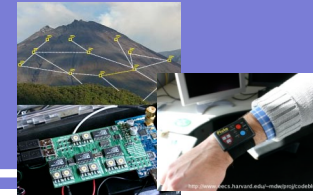# MauveDB: Statistical Modeling inside Database Systems
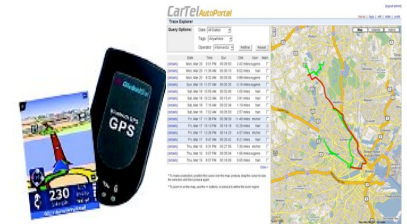
Amol Deshpande, University of Maryland

# Motivation


*Wireless sensor networks*

- Unprecedented, and rapidly increasing, instrumentation of our every-day world

- Huge data volumes generated *continuously* that must be processed in *real-time*


*Distributed measurement networks (e.g. GPS)*

- Typically *imprecise, unreliable and incomplete* data

  - Inherent measurement noises (e.g. GPS)

  - Low success rates (e.g. RFID)

  - Communication link or sensor node failures (e.g. wireless sensor networks)

  - Spatial and temporal biases


*RFID*

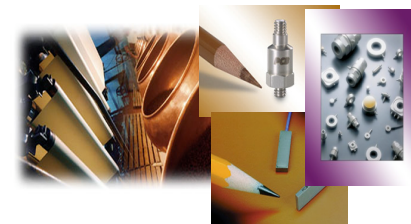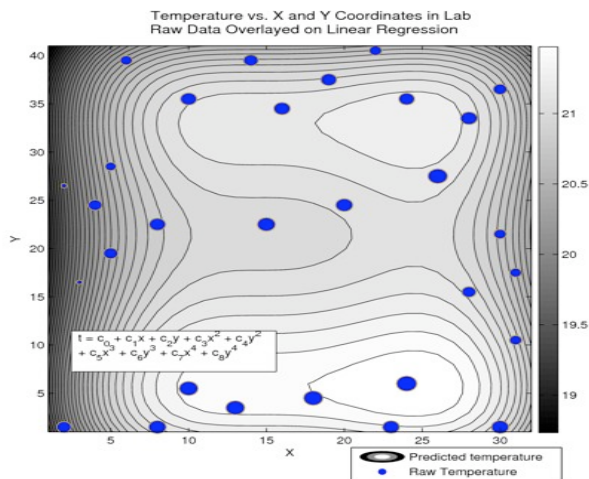- Raw sensed data is not what users want to see/query


*Industrial Monitoring*

# Data Processing Step 1

- Process data using a statistical/probabilistic model
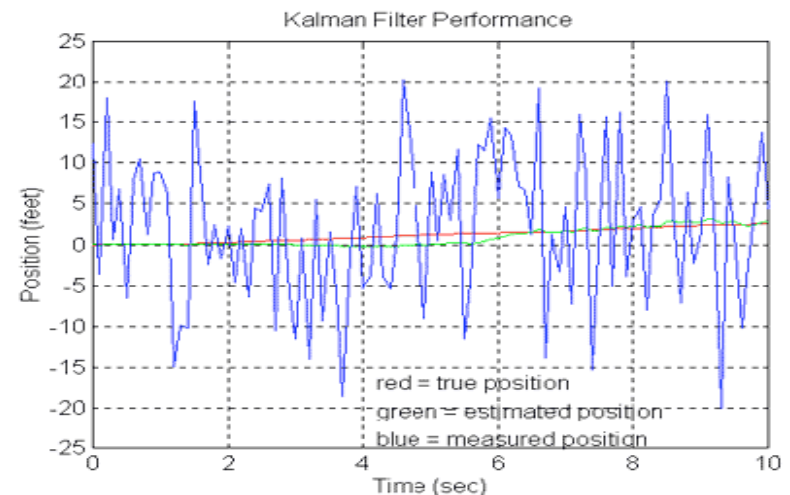    - Regression and interpolation models
        - To eliminate spatial or temporal biases, handle missing data, prediction
    - Filtering techniques *(e.g. Kalman Filters)*, Bayesian Networks
        - To eliminate measurement noise, to infer hidden variables etc

