# Lecture 28

# Knowledge Discovery in Databases (KDD)
# and Data Mining

### Tuesday, December 04, 2001

## William H. Hsu

## Department of Computing and Information Sciences, KSU

http://www.cis.ksu.edu/~bhsu

Readings:

Handout, "Data Mining with MLC++", Kohavi *et al*
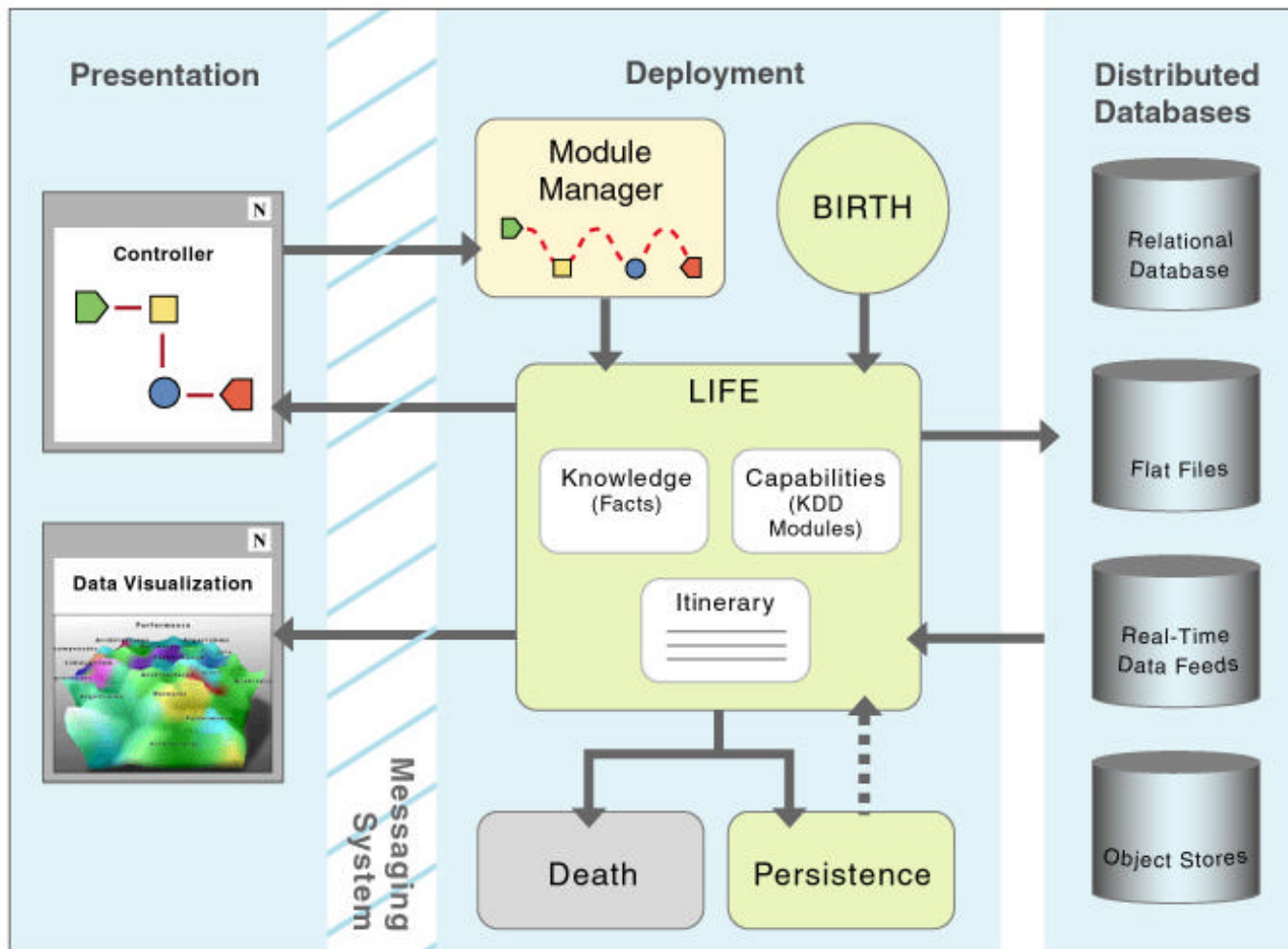
**KSU**

