# Analytic DB Technology
# for the
# Data Enthusiast

**Pat Hanrahan**

**Stanford & Tableau**

**SIGMOD Keynote 2012**

# My 1ˢᵗ Job: Analyzing Data



University of Wisconsin

Experimental Particle Physics

# My 1ˢᵗ Job: Analyzing Data

# Data Analysis is Spreading

**Doctor:**

**Why are my patients returning to the hospital?**

**Call Center Operator:**

**Why does dispatching a tow truck cost so much in ND?**

**Doll Collector:**

**What caused the price of vintage Barbie dolls to increase?**

**Game Producer:**

**What causes players to buy virtual goods?**

# Why are databases so slow?

# Data Scientist



Netflix Prize Team



Data science?

engineering

math

nerds

nerds

nerds

comp sci

nerds

hacking

awesome nerds



THE QUANTS

How a New Breed of
Math Whizzes
Conquered Wall Street
and Nearly
Destroyed It

SCOTT PATTERSON
STAFF REPORTER, *THE WALL STREET JOURNAL*

# Analytical Thinking?

"HOLMES GAVE ME A SKETCH OF THE EVENTS."

# Sherlock Holmes

"It is a capital mistake to theorize before one has data. Insensibly one begins to twist facts to suit theories, instead of theories to suit facts."

"Data, Data, Data! I can't make bricks without clay."

# Definition: Analytical Thinking

**"A structured approach**

**to answering questions**

**and making decisions**

**based on facts and data"**

# My Process

Pose the question

Find or collect the appropriate data

Check and verify

Clean and normalize

Contextualize the data by joining with other data

Explore relationships & patterns in the raw data

Generalize and summarize

Confirm hypotheses and analyze errors

Share findings with others

Decide and act

**Question**

Forage for data

Decide and act

Check and clean

Test hypothesis, analyze errors, discover insight

Show relationships and patterns using visual representations

# "Data Analysis is like doing Experiments," J. Tukey

## Experiments

1. Theorize and hypothesize

2. Experiment

3. Revise theory

4. The craft occupies the experimenter allowing time to think

## Data Analysis

1. Theorize and hypothesize

2. Find trends and relationships

3. Find limitations of the model

4. Provide insight to improve the model

# "State of the Art"



## Spreadsheets

# "State of the Art"



## Title

Server_Hostname: itm64vm13.tivlab.raleigh.ibm.com

| AVG_Used_CPU_MHz | | Feb 14, 2011 | Feb 15, 2011 | Feb 16, 2011 | Feb 17, 2011 | Feb 18, 2011 | Feb 19, 2011 | Feb 20, 2011 | Feb 21 |
|---|---|---|---|---|---|---|---|---|---|
| vsvdash1 | itm64vm13.tivlab.raleigh.ibm.com | 182.73 | 169.17 | 166.52666667 | 163.02 | 161.25 | 154.7 | 152.34 | |
| | **vsvdash1** | **182.73** | **169.17** | **166.52666667** | **163.02** | **161.25** | **154.7** | **152.34** | |
| vsvdash2 | itm64vm13.tivlab.raleigh.ibm.com | 124.45 | 109.04333333 | 155.74333333 | 117.35 | 114.705 | 110.3375 | 109.87666667 | |
| | **vsvdash2** | **124.45** | **109.04333333** | **155.74333333** | **117.35** | **114.705** | **110.3375** | **109.87666667** | |
| vsvdash3 | itm64vm13.tivlab.raleigh.ibm.com | 122.98 | 111.82 | 226.14333333 | 120.425 | 112.77 | 107.215 | 110.56666667 | |
| | **vsvdash3** | **122.98** | **111.82** | **226.14333333** | **120.425** | **112.77** | **107.215** | **110.56666667** | |
| vsvtaddm | itm64vm13.tivlab.raleigh.ibm.com | 123.69 | 123.195 | 122.2 | 121.255 | 121.965 | 125.1475 | 123.55 | |
| | **vsvtaddm** | **123.69** | **123.195** | **122.2** | **121.255** | **121.965** | **125.1475** | **123.55** | |
| win2008vm1 | itm64vm13.tivlab.raleigh.ibm.com | 2,654.675 | 2,653.92666667 | 2,655 | 2,653.35 | 2,654.66 | 2,654.465 | 2,655.15666667 | 2,655.0 |
| | **win2008vm1** | **2,654.675** | **2,653.92666667** | **2,655** | **2,653.35** | **2,654.66** | **2,654.465** | **2,655.15666667** | **2,655.05** |
| **Summary** | | **641.705** | **669.87642857** | **703.90285714** | **635.08** | **633.07** | **630.373** | **666.49428571** | **671.54** |

