

# MCDB and SimSQL: Scalable Stochastic Analysis within the Database

Peter Haas (IBM) and Chris Jermaine (Rice University)

People who have worked on MCDB/SimSQL: Subi Arumugam, Zhuhua Cai, Jacob Gao, Shangyu Luo, Ravi Jampani, Luis Perez, Foula Vagena, Mingxi Wu, Fei Xu

## MCDB/SimSQL: What Is It?

- Complete analytic SQL DB system
- Runs on top of Hadoop
- Scales easily to 100-machine clusters
- Handles terabytes of data
- Special support for *stochastic analytics*

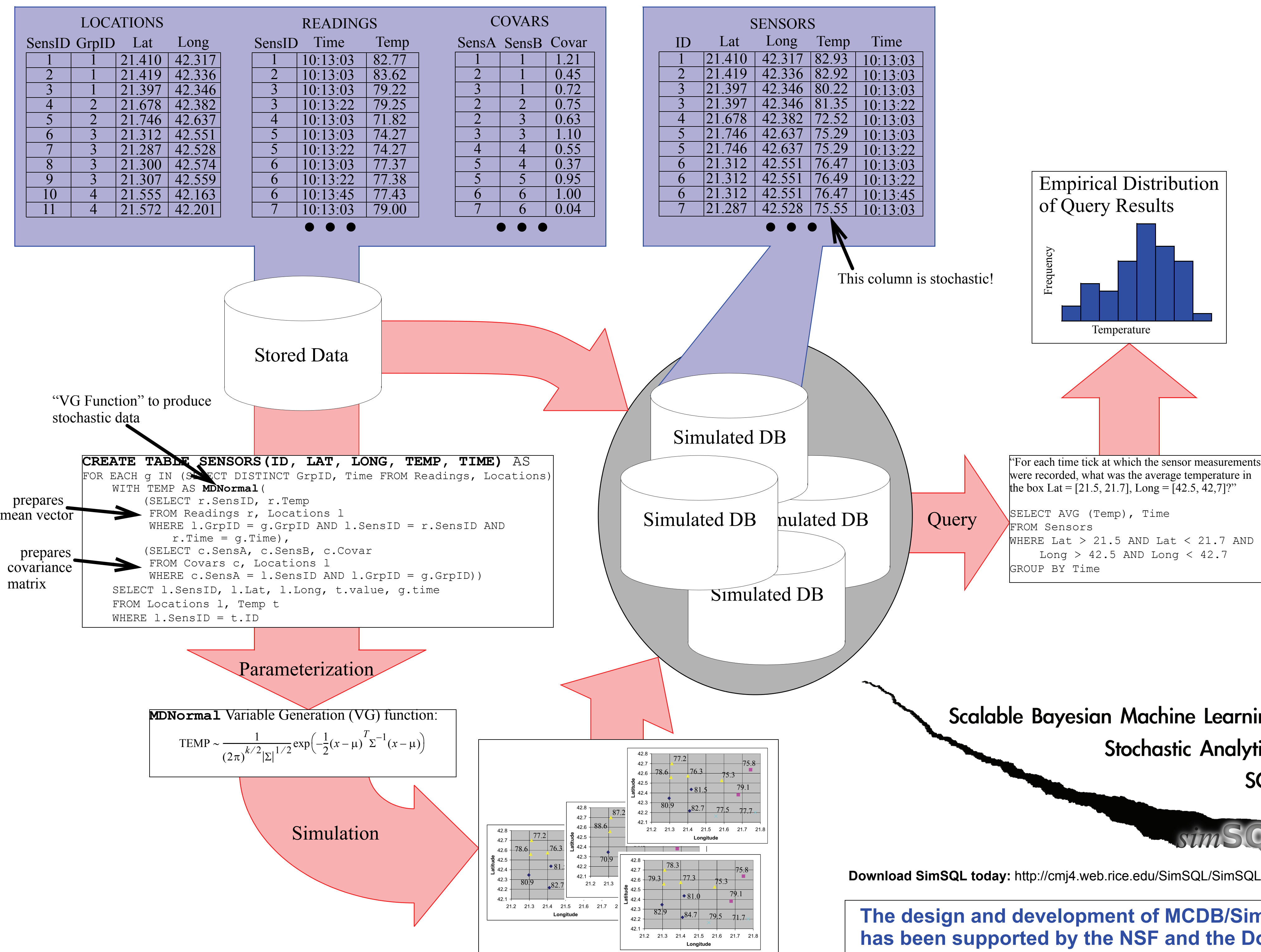
## What Does “Stochastic Analytics” Mean?

- Data analytics based on stochastic simulation
- Model missing/uncertain data via prob. dists
- System uses prob. dists to simulate DB instances
- Can then query those instances via SQL
- Complicated dists supported, even Markov chains

## Special Features of MCDB/SimSQL

- Unique randomized table functions (“VG functions”)
- Multiple concurrent simulations for the price of one
- Special support for recursion (MCMC and Bayesian ML)
- Support for incremental execution of very large plans
- Coming soon: first class vectors and matrices

Simple example: Writing SQL queries that take into account errors on sensor readings.



MDNormal Variable Generation (VG) function:

$$TEMP \sim \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

