#### MauveDB: Supporting Model-based User Views in Database Systems

Amol Deshpande, University of Maryland Samuel Madden, MIT

## Motivation

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



Wireless sensor networks



Distributed measurement networks (e.g. GPS)



RFID



Industrial Monitoring

## **Motivation**

- Unprecedented, and rapidly increasing, instrumentation of our every-day world
  - Overwhelmingly large raw data volumes generated <u>continuously</u>
  - Data must be processed in *real-time*
  - The applications have strong *acquisitional* aspects
    - Data may have to be actively acquired from the environment
  - Typically <u>imprecise</u>, <u>unreliable</u> and <u>incomplete</u> data
    - Inherent measurement noises (e.g. GPS) and low success rates (e.g. RFID)
    - Communication link or sensor node failures (e.g. wireless sensor networks)
    - Spatial and temporal biases because of measurement constraints
- Traditional data management tools are ill-equipped to handle these challenges

#### **Example: Wireless Sensor Networks**



A wireless sensor network deployed to monitor temperature

#### **Example: Wireless Sensor Networks**



A wireless sensor network deployed to monitor temperature

# **Typical Solution**

- Process data using a statistical/probabilistic model before operating on it
  - 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



#### Issues

- Can't exploit commonalities, reuse/share computation
- No easy way to keep the model outputs up-to-date
- Lack of declarative languages for querying the processed data
- Large amount of duplication of effort
- Non-trivial
  - Expert knowledge & MATLAB familiarity required !
- Prevents real-time analysis of the data in most cases
- Why are databases not doing any of this ?
  - We are very good at most of these things

#### **Solution: Model-based User Views**

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



#### MauveDB System

- Supports the abstraction of Model-based User Views
- Provides declarative language constructs for creating such views
- Supports SQL queries over model-based views
- Keeps the models up-to-date as new data is inserted into the database

#### MauveDB System

- Supports the abstraction of  $\underline{M}$  odel-b<u>a</u>sed  $\underline{U}$ ser  $\underline{V}$  iews
- Provides declarative language constructs for creating such views
- Supports SQL queries over model-based views
- Keeps the models up-to-date as new data is inserted into the database



## Outline

- Motivation
- Model-based views
  - Details, view creation syntax, querying
- Query execution strategies
- MauveDB implementation details
- Experimental evaluation

### Linear Regression

 Models a <u>dependent variable</u> as a function of a set of independent variables



## **Grid Abstraction**





## **View Creation Syntax**

Somewhat model-specific, but many commonalities

#### A Interpolation-based View



IntView(t [0::1], sensorid [::1], y[0:100:10], temp)

AS



FOR EACH sensorid M

TRAINING DATA

SELECT temp, time, sensorid

FROM raw-temp-readings

WHERE raw-temp-readings.sensorid = M

## Outline

#### Motivation

- Model-based views
  - Details, view creation syntax, querying
- Query execution strategies
- MauveDB implementation details
- Experimental evaluation

## **Querying a Model-based View**

- Analogous to traditional views
- So:
  - 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
  - . . .

### **Query Processing**

- Two operators per view type that support get\_next() API
  - ScanView
    - Returns the contents of the view one-by-one
  - IndexView (condition)
    - Returns tuples that match a condition
      - e.g. return *temperature* where (x, y) = (10, 20)



## **View Maintenance Strategies**

- Option 1: Compute the view as needed from base data
  - For regression view, scan the tuples and compute the weights
- Option 2: Keep the view materialized
  - Sometimes too large to be practical
    - E.g. if the grid is very fine
  - May need to be recomputed with every new tuple insertion
    - E.g. a regression view that fits a single function to the entire data
- Option 3: Lazy materialization/caching
  - Materialize query results as computed
- Generic options shared between all view types

## **View Maintenance Strategies**

- Option 4: Maintain an efficient *intermediate representation*
- Typically model-specific
- Regression-based Views
  - Say temp =  $f(x, y) = w_1 h_1(x, y) + \dots + w_k h_k(x, y)$
  - Maintain the *weights* for *f*(*x*, *y*) and a *sufficient statistic* 
    - Two matrices ( $O(k^2)$  space) that can be incrementally updated
  - ScanView: Execute *f*(*x*, *y*) on all grid points
  - IndexView: Execute f(x, y) on the specified point
  - InsertTuple: Recompute the coefficients
    - Can be done very efficiently using the sufficient statistic
- Interpolation-based Views
  - Build and maintain a tree over the tuples in the TRAINING DATA

## Outline

#### Motivation

- Model-based views
  - Details, view creation syntax, querying
- Query execution strategies
- MauveDB implementation details
- Experimental evaluation

#### **MauveDB: Implementation Details**

- Written in the Apache Derby Java open source database system
- Support for *Regression-* and *Interpolation-based views*
- Minimal changes to the main codebase
- Much of the additional code (approx 3500 lines) fairly generic in nature
  - A view manager (for bookkeeping)
  - Query processing operators
  - View maintenance strategies
- Model-specific code
  - Intermediate representation
  - Part of the view creation syntax

#### MauveDB: Experimental Evaluation

- Intel Lab Dataset
  - 54-node sensor network monitoring temperature, humidity etc
  - Approx 400,000 readings
  - Attributes used
    - Independent time, sensorid, x-coordinate, y-coordinate
    - Dependent temperature



## **Spatial Regression**

Contour plot over the data obtained using:

select \*

from reg-view

where time = 2100

Temperature vs. X and Y Coordinates in Lab Raw Data Overlayed on Linear Regression



## Interpolation



#### **Comparing View Maintenance Options**

- 50000 tuples initially
- Mixed workload:
  - insert 1000 records
  - issue 50 point queries
  - issue 10 average queries
- Brief summary:
  - Intermediate representation typically the best
  - Among others, dependent on the view properties, and query workload



#### Regression, per time



Interpolation, per sensor

## **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
- Query processing, query optimization, and view maintenance issues
- Much research still needs to be done

#### Conclusions

- Proposed the abstraction of model-based views
  - Poweful abstraction that enables declarative querying over noisy, imprecise data
- Exploit commonalities to define, to create, and to process queries over such views
- MauveDB prototype implementation
  - Using the Apache Derby open source DBMS
  - Supports Regression- and Interpolation-based views
  - Supports many different view maintenance strategies

# Thank you !!

• Questions ?