

# From MUD to MIRE: Managing Inherent Risk in the Enterprise

Peter J. Haas

IBM Almaden Research Center  
San Jose, CA



# The Two Perpetual Questions

- “Where do the probabilities come from?”
- “Who is going to use this stuff in the real world?”



# My background in probabilistic DB

# RAQA: Resolution-Aware Query Answering for Business Intelligence

(Sismanis et al. 2009)

- OLAP querying (datacubes: roll-up, drill-down)
- Uncertainty due to entity resolution
- **Bounds** on query answers
- Implemented via SQL queries
- Conservative approach

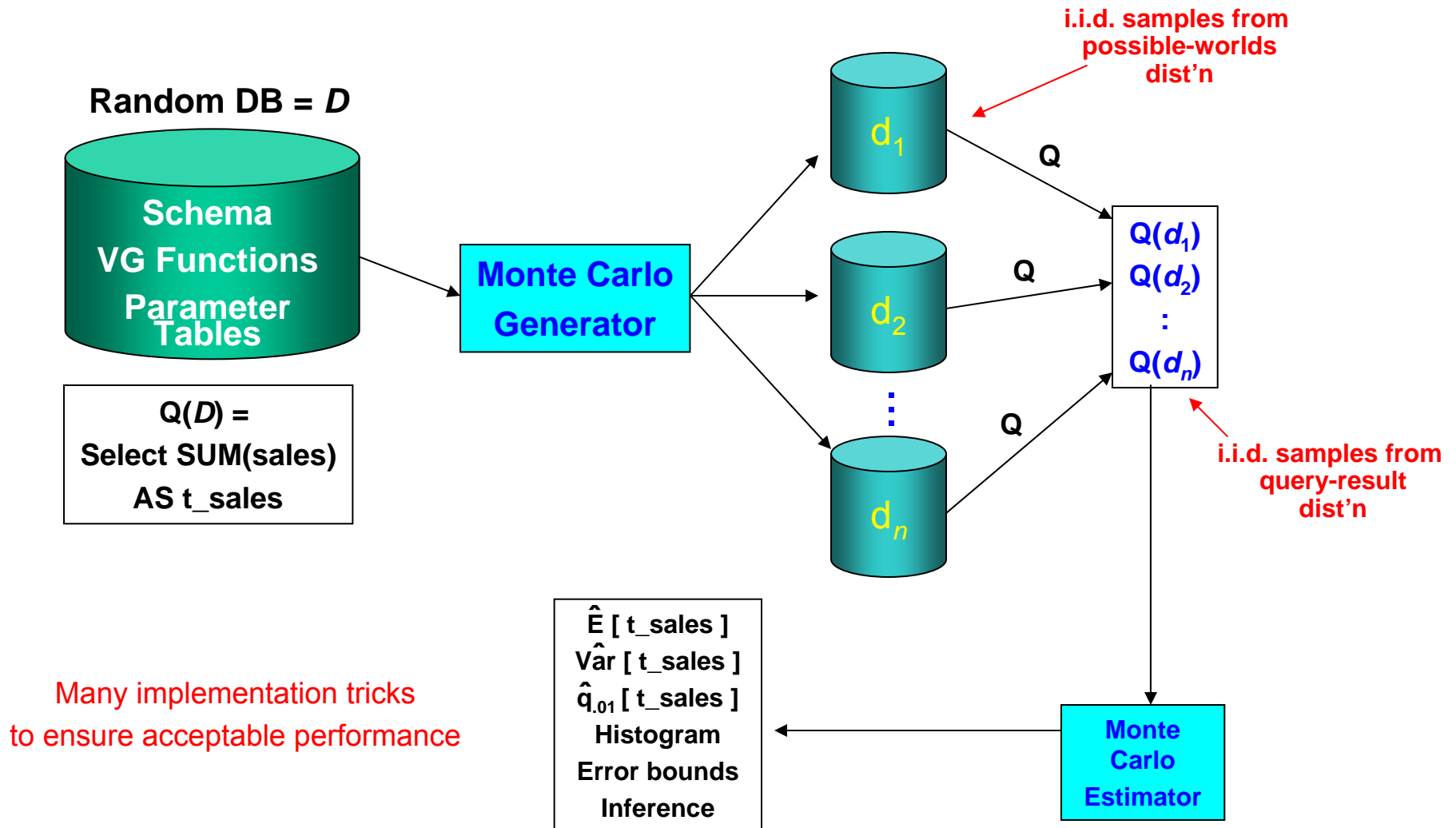
City	State	Strict range	Status
San Francisco	CA	[\$30,\$230]	guaranteed
San Jose	CA	[\$70,\$200]	non-guaranteed

Sum(Sales) group by City,State

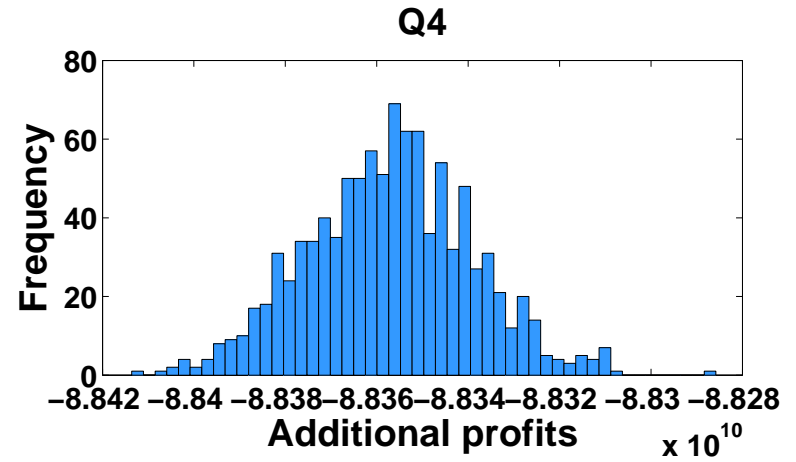
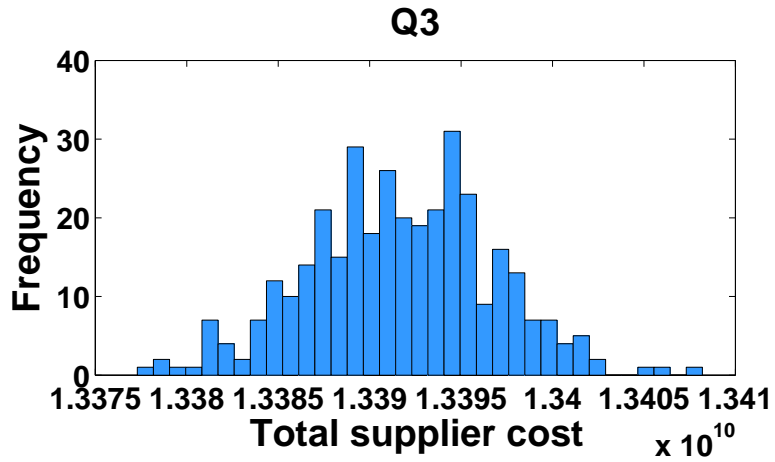
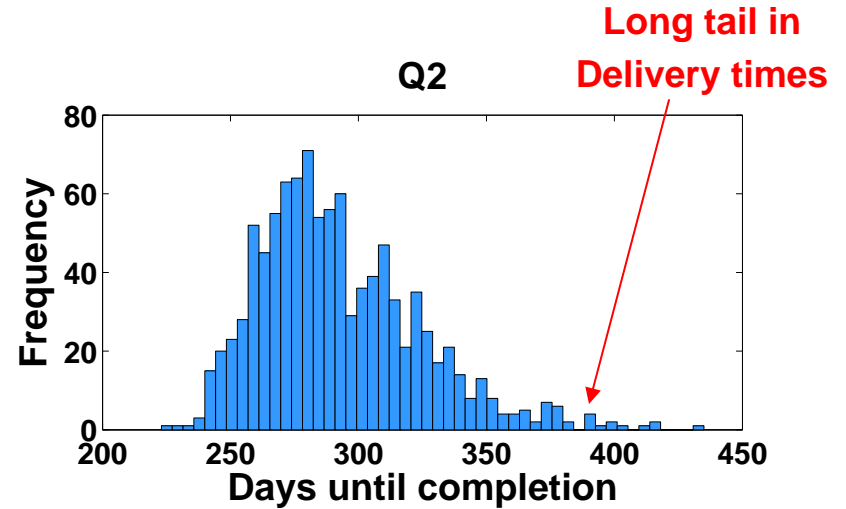
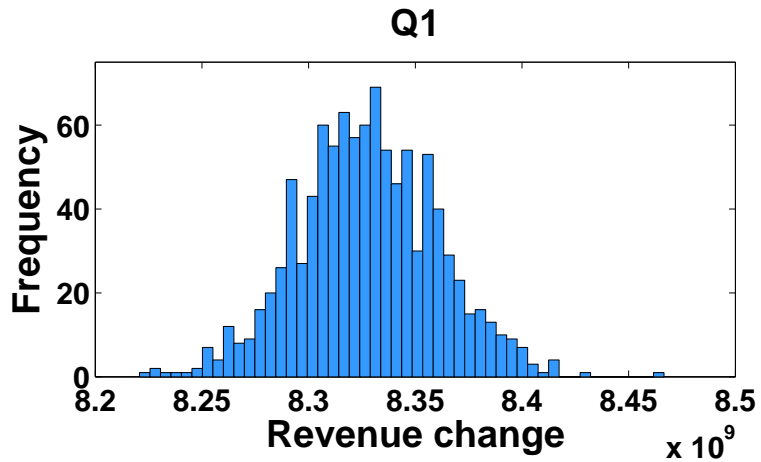
State	Strict range	Status
CA	[\$230,\$230]	guaranteed

