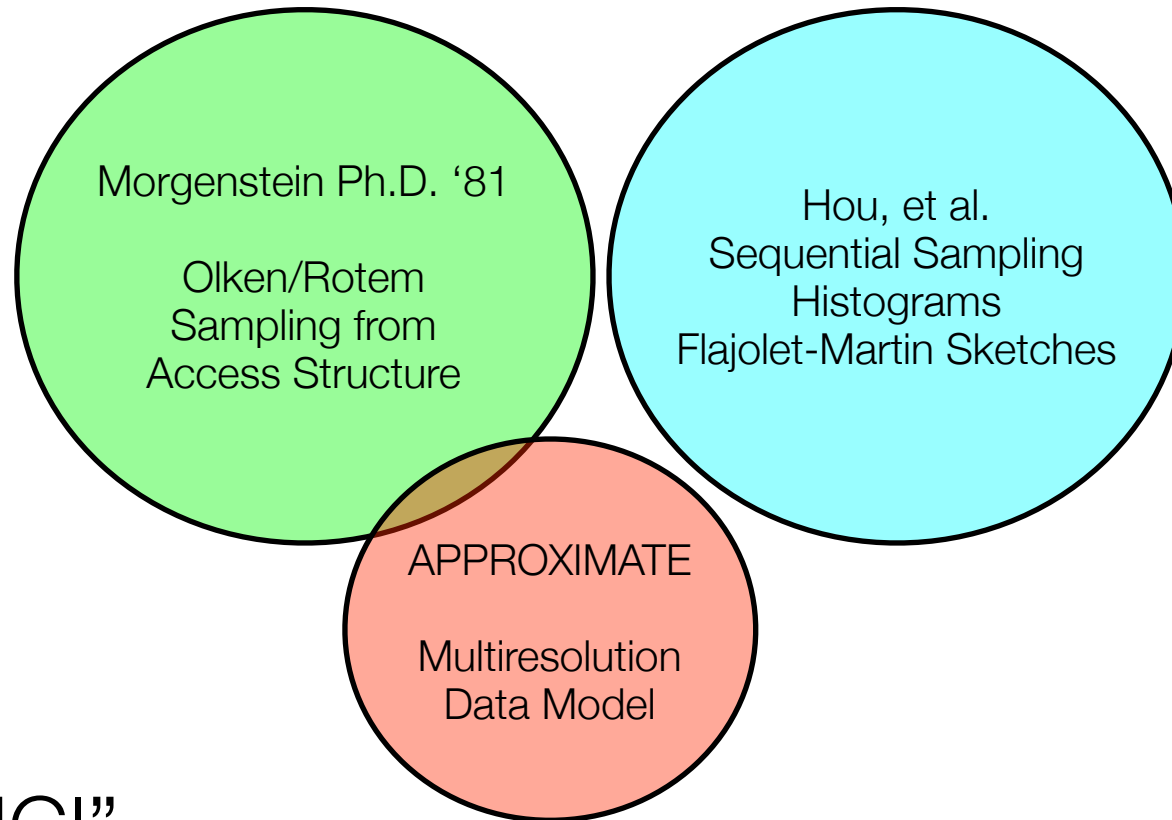  - Visual-basic + Excel interface

# Before

Morgenstein Ph.D. '81

Olken/Rotem
Sampling from
Access Structure

Hou, et al.
Sequential Sampling
Histograms
Flajolet-Martin Sketches

APPROXIMATE
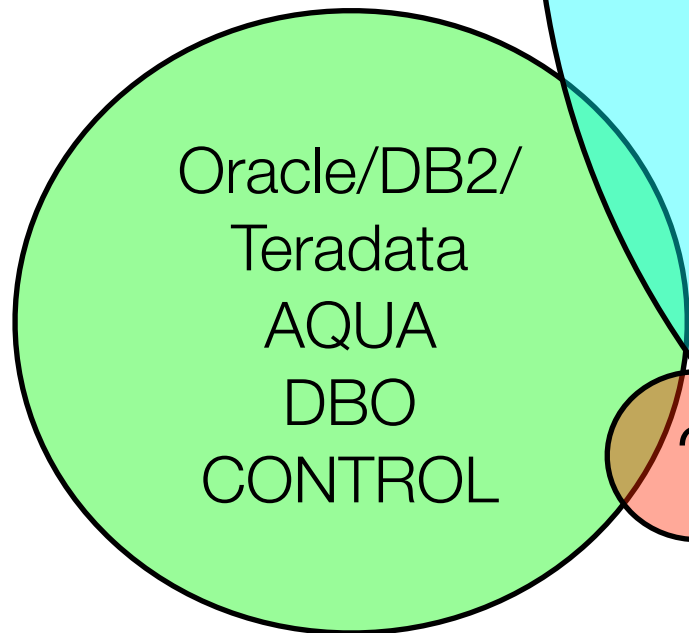
Multiresolution
Data Model

"HCI"