# Lecture Outline

- **Readings: "Data Mining with *MLC++*", Kohavi *et al***

- **Final Exam**

  - **Format**

    - **Open book**

    - **110 minutes**

    - **10 questions (see format online)**

  - **Sample questions online**

- **Knowledge Discovery in Databases (KDD) and Data Mining**

  - **Problem framework (stages)**

  - **Design and implementation issues**

- **Role of Machine Learning and Inference in Data Mining**

  - **Unsupervised learning**

  - **Supervised learning**

  - **Decision support (information retrieval, prediction, policy optimization)**

- **Next Lecture: Final Review Session**
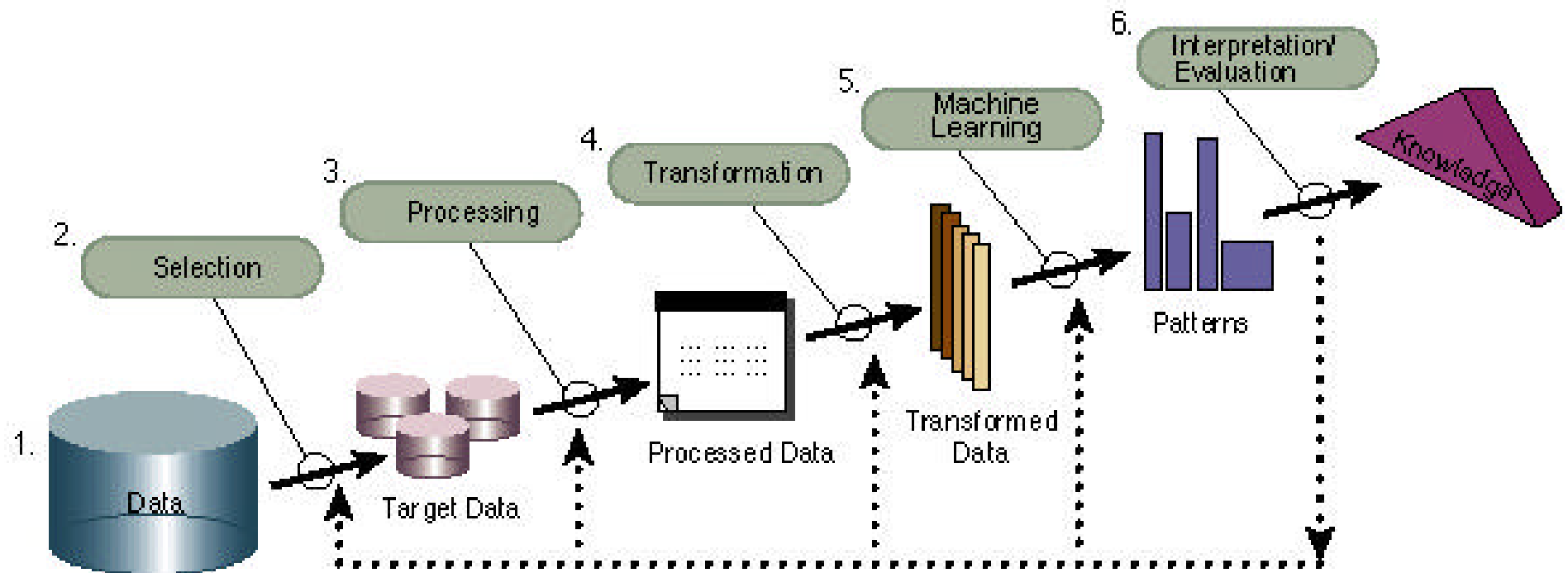
# What Is Data Mining?

- **Two Definitions (FAQ List)**
  - The process of automatically extracting valid, useful, previously unknown, and ultimately comprehensible information from large databases and using it to make crucial business decisions
  - *"Torturing the data until they confess"*

- **Data Mining: An Application of Machine Learning**
  - Guides and integrates learning (model-building) processes
    - Learning methodologies: supervised, unsupervised, reinforcement
    - Includes preprocessing (data cleansing) tasks
    - Extends to pattern recognition (inference or *automated reasoning*) tasks
  - Geared toward such applications as:
    - Anomaly detection (fraud, inappropriate practices, intrusions)
    - Crisis monitoring (drought, fire, resource demand)
    - Decision support

- **What Data Mining Is *Not***
  - Data Base Management Systems: *related but not identical field*
  - "Discovering objectives": still need to *understand performance element*

**KSU**

# KDD and Software Engineering



**Rapid KDD Development Environment**

# Stages of Data Mining



An Overview of the Steps That Compose the KDD Process

# Databases and Data Mining

- **Database Engineering ? Data Mining!**
  - <u>Database design and engineering</u>
    - <u>D</u>ata <u>B</u>ase <u>M</u>anagement <u>S</u>ystem (DBMS): computational system that supports efficient *organization*, *retrieval*, and *processing* of data
    - <u>Data warehouse</u>: repository of integrated information for queries, analysis
  - <u>Data mining</u>
    - Often an *application* of DBMS and data warehousing systems
    - Includes inductive model building (learning), pattern recognition, inference
- **Selection**
  - Guides and integrates learning (model-building) processes
  - Learning methodologies: supervised, unsupervised, reinforcement
  - Includes preprocessing (data cleansing), pattern recognition and inference
- <u>On</u>line <u>A</u>nalytical <u>P</u>rocessing (<u>OLAP</u>)
  - Efficient collection, storage, manipulation, reproduction of multidimensional data
  - Objective: analysis (e.g., for decision support)
  - See: http://perso.wanadoo.fr/bernard.lupin/english/glossary.html