Sum(Sales) group by State

# The MCDB System (with Chris Jermaine & students)

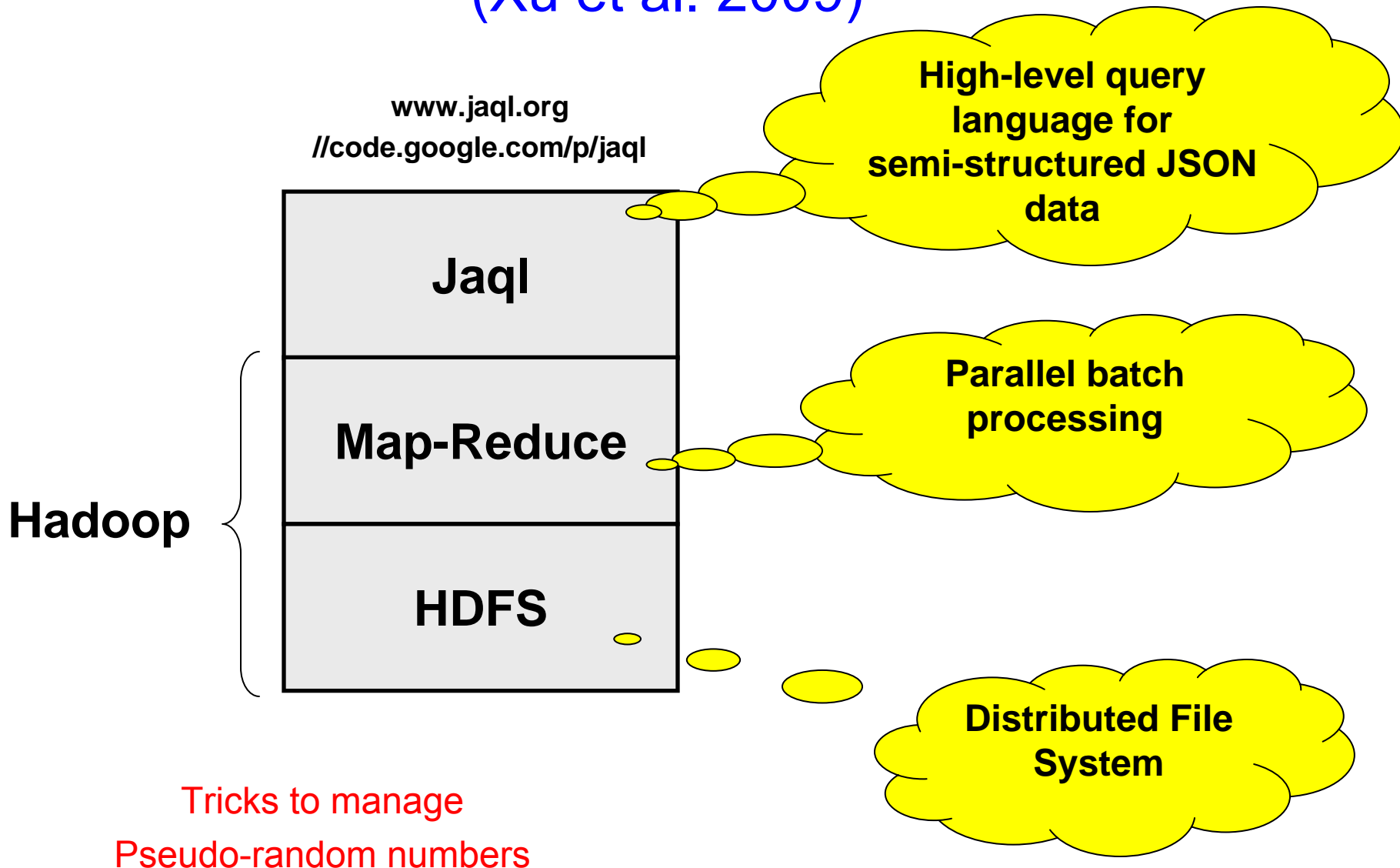


# Query-Result Distributions



# MC<sup>3</sup>: MapReduce + MCDB

(Xu et al. 2009)



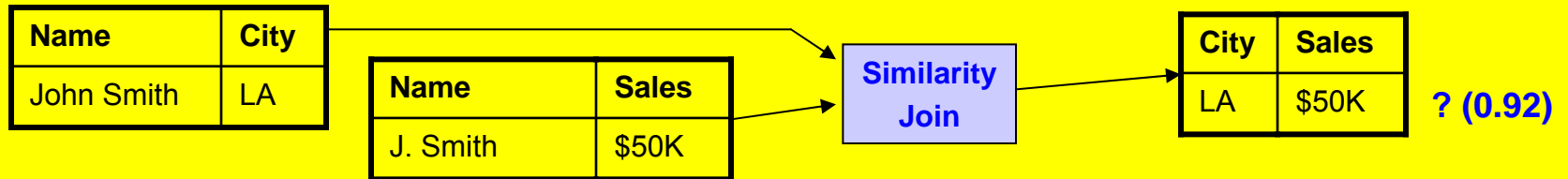
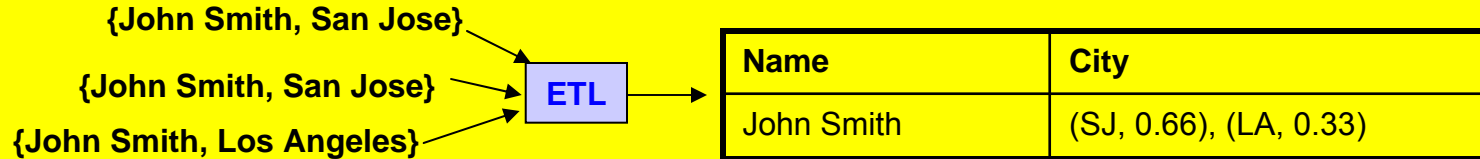
Where do the probabilities come from?



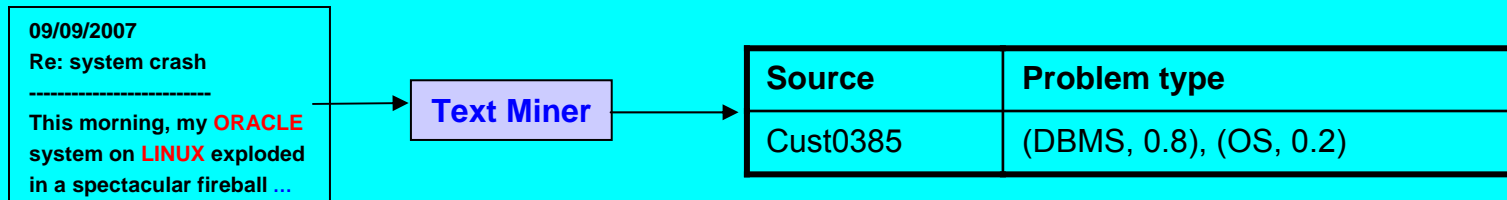


# Data-Warehouse Uncertainty

## Data Integration

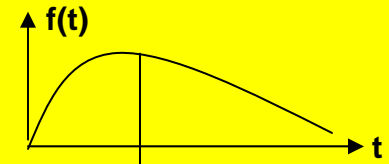
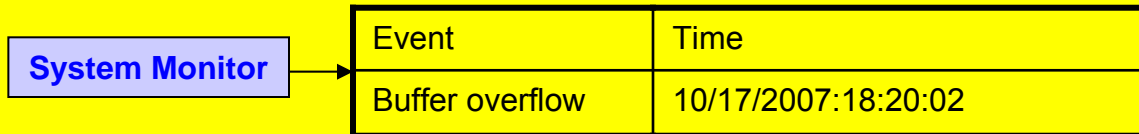
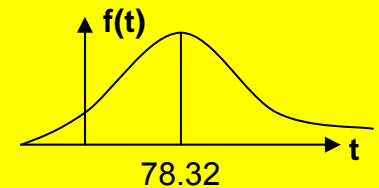
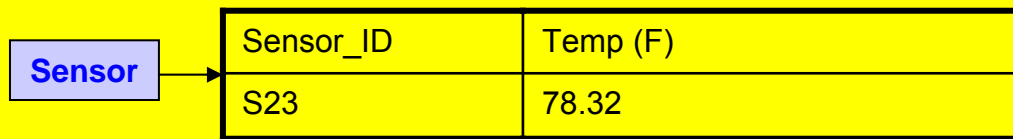