## Crosstabs and Pivot Tables

# Idea: Visual Analysis

"Analytical Reasoning

Facilitated by

Interactive Visualization"

# Polaris / Tableau Demo

# C. Stolte's PhD Thesis

# 1. Best Visualization Depends on the Question/Task

## Summary of Financial Performance

| | | Central | | East | | South | | West | |
|---|---|---|---|---|---|---|---|---|---|
| | | Profit | Sales | Profit | Sales | Profit | Sales | Profit | Sales |
| Coffee | Amaretto | $5,104 | $14,012 | $1,010 | $2,994 | | | -$1,224 | $9,263 |
| | Columbian | $8,525 | $28,911 | $27,256 | $47,385 | $8,767 | $21,663 | $11,256 | $30,352 |
| | Decaf Irish Cream | $9,635 | $26,157 | $2,726 | $6,262 | $2,935 | $11,596 | -$1,307 | $18,233 |
| Espresso | Caffe Latte | | | | | $3,873 | $15,443 | $7,502 | $20,456 |
| | Caffe Mocha | $14,642 | $35,218 | -$6,232 | $16,646 | $5,202 | $14,166 | $4,066 | $18,874 |
| | Decaf Espresso | $8,859 | $24,483 | $2,411 | $7,720 | $5,930 | $15,381 | $12,302 | $30,578 |
| | Regular Espresso | | | $10,065 | $24,031 | | | | |
| Herbal Tea | Chamomile | $14,435 | $36,571 | $764 | $2,193 | $3,178 | $11,183 | $8,854 | $25,631 |
| | Lemon | $6,253 | $21,982 | $7,902 | $27,177 | $2,593 | $14,494 | $13,121 | $32,273 |
| | Mint | $4,069 | $9,335 | -$2,243 | $11,991 | | | $4,328 | $14,384 |
| Tea | Darjeeling | $10,769 | $30,284 | $6,500 | $14,094 | | | $11,784 | $28,773 |
| | Earl Grey | $10,334 | $32,883 | $3,404 | $6,507 | | | $10,426 | $27,382 |
| | Green Tea | $1,227 | $5,209 | $5,654 | $11,576 | | | -$7,112 | $16,065 |

# How much mint tea was sold in the west?

## Summary of Financial Performance

| | | Central | | East | | South | | West | |
|---|---|---|---|---|---|---|---|---|---|
| | | Profit | Sales | Profit | Sales | Profit | Sales | Profit | Sales |
| Coffee | Amaretto | $5,104 | $14,012 | $1,010 | $2,994 | | | -$1,224 | $9,263 |
| | Columbian | $8,525 | $28,911 | $27,256 | $47,385 | $8,767 | $21,663 | $11,256 | $30,352 |
| | Decaf Irish Cream | $9,635 | $26,157 | $2,726 | $6,262 | $2,935 | $11,596 | -$1,307 | $18,233 |
| Espresso | Caffe Latte | | | | | $3,873 | $15,443 | $7,502 | $20,456 |
| | Caffe Mocha | $14,642 | $35,218 | -$6,232 | $16,646 | $5,202 | $14,166 | $4,066 | $18,874 |
| | Decaf Espresso | $8,859 | $24,483 | $2,411 | $7,720 | $5,930 | $15,381 | $12,302 | $30,578 |
| | Regular Espresso | | | $10,065 | $24,031 | | | | |
| Herbal Tea | Chamomile | $14,435 | $36,571 | $764 | $2,193 | $3,178 | $11,183 | $8,854 | $25,631 |
| | Lemon | $6,253 | $21,982 | $7,902 | $27,177 | $2,593 | $14,494 | $13,121 | $32,273 |
| | Mint | $4,069 | $9,335 | -$2,243 | $11,991 | | | $4,328 | $14,384 |
| Tea | Darjeeling | $10,769 | $30,284 | $6,500 | $14,094 | | | $11,784 | $28,773 |
| | Earl Grey | $10,334 | $32,883 | $3,404 | $6,507 | | | $10,426 | $27,382 |
| | Green Tea | $1,227 | $5,209 | $5,654 | $11,576 | | | -$7,112 | $16,065 |

# What product in what region sold the most?

**What product in what region sold the most?**

# 2. Formulate Any Query

Q2. Find the department(s) that sells an item(s) supplied by the supplier Parker.