# Data Integrity and Data Modeling: Ontologies



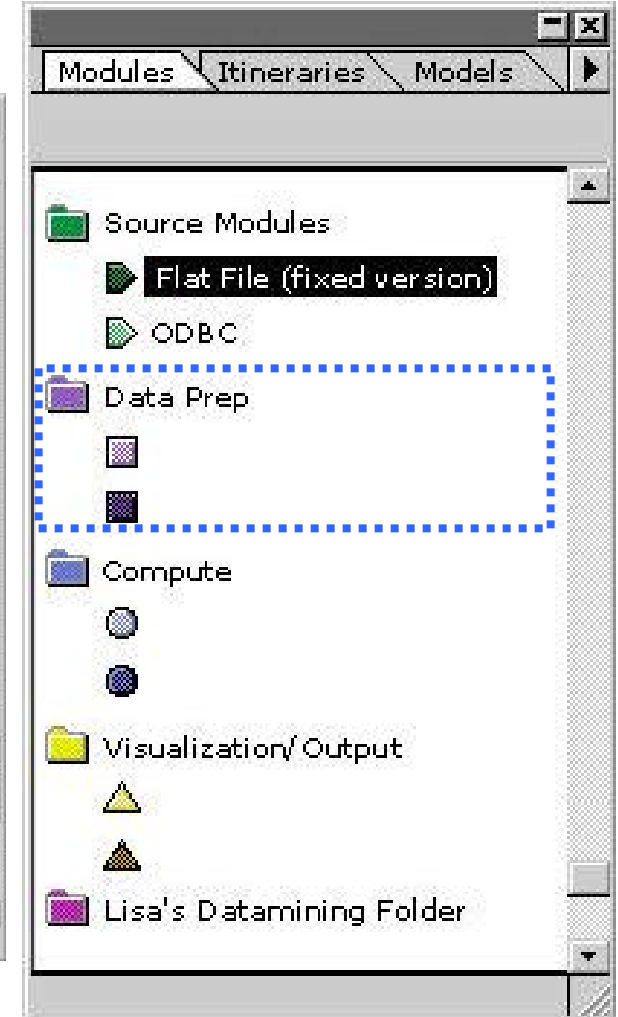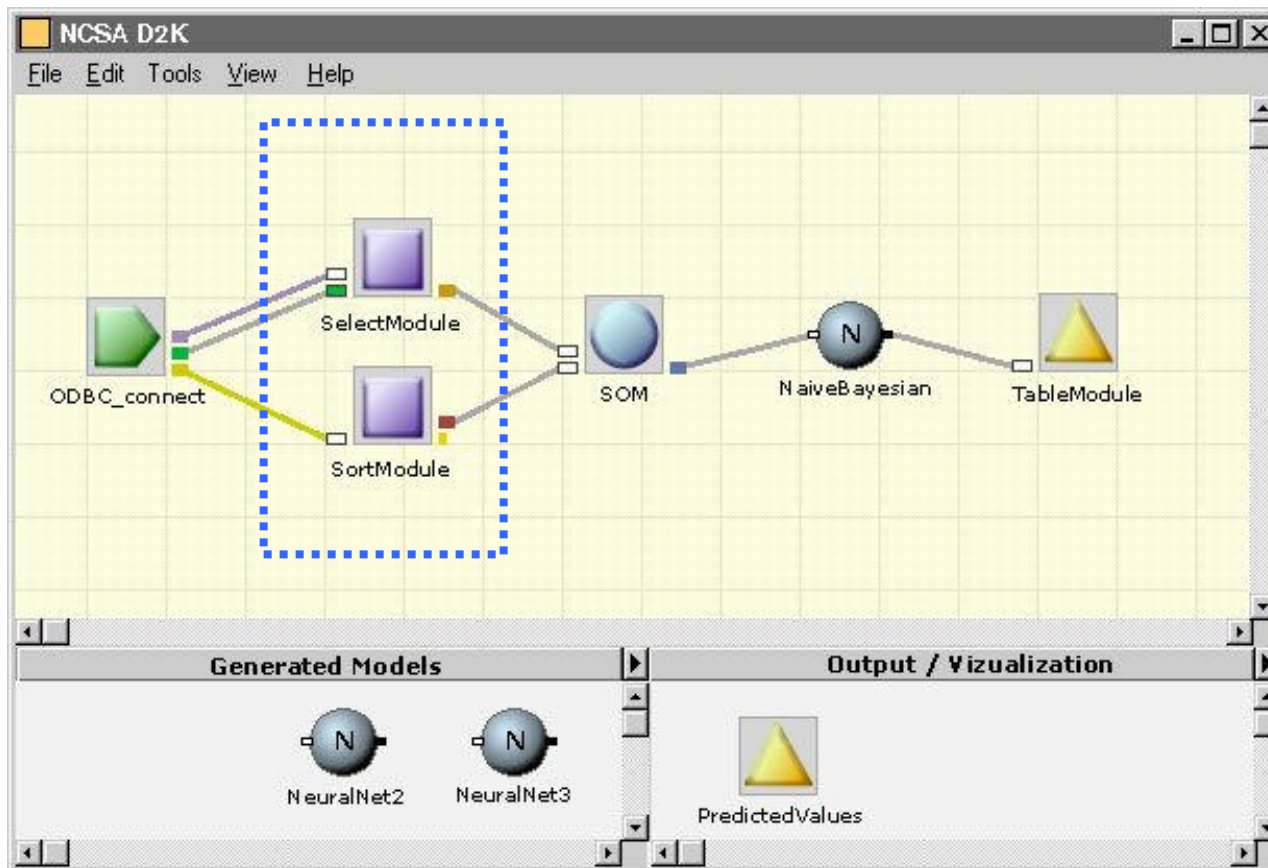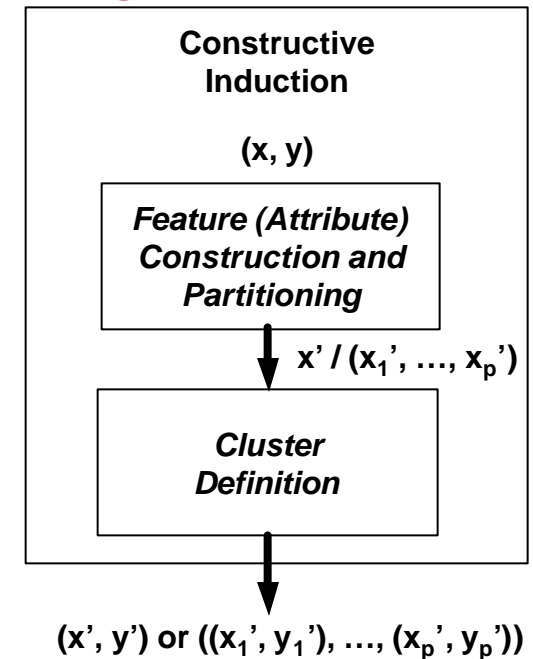| Caution/Warning | Profilometer | Fuel Systems | Timing | Spatial/GPS/ Navigation |
|---|---|---|---|---|
| Hydraulics | Data Bus/Control/ Diagnostics | Ballistics | Electrical | Unused |