## Information extraction



# Data-Warehouse Uncertainty – Cont'd

## Measurement Uncertainty



# Real-World Challenges with Data-Warehouse Uncertainty

- People don't like to admit that it exists!
  - **Retailers** view uncertainty as failure of security, supply chain management
    - IBM research relationship manager for retail
  - **Law enforcement**
    - Photo ID in meth dealer trial
  - **Scientists** pretend data is perfect: uncertainty undermines results
    - Hans-Joachim Lenz
  - **Database vendors**
    - Data “cleaning” products
- Data warehouse may not even exist!
  - Ex: cancer data at medical center
  - Ex: tomato soup supply chain data



# Stochastic Predictive Analytics on Big Data

- Uncertain data describes **future** or **hypothetical** events
  - Based on complex, fine-grained stochastic model over big data
  - Minimizes denial problem
- Intense recent interest in “business analytics” driven by
  - Need for low risk, quick payback projects (flexibility, low cost, fine data granularity)
  - Technical advances
    - Cloud computing
    - Software as a Service (SaaS)
    - Next generation tools, portals, visualization
- Often with a spreadsheet front end
  - \$8 Billion of such tools [Gnatovich06]
  - IBM services pricing
- Lots of prototype activity
  - Fox/GreenPlum [Cohen09 MAD analytics paper]
  - VISA/IBM [Das10 SIGMOD paper]



# Ex. 1: Portfolio Values

**Customer**

CustID	OptionID	NumShares	...
John Smith	23	50	...
...	...	...	...

**EuroCallOptions**

OptionID	InitVal	...	StrikeP	OVal
23	\$2.35	...	\$4.00	?
...	...	...	...	...

```
SELECT SUM (c.NumShares * o.Val)
FROM Customer c, EuroCallOptions o
WHERE c.OptionID = o.OptionID
      AND c.CustType = 'Institutional'
```

Option value  
one month from now  
(exercise date)

Modified Black-Scholes model for European call option:

$$dV = rVdt + (a\sqrt{V})VdW \quad OVal = \max(V(t_{\text{final}}) - S, 0)$$

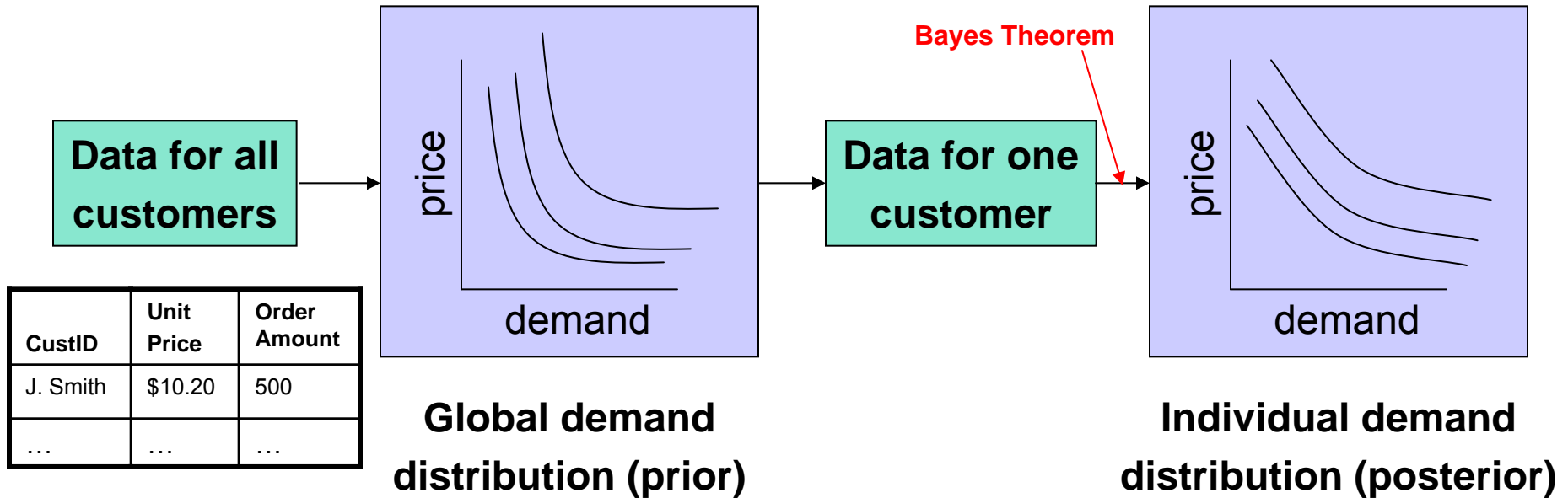
Simulation approximation (Euler approach):

$$V(t + \Delta t) = V(t) + rV(t)\Delta t + (a\sqrt{V(t)})V(t)\sqrt{\Delta t}Z_j$$

Sample from  
Normal dist'n

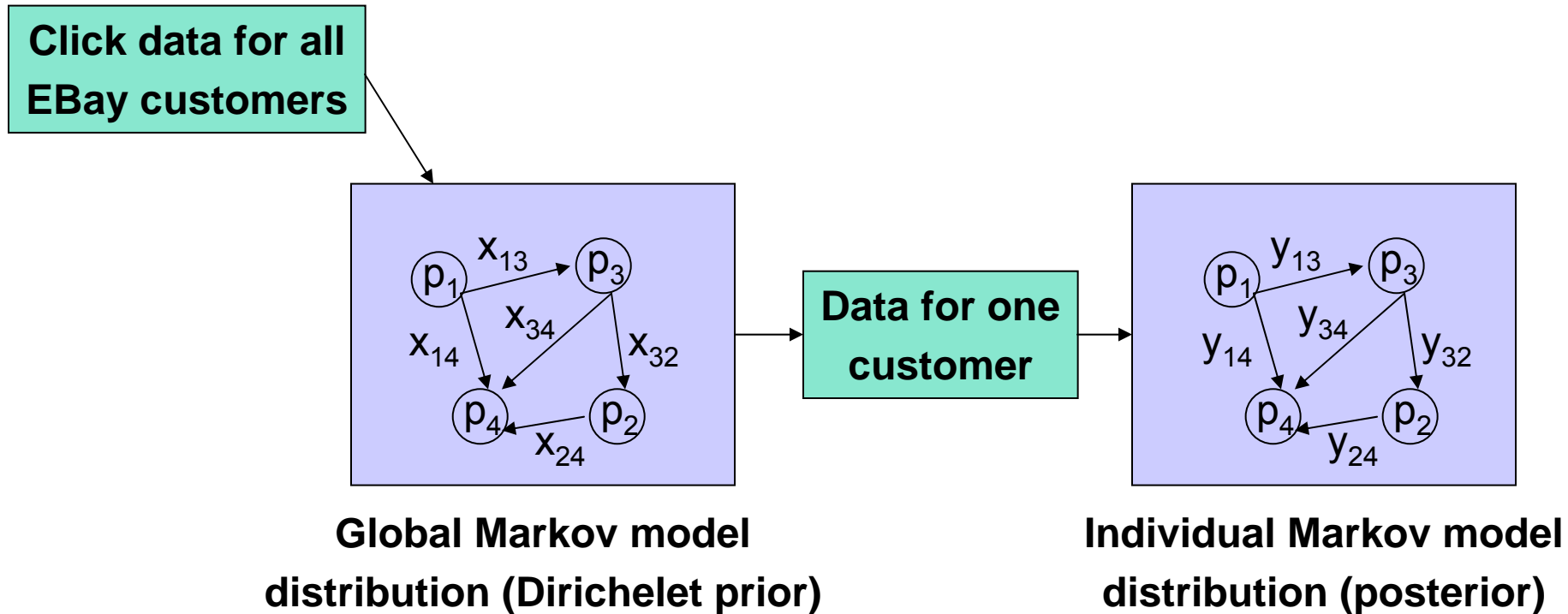
Also CMOs, etc.