### *Temperature monitoring*



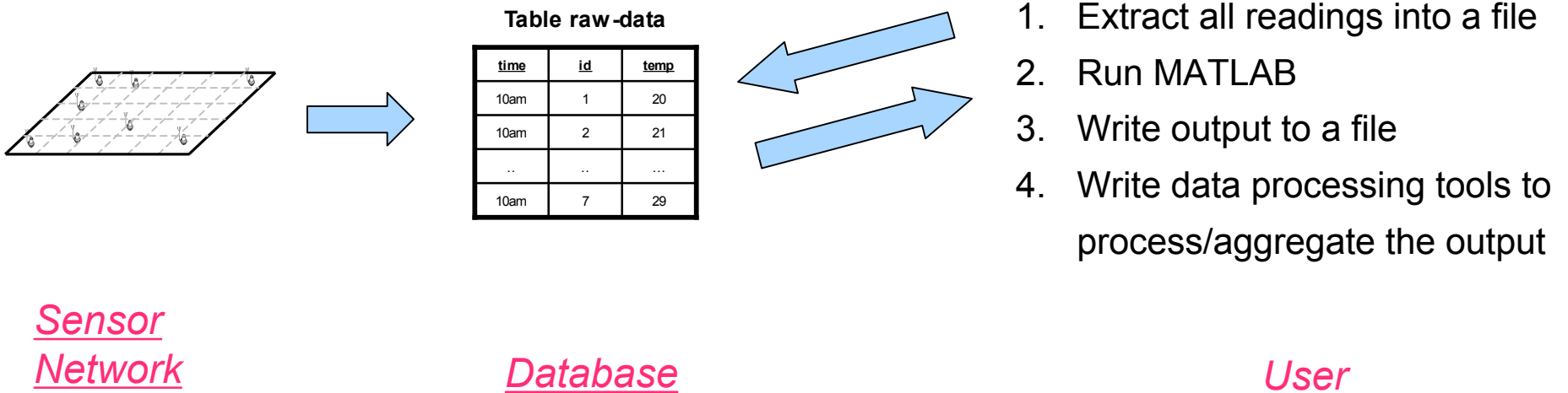*Regression/interpolation models*

### *GPS Data*



*Kalman Filters et*

# Statistical Modeling of Sensor Data

- No support in database systems --> Database ends up being used as a backing store
  - With much replication of functionality
  - Very inefficient, not declarative…
- How can we push statistical modeling inside a database system ?

**Table raw-data**

| time | id | temp |
|------|-----|------|
| 10am | 1 | 20 |
| 10am | 2 | 21 |
| .. | .. | … |
| 10am | 7 | 29 |