# Data Aggregation and Sampling

# Unsupervised Learning

- **Unsupervised Learning in Support of Supervised Learning**
    - Given: *D* **?** labeled vectors (*x, y*)
    - Return: *D'* **?** <u>new training examples</u> (*x', y'*)
    - Constructive induction: <u>transformation</u> step in KDD
        - Feature "construction": generic term
        - Cluster definition

- **Feature Construction: Front End**
    - Synthesizing new attributes
        - Logical: $x_1$ **?** **?** $x_2$, arithmetic: $x_1 + x_5 / x_2$
        - Other synthetic attributes: $f(x_1, x_2, \ldots, x_n)$, etc.
    - Dimensionality-reducing projection, feature extraction
    - <u>Subset selection</u>: finding *relevant attributes* for a given target *y*
    - <u>Partitioning</u>: finding relevant attributes for given targets $y_1, y_2, \ldots, y_p$

- **Cluster Definition: Back End**
    - Form, segment, and label clusters to get <u>intermediate</u> targets *y'*
    - <u>Change of representation</u>: find good (*x', y'*) for learning target *y*

**Constructive Induction**

(x, y)

*Feature (Attribute) Construction and Partitioning*

$x' / (x_1', \ldots, x_p')$

*Cluster Definition*

$(x', y')$ or $((x_1', y_1'), \ldots, (x_p', y_p'))$

---

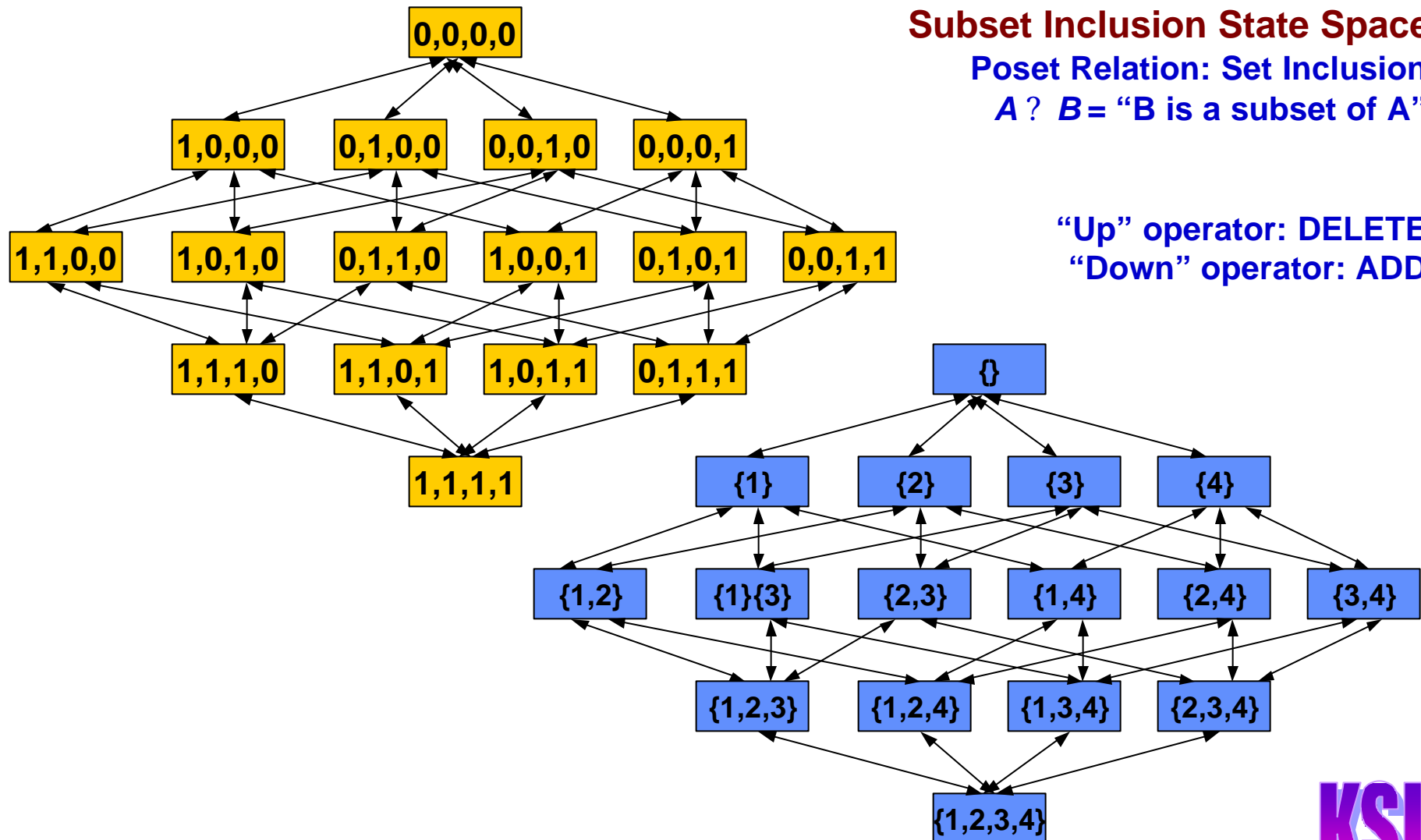**CIS 732: Machine Learning and Pattern Recognition**