# Ex. 2: Pricing Decisions



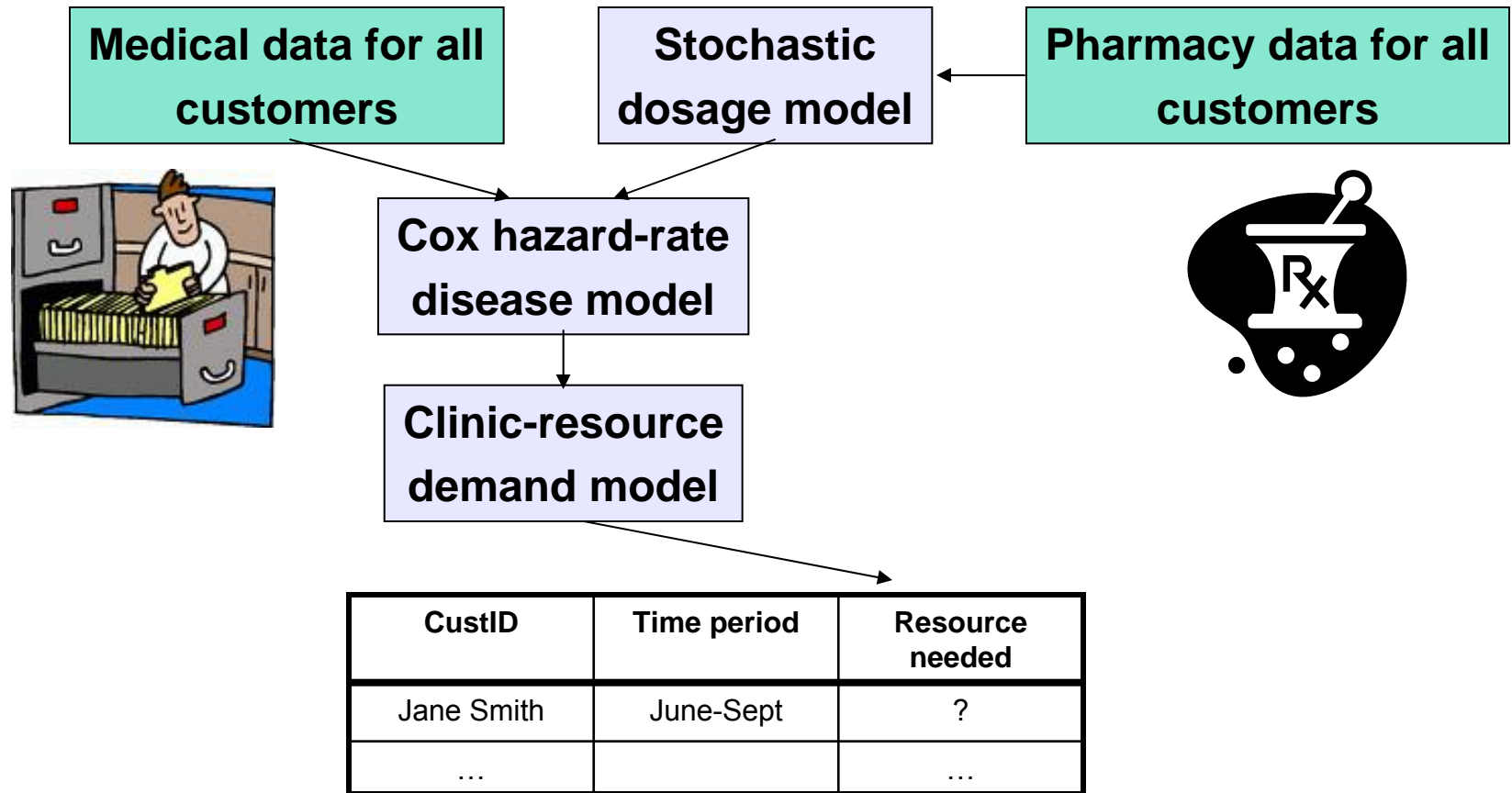
- Can analyze arbitrary dynamically-defined customer segments when determining effect of price increase

# Ex. 3: Individual Click Behavior (EBay)



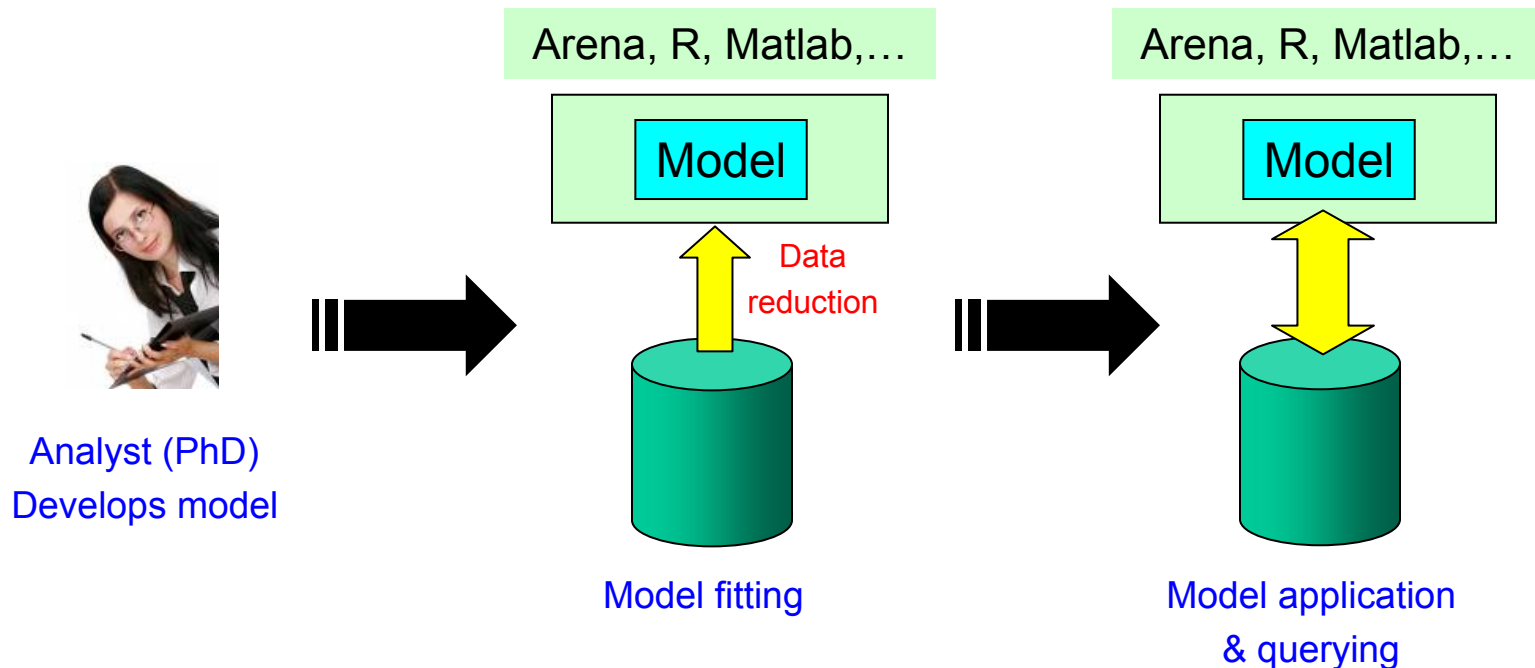
- Can analyze **arbitrary dynamic** customer segments when determining effect of changing EBay pages

# Ex. 4: Clinic-Capacity Risk



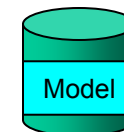


# MCDB: Improvement of Traditional Analytics Workflow



- Data extraction slow and bug-prone
- Only coarse-grained modeling
- No encapsulation for user
- Hard to re-link model results to DB
- Hard to deal with data updates
- Sensitivity, what-if analysis are hard

Goal: Integrate model with Database



# Where do the probabilities come from?



From stochastic predictive models over big data

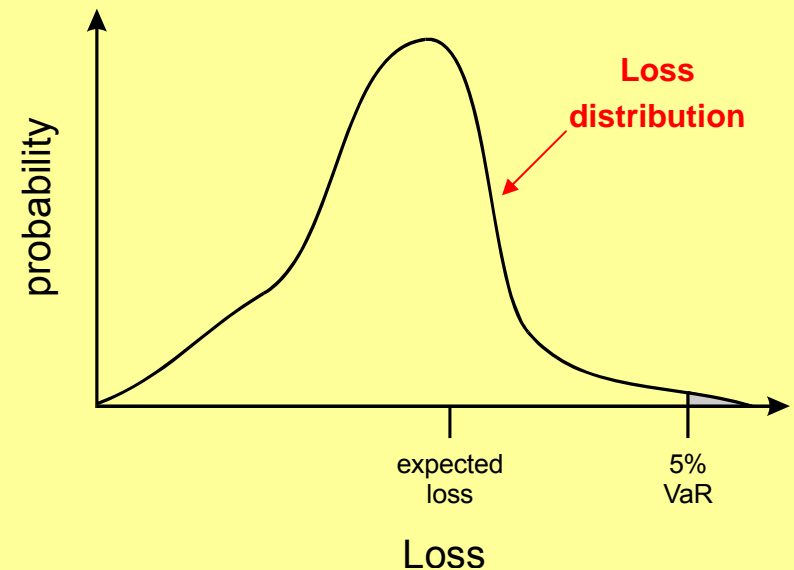
Who is going to use this stuff  
in the real world?