1. Extract all readings into a file
2. Run MATLAB
3. Write output to a file
4. Write data processing tools to process/aggregate the output
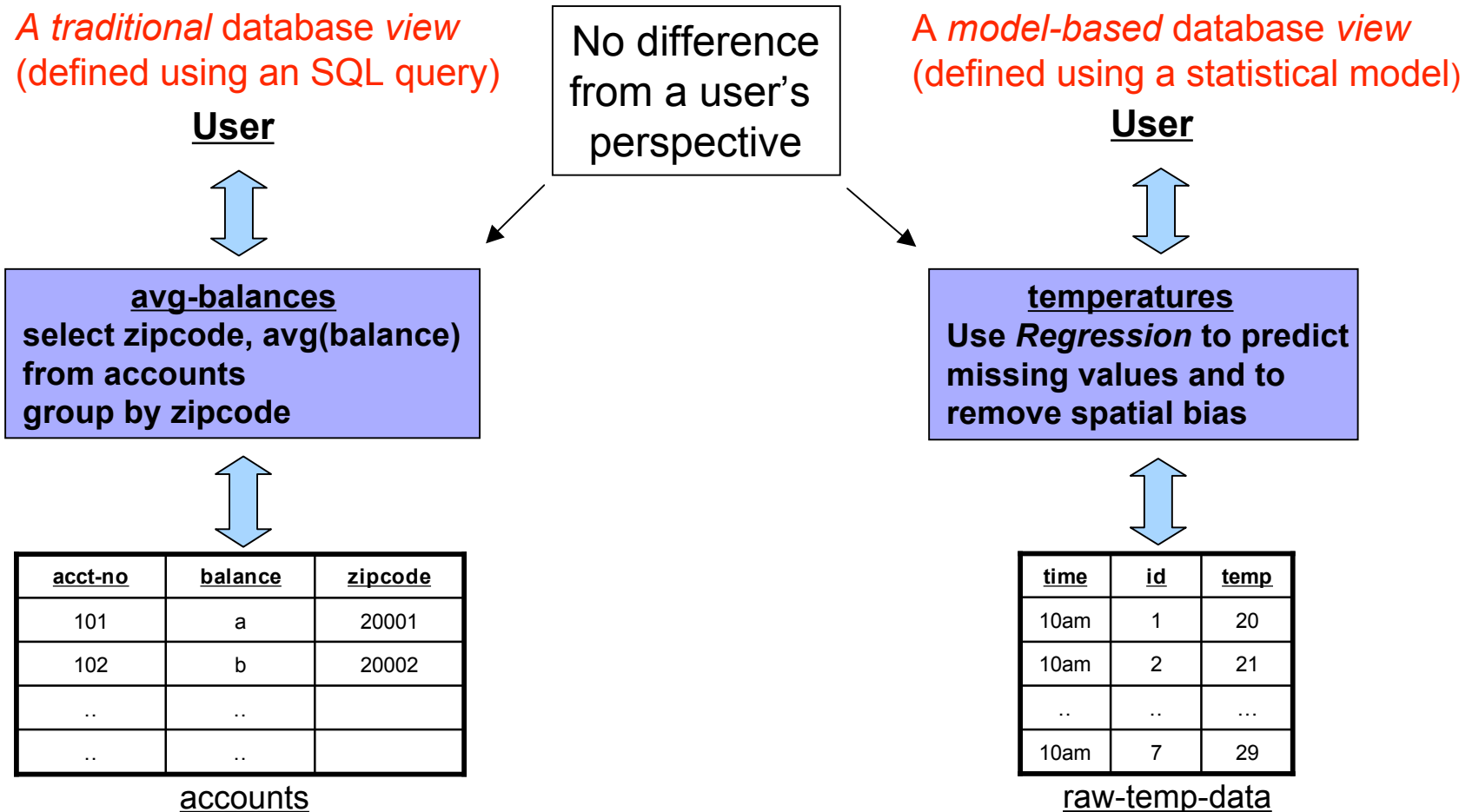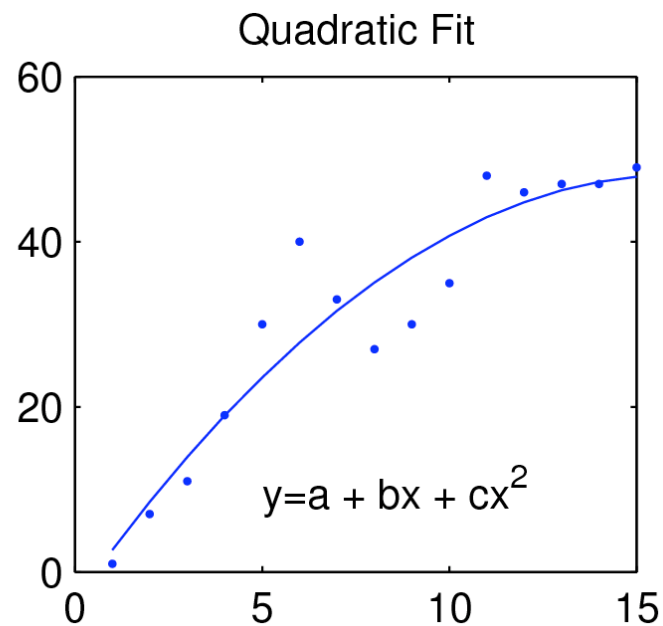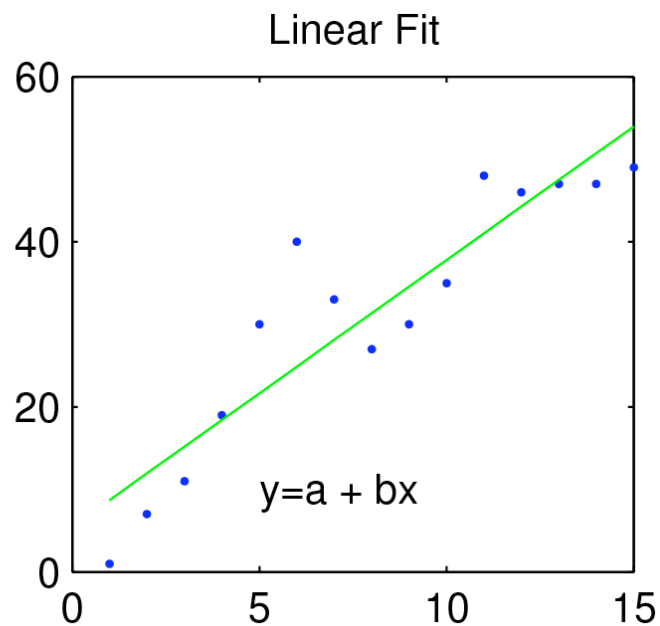
*Sensor Network*