# The Last Decade

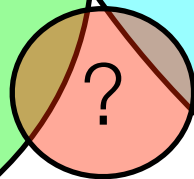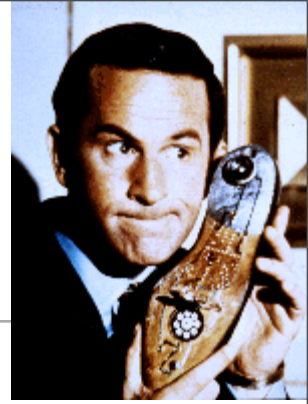## Statistics

## Systems

Oracle/DB2/
Teradata
AQUA
DBO
CONTROL

**?**

**SQL sampling standard**
**Synopses/Sketches Galore**
AMS, wavelets, DCT, DV, samples …
**"Robust" sampling**
exploit indexes, workload, precomputation
**Stream samples & sketches**
General and special purpose
**The Florida Renaissance**
Large-sample maintenance
Scalable online joins
Monte Carlo methods (bootstrap)
**Probabilistic DB**
… and so on!

## HCI

# The CONTROL Project

- Online Aggregation

    - This paper, Ripple Joins, new confidence intervals
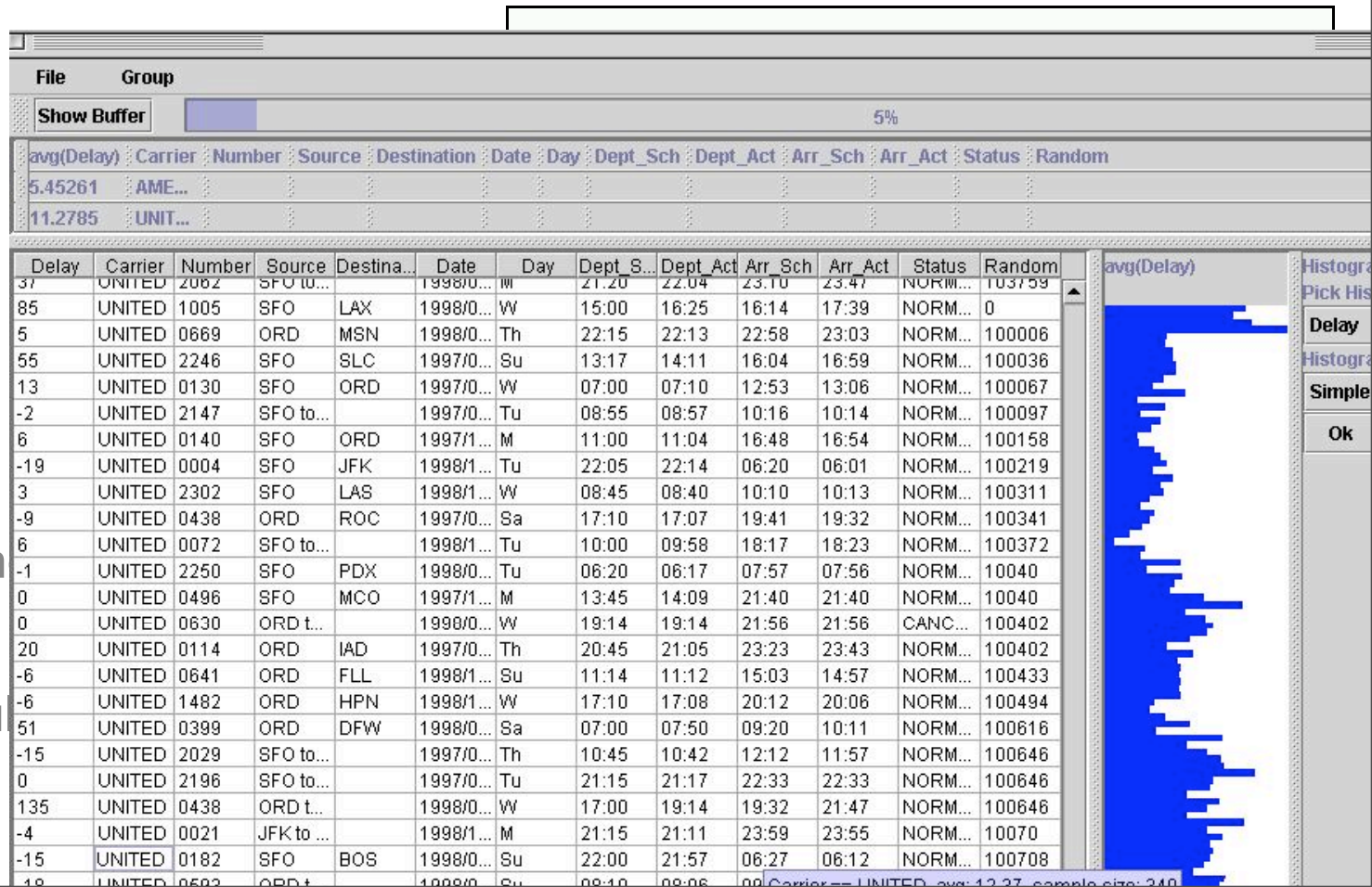
- Clouds