Here the user fills in both the SALES and the SUPPLY Tables as follows.

| SALES | DEPT | ITEM |
|-------|------|------|
|       | P.TOY | PEN |

| SUPPLY | ITEM | SUPPLIER |
|--------|------|----------|
|        | PEN  | PARKER   |

ANS:

| DEPT |
|------|
| HOUSEHOLD |
| TOY |
| STATIONARY |
| HARDWARE |

**Query-By-Example [Zloof, 1975]**

2008 Presidential Election: Where the Donors Are

2008 Presidential Election: Where the Donors Are

Obama, Barack | McCain, John S

**SELECT AS CIRCLE**

**Candidate * Longitude ON X**

**Latitude ON Y**

**Zipcode IN PANES**

**Party ON COLOR**

**Sum(Amount) ON SIZE**

**FROM ContributionsDatabase**

SELECT AS SHAPE

    Market * Sales ON COLS

    Quarter * Profit ON ROWS

    State * Product IN PANES

    ProductType ON COLOR

    Year ON SHAPE

FROM RetailDatabase

# Four Main Ideas

1. Support cycle of analysis

2. Answer a question by composing a picture

3. Best visualization depends on question/task

4. Must be able to generate any query


+ Easy to use

# Analysis at the
# Speed of Thought

# Transactional Databases are Slow!!

TPC-H, 1 GB, Query 1*:

| C Program | 0.2 s |
|---|---|
| mysql | 26.2 s |
| DBMS "X" | 28.4 s |

*Boncz et al., CIDR 2005

# In-Memory Column Stores: 100x

Columns are efficient (MonetDB/X100, C-store, ...)

- Only access needed columns

- Well-matched to processor+memory architecture

- Columns compress better than records

- Optimized for read/append

- Vector semantics instead of set semantics

In-Memory reduces latency enabling interaction

- Memory is cheap, memory hierarchy is expending

- Median business database fits in memory

# kx.com



EVENT STREAMS

Market Data Feeds

Other Events (News Feeds)

MEMORY

Feed Handlers / Adapters

Ticker Plant

**2 B trades per day**

External Interfaces

Java

.NET

C

kdb+ In Memory Database

kdb+ Event Processing

NBBO | Filter Logic | VWAP

Depth of Book | Business Logic

Publish & Subscribe

**STAC-M3**

**~10-20 msec latency**

Application Clients

Other Event Stream | Order Mgmt. System | Research Applications

Intraday Views | Historical Analysis | Back Test

DISK

kdb+ Log File/ Journals

kdb+ Historical Database

Legacy Databases

Reference

Compliance

Corporate

©2008 Kx Systems

# High Frequency Trading

# Fully Utilize *All* Resources

Intel Ivy Bridge Processor  (Core i7 3770K)

    22 nm, 1.4B 3D transistors

    4 3.5 Ghz cores

    256-bit AVX vector instructions (16 FADD/FMUL)

Resource limits

    Bandwidth limited: 2 DDR3 2133 = 34 GB/s

    Compute limited: 4 * 3.5 Ghz * 16 = 224 GFLOPS


Theoretical: 1B values can be summed in 125/5 msecs

**2018 Laptop ~ CPU+GPU**

**10 Teraflops**

# Supporting Data Enthusiasts

"Although we often hear that data speak for themselves, their voices can be soft and sly.

We need statistics to help them tell their story"

Beginning Statistics with Data Analysis
Mosteller, Fienberg, Rourke

# Data Integration

**Provides context for analysis**

**Semantic integration => people**

**Promising tools**

- **Potters wheel**
- **Google fusion tables**
- **Data wrangler**
- **Data blending**



Dynamic Workload Driven Data Integration in Tableau

K. Morton, R. Bunker, J. Mackinlay, R. Morton, C. Stolte

# Wrap Up

# Summary

**Large number of data enthusiasts**

- Business users, with the questions, on a mission

- Excellent analytical thinkers

- Not DBAs, not programmers, not statisticians

**You can help them**

- Current tools support only basic visual analysis

- … not the entire process of analysis in the large

# Thank You