# Key Driver: Risk Management

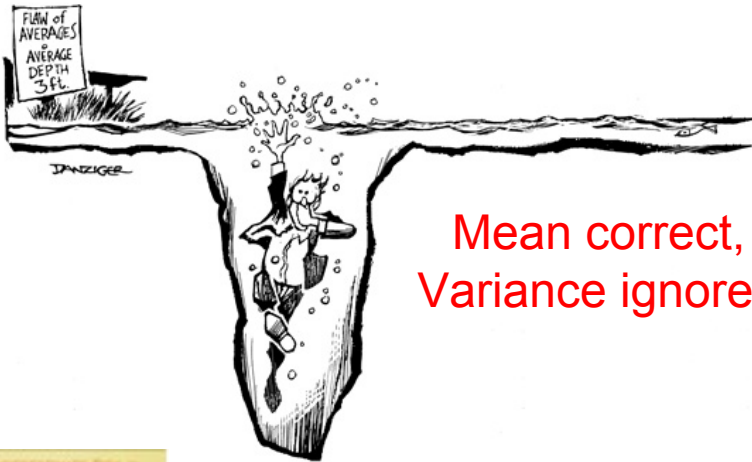
- Ex: Projected sales under micromarketing campaign
- Ex: ERP
  - # OS experts for help desk
  - Demand projected from historical text data (2x uncertainty)
  - Provide principled **safety factor**
- Regulatory pressure
  - Basel II, Solvency II
- Business pressure
  - Ex.: Energy Risk Professionals

```
SELECT SUM (s.amount)
FROM SALES s, CUST c
WHERE s.ID = c.ID
      AND c.city = 'Los Angeles'
```



# Challenge: Decision-makers' Poor Intuition About Risk

Flaw of averages (weak form):

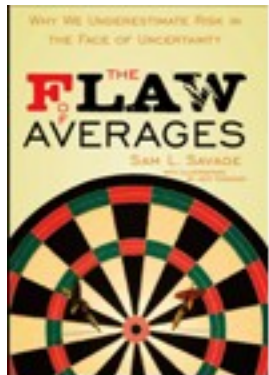


Mean correct,  
Variance ignored

Flaw of averages (strong form):



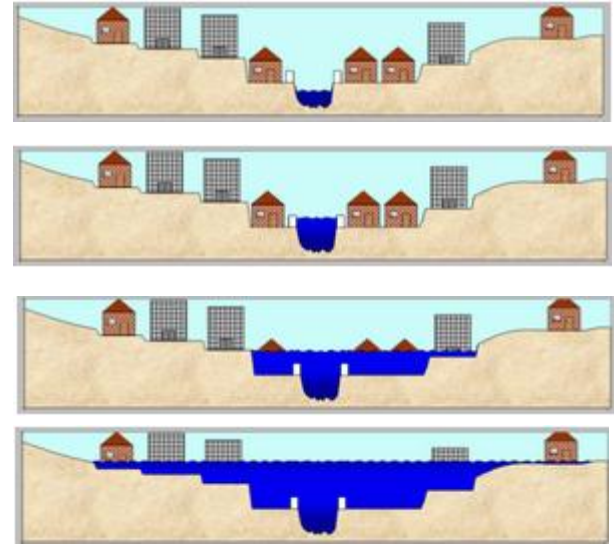
Wrong value of mean:  
 $f(E[X]) \neq E[f(X)]$



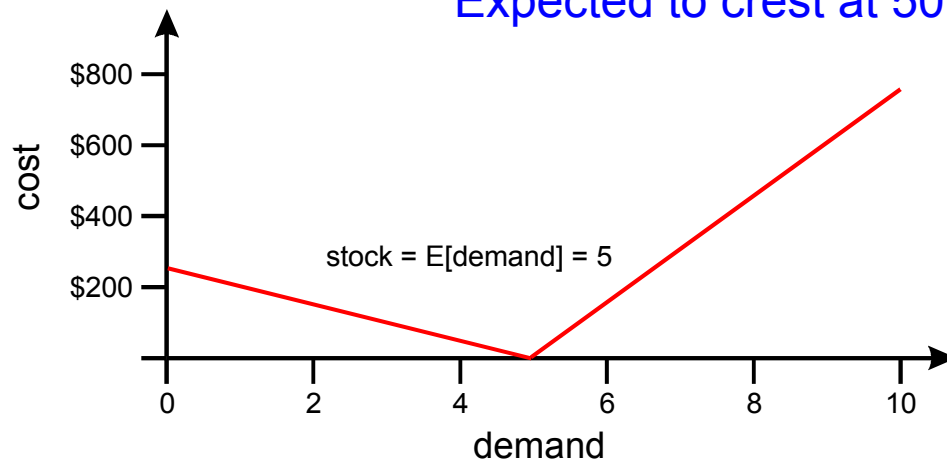
Sam Savage's book  
(why we underestimate risk)

# Examples

- Red River (ND) flooding
- Perishable Inventory (Red Lobster)
- U.S. accounting standards (FASB)
- Project completion time: 10 parallel tasks,  $E[T_i] = 6$  mo.
- Data cleansing
- Machine learning
- Trio agg. paper (MUD 2008)
- Basic probability



“Expected to crest at 50 feet”



# Probability Management and Interactive Spreadsheets

- DIST 1.1 standard
  - DIST = distribution string
  - IID Monte Carlo (multivariate) samples
  - Compressed, with metadata
- Ensures correct, coherent risk computations throughout enterprise and beyond
  - E.g., Royal Dutch Shell
- “Electricity network” for probability
  - Royal Dutch Shell, Merck Pharmaceutical, Oracle, Wells Fargo Bank, Bessemer Trust, and IBM
- DISTs can be manipulated like numbers
  - Facilitates **interactive spreadsheets** (demo)



Audit seal of  
approval

# Demo 1

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	BUSINESS PLAN: Requiring Capital Investment with Uncertain Demand												
2													
3													
4	INVESTMENT	\$1,600,000											
5													
6	CAPACITY	2000000											
7													
8	DEMAND	2507670											
9													
10	PRICE	\$1											
11													
12	SALES	2000000											
13													
14	REVENUE	\$2,000,000											
15													
16	PROFIT	\$400,000											
17	Average	\$200,496											
18	Inputs												
19	Decisions												
20	Certainties												
21	DIST String												
22	Dst(DIST)												
23	Outputs												
24	Formulas												
25													
26													

NOTE: You must activate the model with the Simulate Activate command or the button the the tool bar.

Bins	Frequency	Investment in \$10
<dist name=" Unnamed---		16
-\$1,400,000	0	
-\$1,200,000	0.003	
-\$1,000,000	0.007	
-\$800,000	0.022	
-\$600,000	0.027	
-\$400,000	0.05	
-\$200,000	0.1	
\$0	0.132	
\$200,000	0.659	
\$400,000	0	
\$600,000		