- CARMA

- Online "Enumeration"

    - Scalable Spreadsh...

    - Partial Query Resul...

# If you're so smart, why ain't you rich?

- Marketplace Challenge

  - Apps + Engine

  - Customer aversion to statistics (see no evil)

  - OLAP postponed this by a decade

  - If you want to rewrite the DB2 engine, you better have a VERY good business case
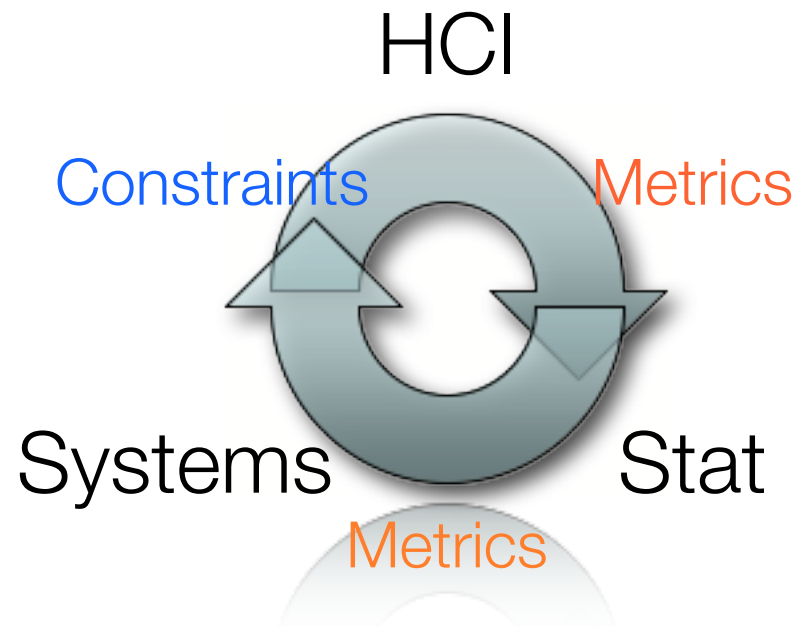
- Technical Challenge

  - How many DBMS engineers does it take to screw in a confidence interval?

# Why the Renaissance?

- The rise of stat in CS

  - KDD, ML, WebSearch

- Market forces

  - DBMS Market Consolidation

  - Software Appliances/Services

  - New "niche" opportunities

- Web expectations

  - Speed, data-rich, rough-and-ready answers

  - App/Engine integration not a barrier

# What Next?

- Clearly, tons of community energy on approximation

  - Especially algorithmics

- The fun (for us) is in the integrative work

HCI

Constraints          Metrics

Systems                    Stat

Metrics

"What unlike thing must meet and mate"     -- Melville

# With Thanks…

Joe:

Jeff Naughton
Mike Stonebraker

The CONTROL Freaks:
Ron Avnur
Andy Chou
Christian Hidber
Bruce Lo
Chris Olston
Vijayshankar Raman
Tali Roth

Peter:

Jeff Naughton
Pat Selinger
Bill Cody
S. Seshadri
Ashutosh Singh
Guy Lohman
Vijayshankar Raman
Gang Luo