# Relevance Determination
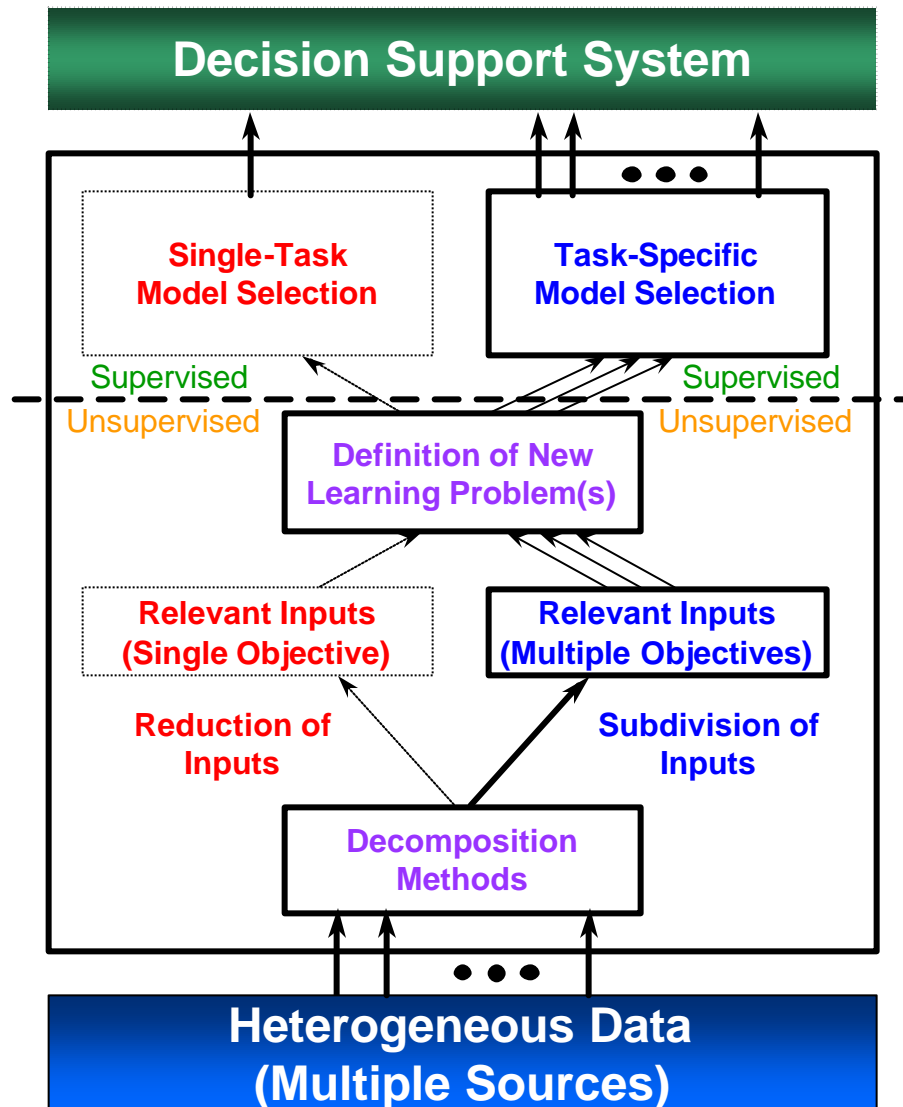
**Subset Inclusion State Space**

Poset Relation: Set Inclusion

$A$ ? $B$ = "B is a subset of A"