*Database*

*User*

# Abstraction: Model-based Views

- An abstraction analogous to *traditional database views*
- Provides *independence from the messy measurement details*

*A traditional* database *view*
(defined using an SQL query)

**User**

No difference
from a user's
perspective

*A model-based* database *view*
(defined using a statistical model)

**User**

**avg-balances**
**select zipcode, avg(balance)**
**from accounts**
**group by zipcode**

**temperatures**
**Use *Regression* to predict**
**missing values and to**
**remove spatial bias**

| acct-no | balance | zipcode |
|---------|---------|---------|
| 101 | a | 20001 |
| 102 | b | 20002 |
| .. | .. | |
| .. | .. | |

accounts

| time | id | temp |
|------|-----|------|
| 10am | 1 | 20 |
| 10am | 2 | 21 |
| .. | .. | ... |
| 10am | 7 | 29 |

raw-temp-data

# Example: Regression-based Views

*Regression:*

*Model a dependent variable as a function of independent variables*

### Linear Fit

$$y = a + bx$$

### Quadratic Fit

$$y = a + bx + cx^2$$
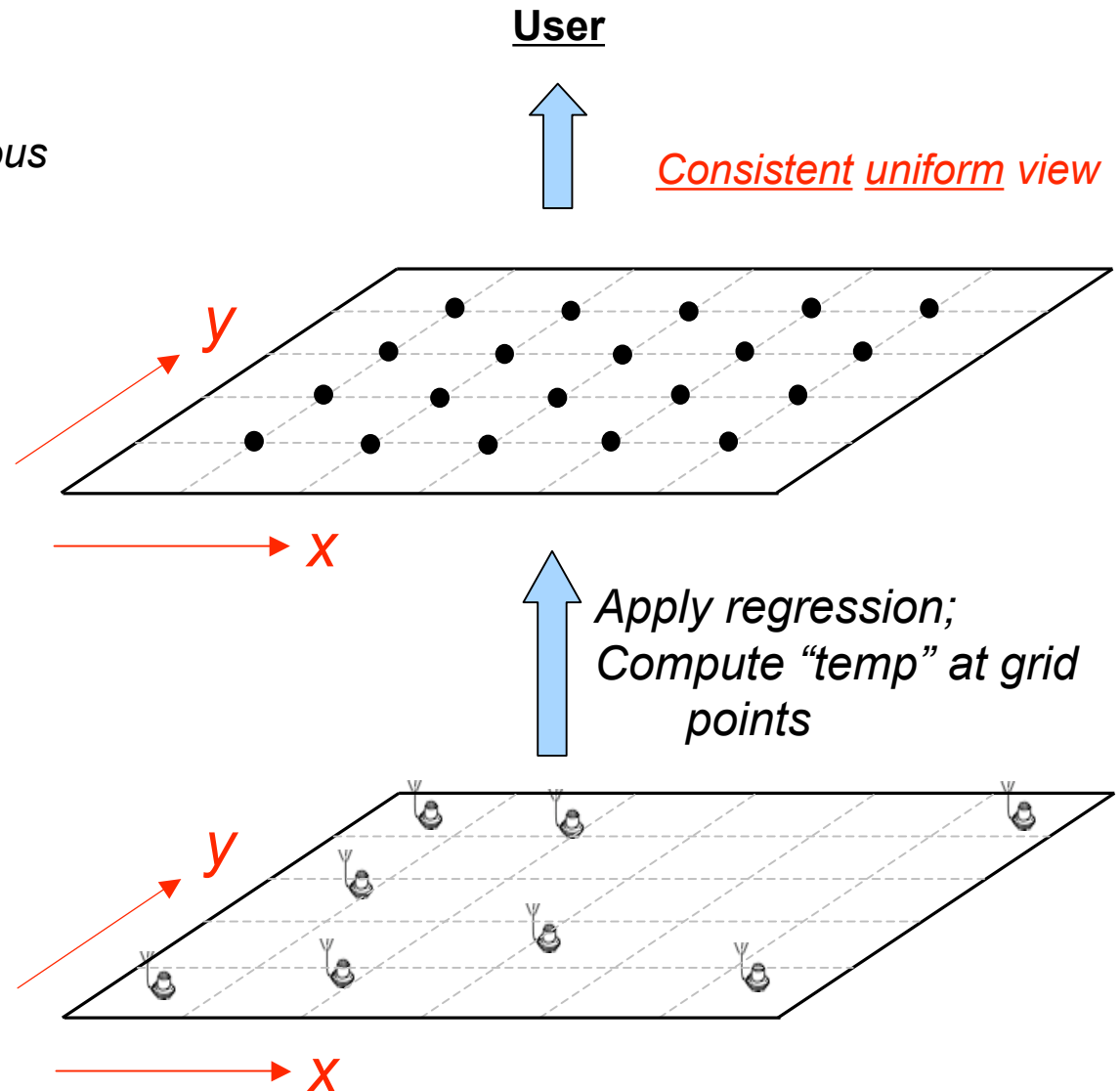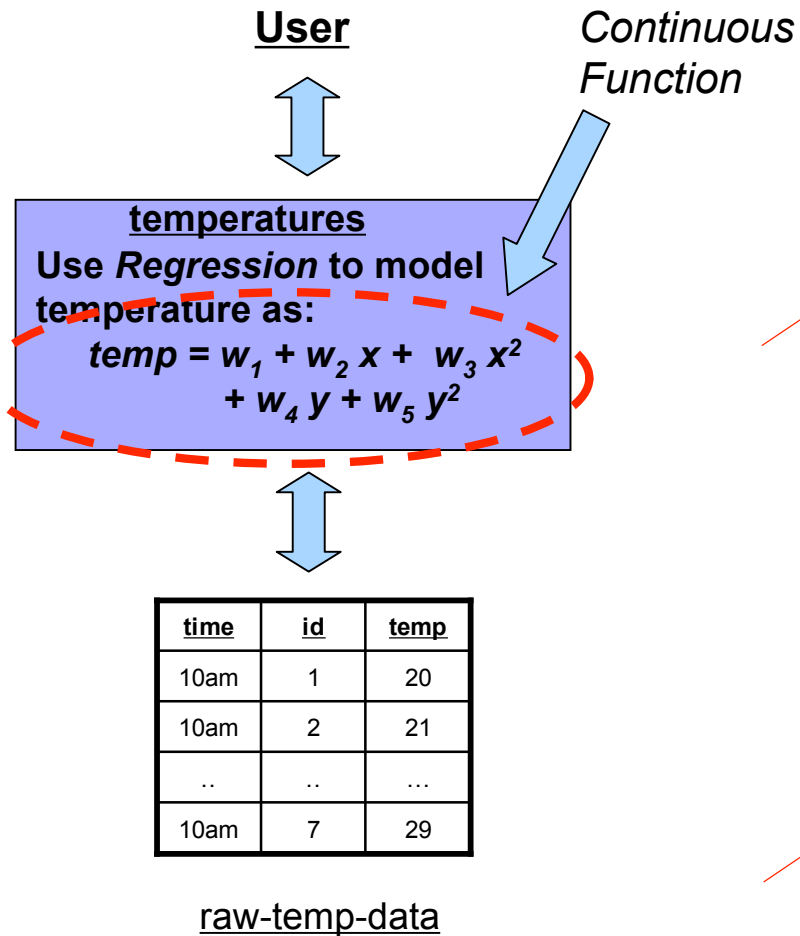
# Example: Regression-based Views