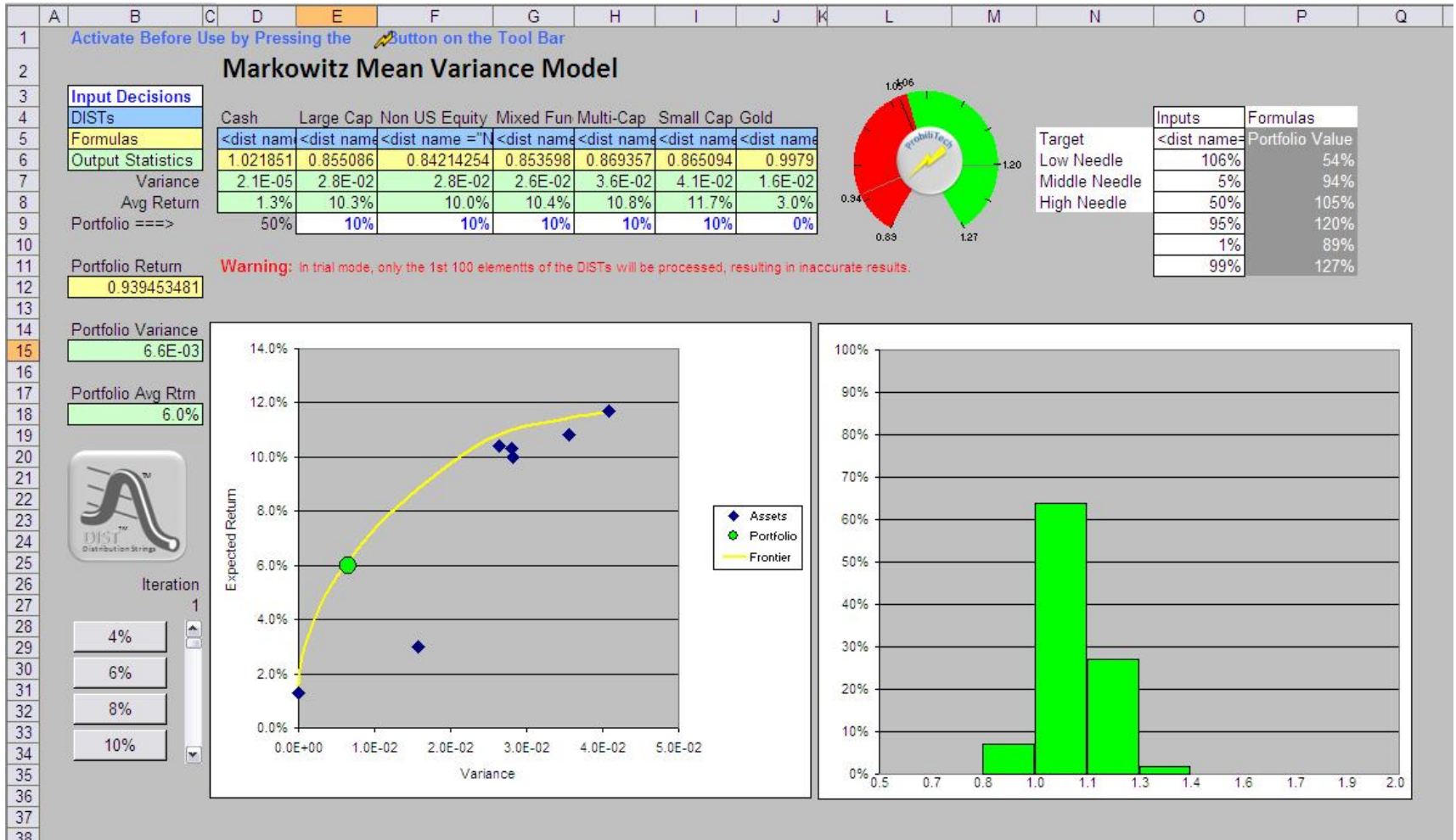
  

DEMAND DIST

<dist name="Demand DIST" avg="2.0000000E+006" min="5.03593565E+005" max="3.4436123E+006" count="1000" type="Single" origin="Normal(2m,0.5m)">rm24frMeuC3qtcv49UKSLmbV4sGamj3vbkv-8iJnryqqDUmmDlzyEeKCQe5mwVYJZo



# Demo 2



# Probability Management and Probabilistic Databases

- ProbDBs can be a source of DISTs
  - Directly from MCDB
  - Can sample from
    - exact distributions
    - approximate empirical distributions
    - Fitted distributions (e.g., compute mean, var)for **aggregation query** or **loss function**
- Greater impact on decision-makers

Who is going to use this stuff  
in the real world?



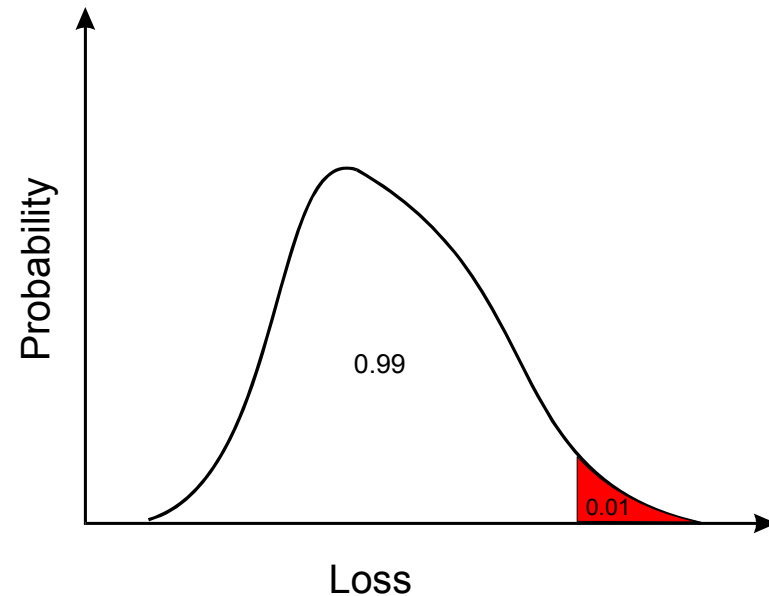
Decision-makers who care about risk  
(Probability Management framework)

# Risk-Orientation Leads to Interesting Research

- Ex 1: MCDB-R
- Ex 2: Risk in top-K queries

# Ex. 1: MCDB-R

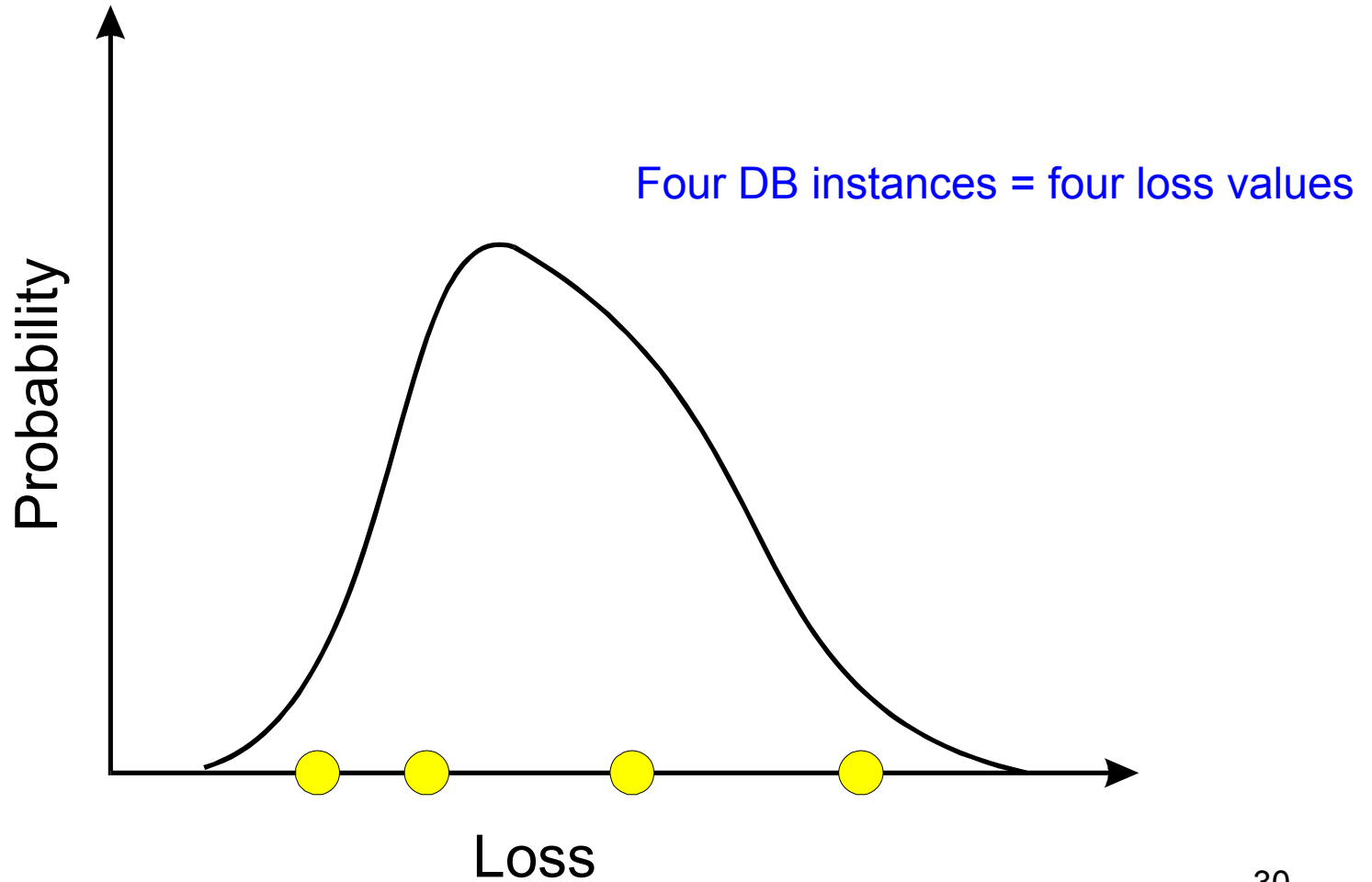
- Goals
  - Determine tail of query-result dist'n (e.g., 0.99-quantile =  $\text{VaR}_{0.01}$ )
  - Generate samples from tail\*
- Challenge for naïve MCDB
  - Huge # of replications needed