"Up" operator: DELETE
"Down" operator: ADD

# Wrappers for Performance Enhancement

**Decision Support System**

| Single-Task Model Selection | Task-Specific Model Selection |

Supervised — — — Supervised
Unsupervised — — — Unsupervised

**Definition of New Learning Problem(s)**

**Relevant Inputs (Single Objective)**

**Relevant Inputs (Multiple Objectives)**

**Reduction of Inputs**

**Subdivision of Inputs**

**Decomposition Methods**

**Heterogeneous Data (Multiple Sources)**

- **Wrappers**
  - "Outer loops" for improving inducers
  - Use inducer performance to optimize
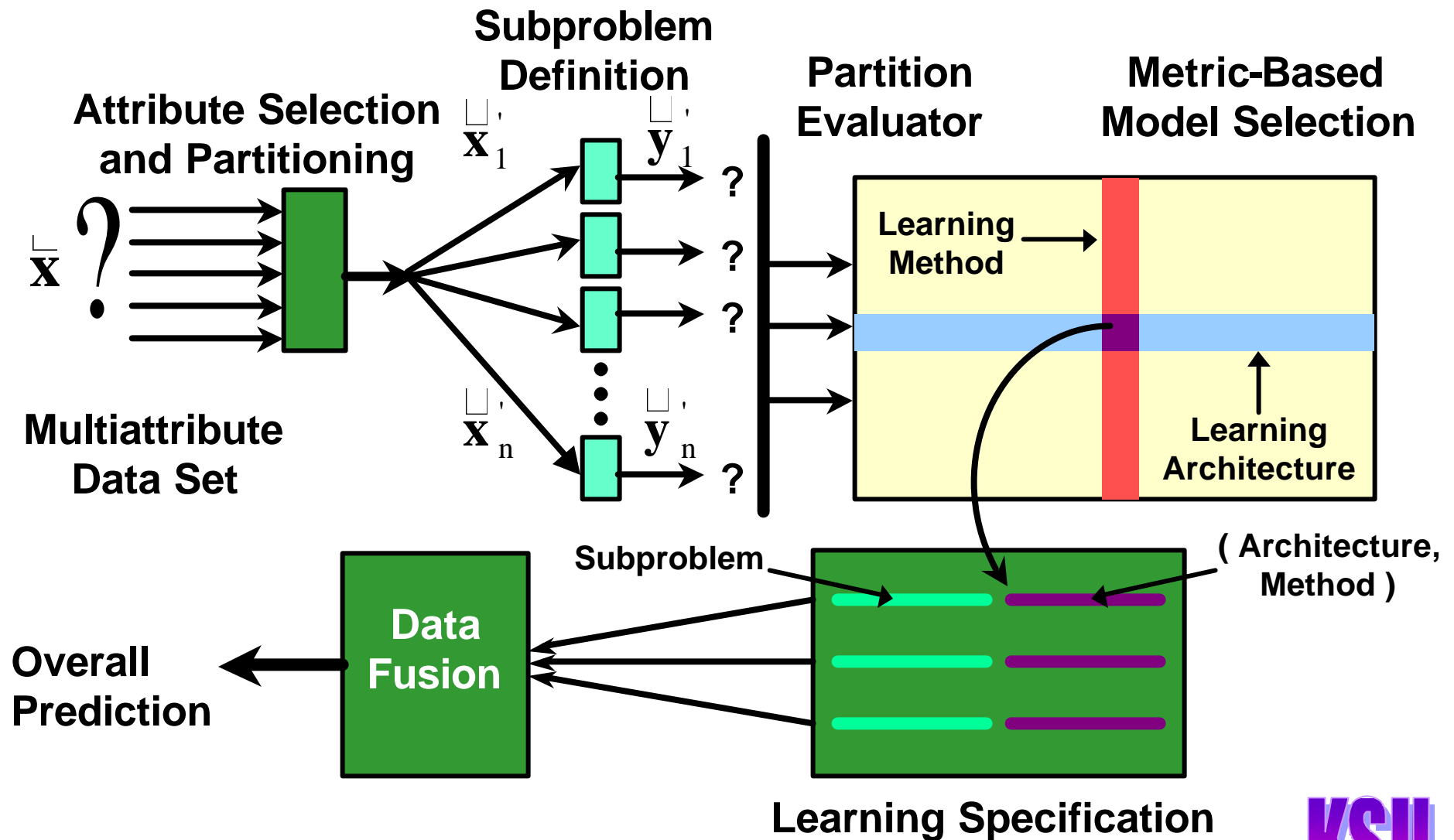
- **Applications of Wrappers**
  - Combining knowledge sources
    - Committee machines (static): bagging, stacking, boosting
    - Other sensor and data fusion
  - Tuning hyperparameters
    - Number of ANN hidden units
    - GA control parameters
    - Priors in Bayesian learning
  - Constructive induction
    - Attribute (feature) subset selection
    - Feature construction