Model *temperature* as a function of *(x, y)*

*E.g.*
$$temp = w_1 + w_2 * x + w_3 * x^2 + w_4 * y + w_5 * y^2$$

# Grid Abstraction

**A *Regression-based View***

**User**

↕

Continuous Function

**temperatures**
**Use *Regression* to model temperature as:**
$$temp = w_1 + w_2\,x + w_3\,x^2 + w_4\,y + w_5\,y^2$$

↕

| time | id | temp |
|------|-----|------|
| 10am | 1 | 20 |
| 10am | 2 | 21 |
| .. | .. | … |
| 10am | 7 | 29 |

raw-temp-data

**User**

↑

*Consistent uniform view*

y

x

↑ *Apply regression; Compute "temp" at grid points*

y

x

# Creating a Regression-based View

CREATE VIEW

   RegView(time [0::1], x [0:100:10], y[0:100:10], temp)

AS

   FIT temp USING time, x, y

   BASES 1, x, $x^2$, y, $y^2$

   FOR EACH time T

   TRAINING DATA

      SELECT temp, time, x, y

      FROM raw-temp-data

      WHERE raw-temp-data.time = T

Fit as:
$$temp = w_1 + w_2 * x + w_3 * x^2 + w_4 * y + w_5 * y^2$$

# Query Processing

- Analogous to querying database tables
  - *select * from reg-view*
    - Lists out temperatures at all grid-points
  - *select * from reg-view where x = 15 and y = 20*
    - Lists temperature at (15, 20) at all times
  - …
- How are queries evaluated ?
  - Different options
    - Do the statistical modeling it as soon as new data arrives
    - *or* when the queries are asked (on demand)
    - *or* …
  - Optimization opportunities that the database system can exploit
    - Without bothering the user

# MauveDB: Status

- Written in the Apache Derby Java open source database system

- Support for *Regression-* and *Interpolation-based views*
  - Declarative constructs for defining and querying views
  - Several update and materialization strategies
  - SIGMOD 2006 (w/ Sam Madden)

- Currently building support for views based on *dynamic Bayesian networks*
  - *Kalman Filters, HMMs* etc

# Ongoing and Future Work

- Adding support for views based on *dynamic Bayesian networks (e.g. Kalman Filters)*
  - A very general class of models with wide applicability
  - Generate *probabilistic* data
- Developing APIs for adding arbitrary models
  - Minimize the work of the model developer
- *Probabilistic databases*
  - Uncertain data with complex correlation patterns
- Query processing, query optimization
- View maintenance in presence of high-rate measurement streams