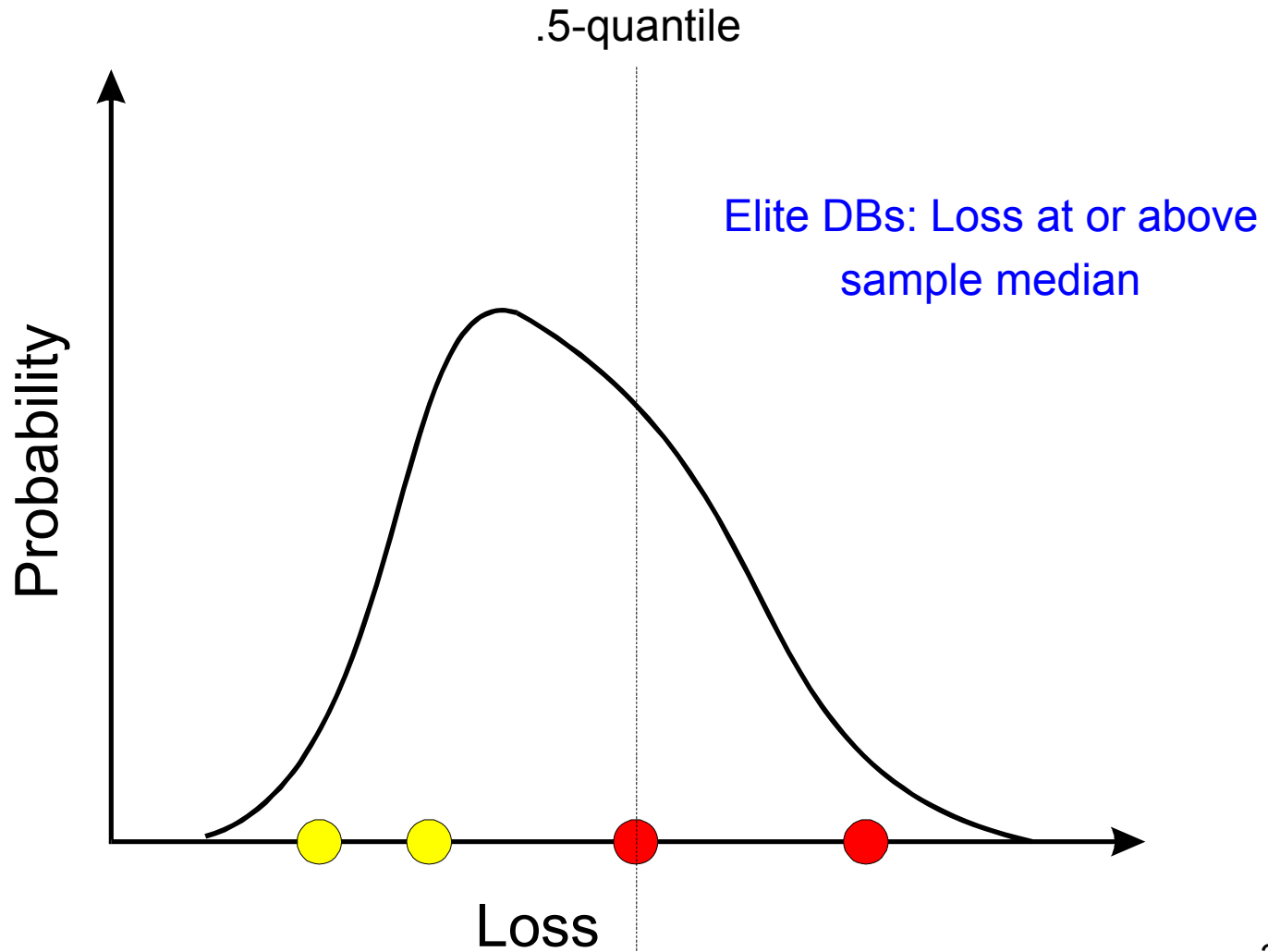
Normal(\$10M,\$1M) loss:  
On average,  $3.5 \times 10^6$  reps before even one \$15M loss is observed!

\*Degen, M., Lambrigger, D.D., Segers, J.: Risk Concentration and Diversification - Second-Order Properties.  
*Insurance: Mathematics and Economics* 46(3), 2010