- **Implementing Wrappers**
  - Search [Kohavi, 1995]
  - Genetic algorithm

**CIS 732: Machine Learning and Pattern Recognition**

# Supervised Learning Framework



**Subproblem Definition**

**Attribute Selection and Partitioning**

$\mathbf{x}'_1$   $\mathbf{y}'_1$   **?**

$\mathbf{x}$ **?**

$\mathbf{x}'_n$   $\mathbf{y}'_n$   **?**

**Multiattribute Data Set**

**Partition Evaluator**

**Metric-Based Model Selection**

Learning Method

Learning Architecture

**Subproblem**

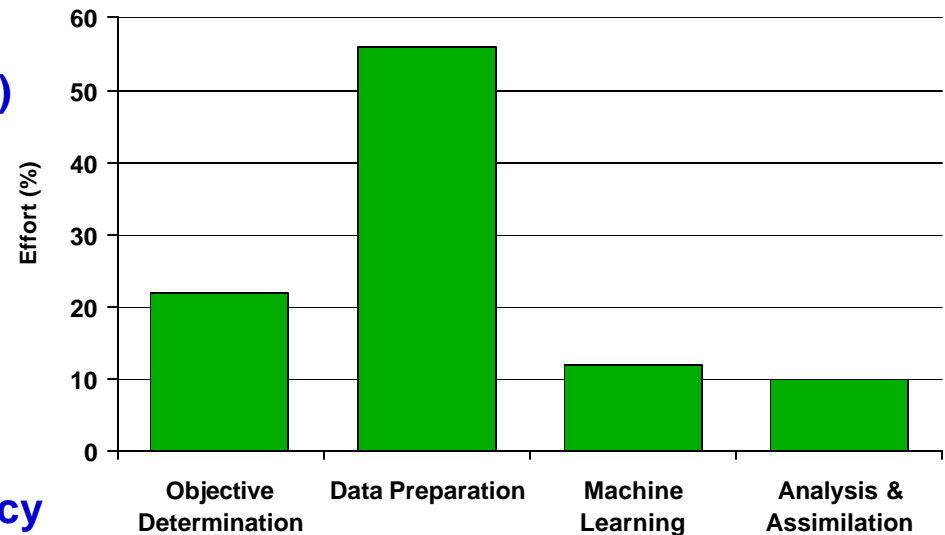( Architecture, Method )