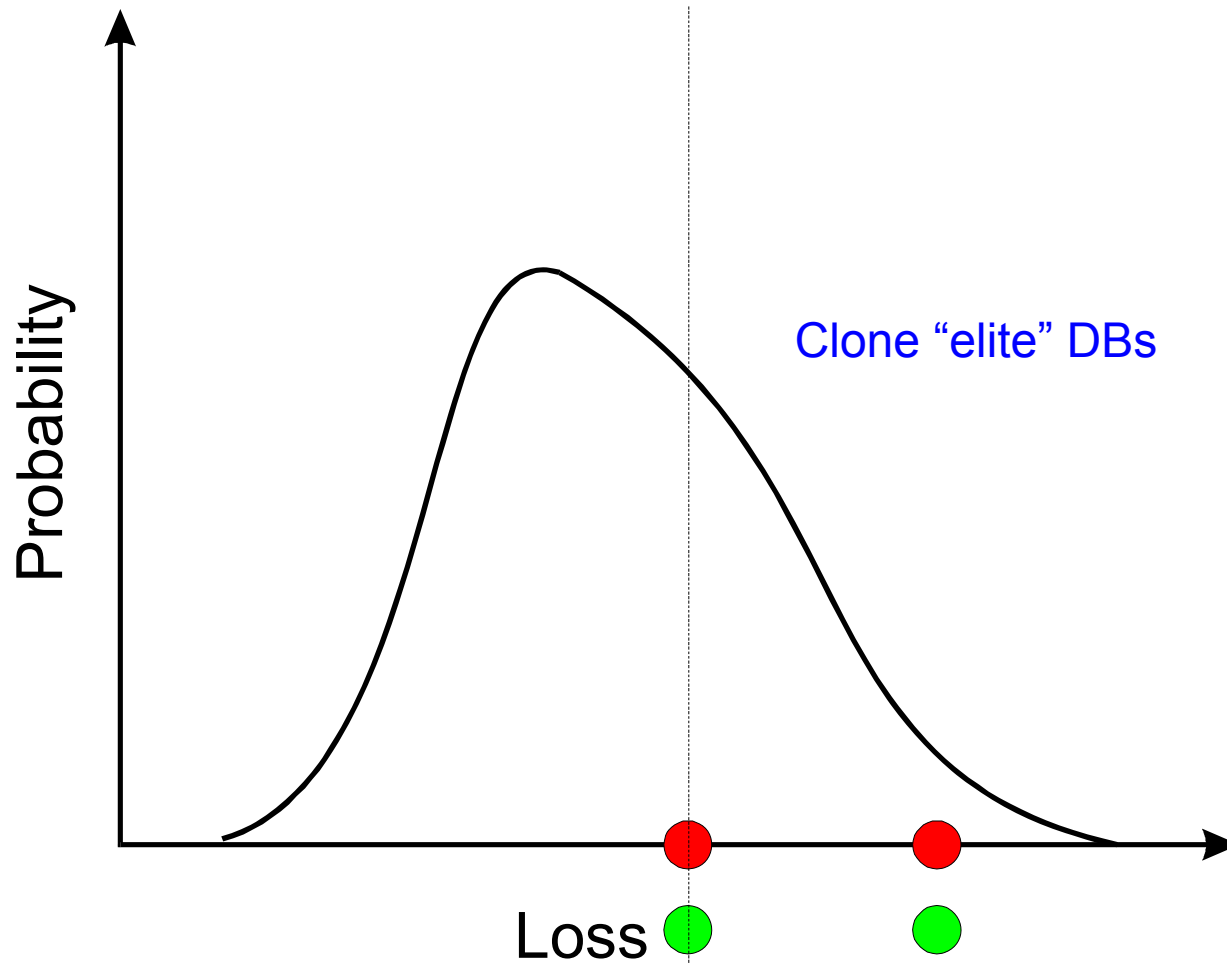
# Gibbs-Cloner Approach



# Gibbs-Cloner Approach

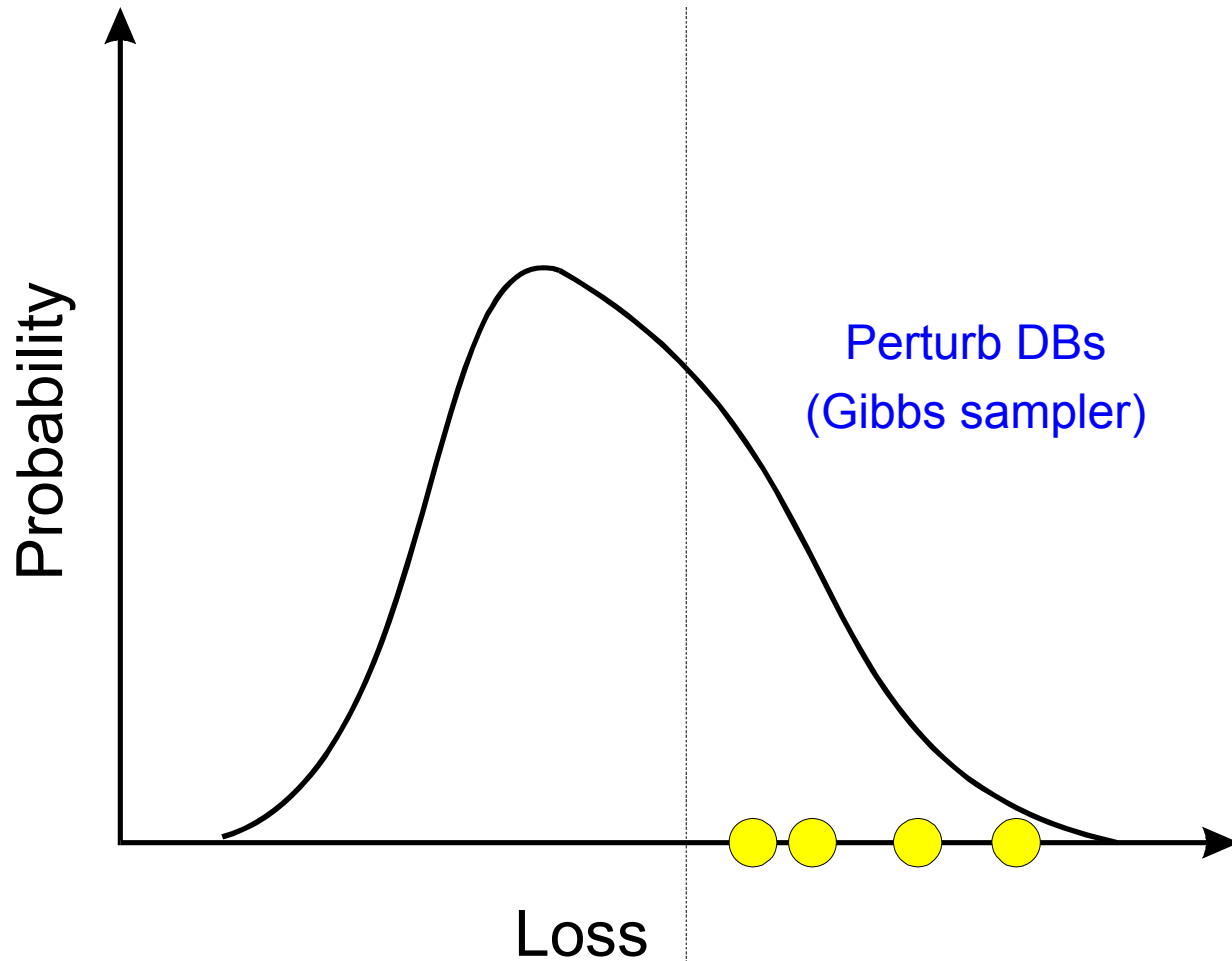


# Gibbs-Cloner Approach

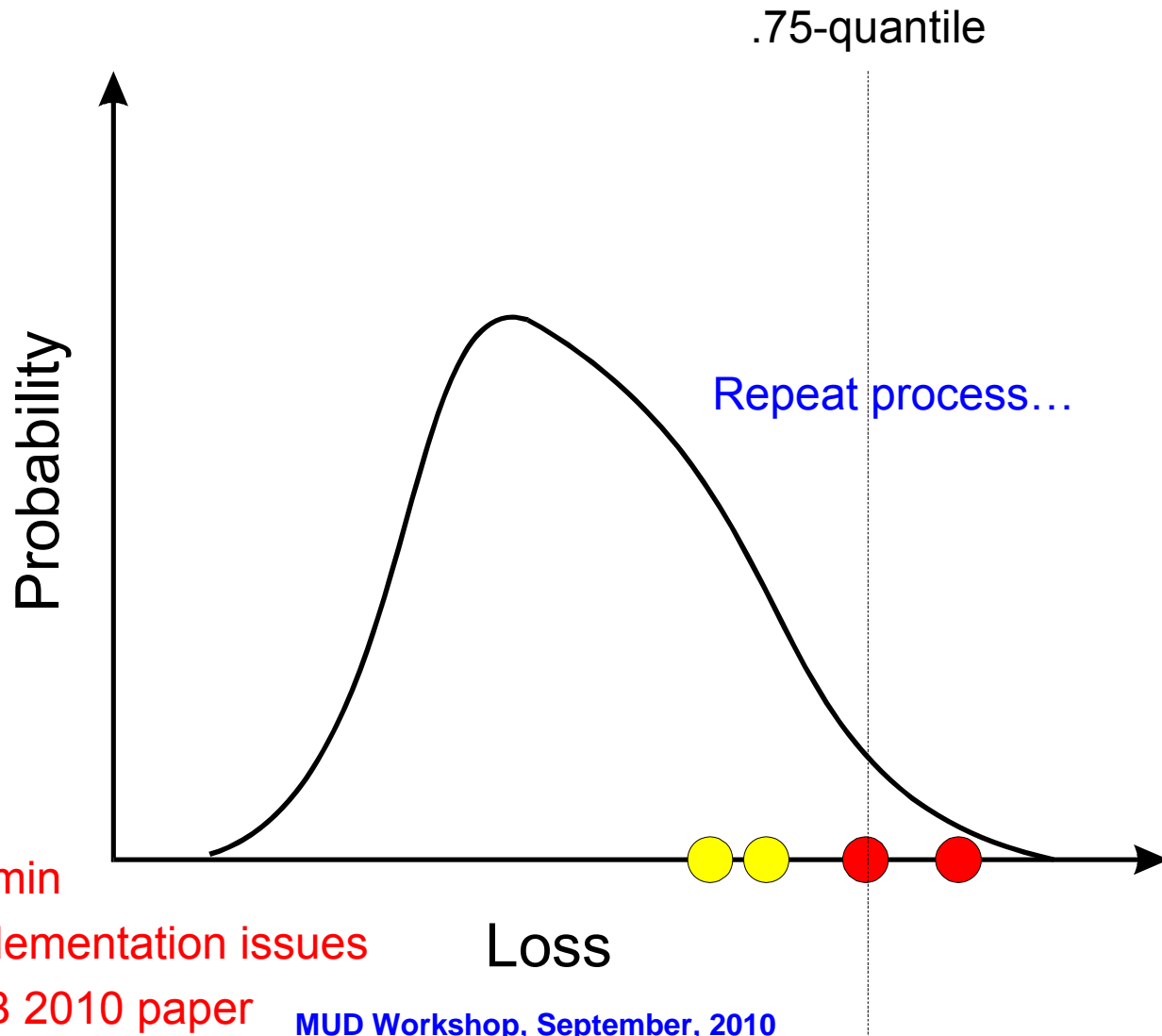




# Gibbs-Cloner Approach



# Gibbs-Cloner Approach



- 18 hrs  $\Rightarrow$  11 min

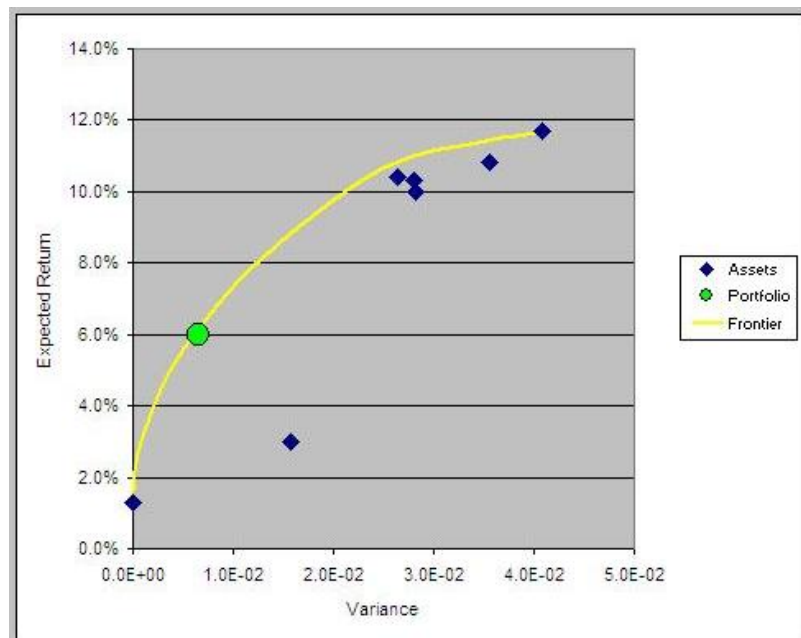
- Complex implementation issues

- Details: VLDB 2010 paper

MUD Workshop, September, 2010

# Ex. 2: Portfolio Theory of IR

- Wang and Zhu [SIGIR 2009]
  - Uncertain relevance (score)
  - Balance mean/variance of “overall relevance” of document group =  $\sum_i (R_i \times w_i)$
  - Diversification of documents
  - Q: Other loss functions?



# Summary

- Easier to sell “stochastic predictive analytics over big data” than “data warehouse uncertainty” to real-world clients
- Risk management is a key driver in this setting but decision-makers are surprisingly clueless
- Probability-management ecosystem: a channel from ProbDBs to decision-makers?
- Risk-orientation leads to interesting research questions as well as potential impact

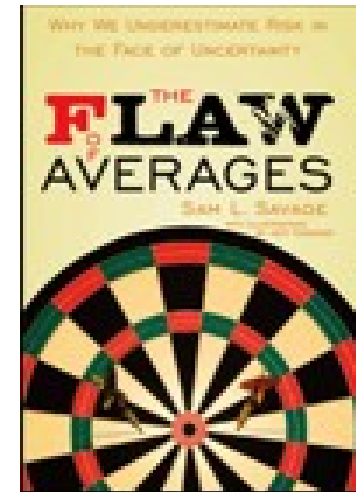


# Special Thanks

- Sam Savage
- Amol Deshpande
- Chris Jermaine and students
- Yannis Sismanis

# Further Details:

- RAQA: ICDE 2009
- MCDB: SIGMOD 2008
- MC<sup>3</sup>: SIGMOD 2009
- ProbIE: SIGMOD 2009
- MCDB-R: VLDB 2010



<http://probabilitymanagement.org>

[www.almaden.ibm.com/cs/people/peterh](http://www.almaden.ibm.com/cs/people/peterh)

[peterh@almaden.ibm.com](mailto:peterh@almaden.ibm.com)

## Thank You!