**Data Fusion**

**Overall Prediction**
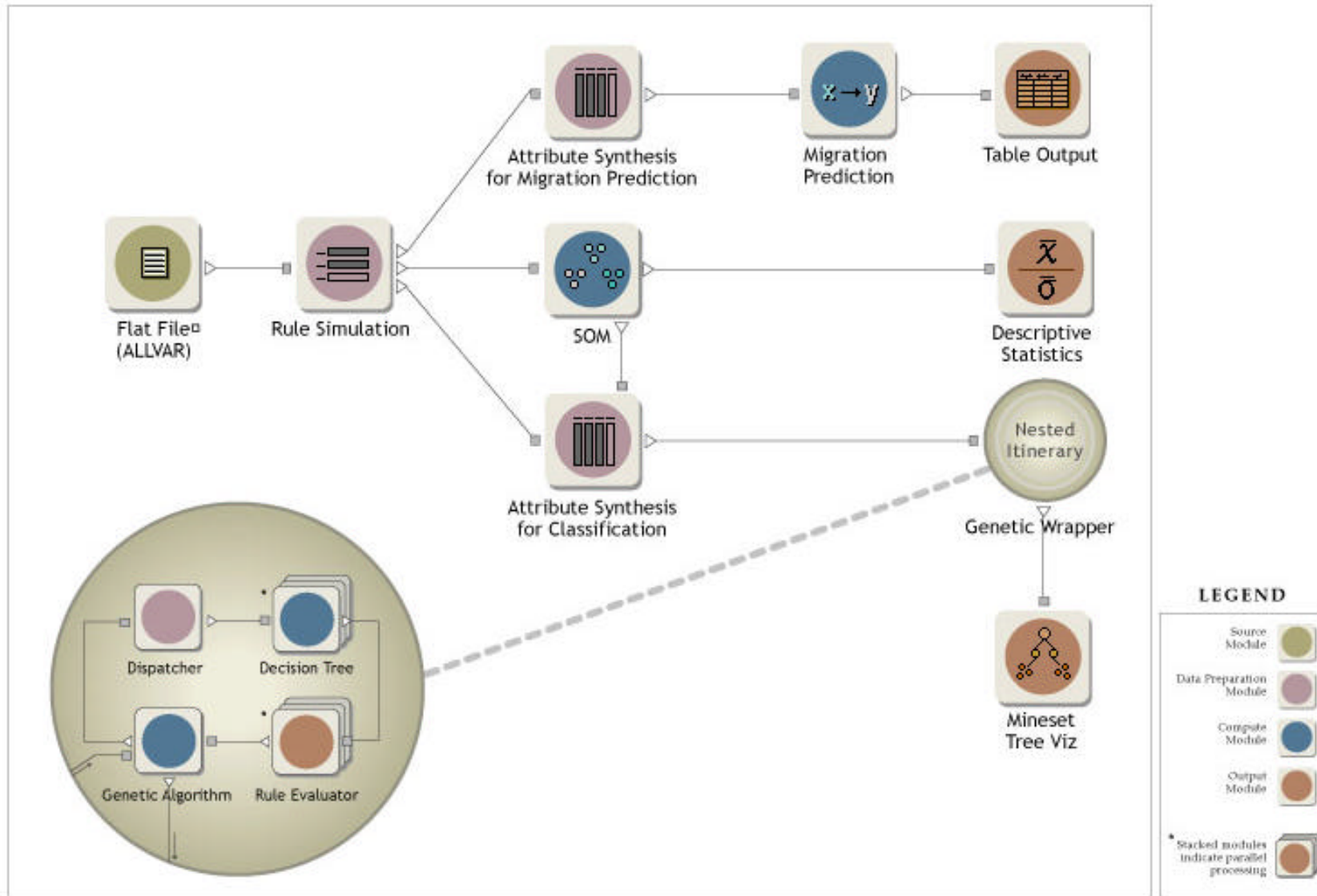
**Learning Specification**

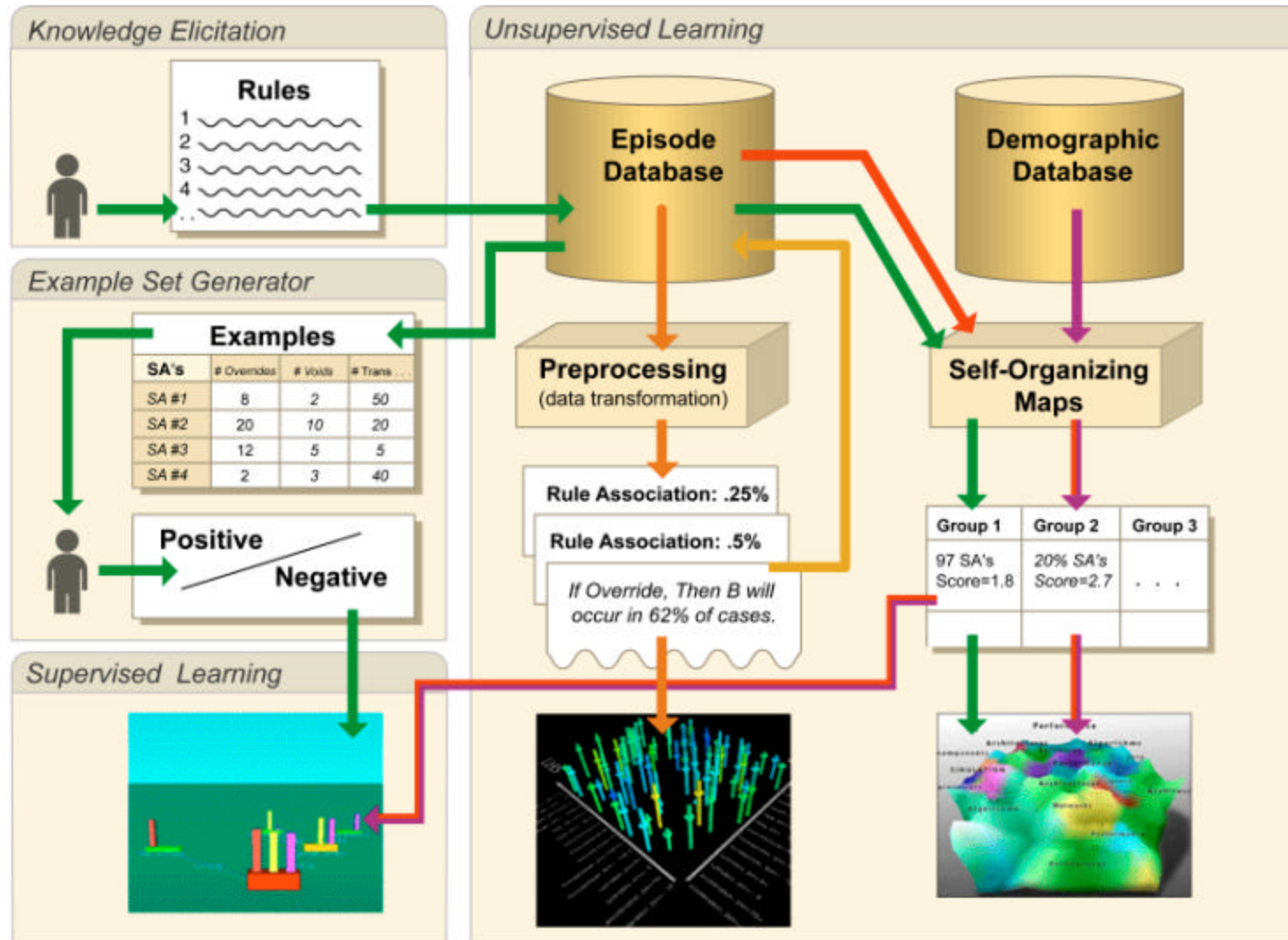# Performance Element: Decision Support Systems (DSS)

- **Model Identification (Relational Database)**
  - **Specify data model**
  - **Group attributes by type (dimension)**
  - **Define queries**
- **Prediction Objective Identification**
  - **Identify target function**
  - **Define hypothesis space**
- **Transformation of Data**
  - <u>**Reduce**</u> **data: e.g., decrease frequency**
  - <u>**Select**</u> ***relevant*** **data channels (given prediction objective)**
  - <u>**Integrate**</u> **models, sources of data (e.g.,** ***interactively elicited rules***)
- **Supervised Learning**
- **Analysis and Assimilation: Performance Evaluation using DSS**



Bar chart — Effort (%) vs. Objective Determination, Data Preparation, Machine Learning, Analysis & Assimilation



Environment (Data Model) → Learning Element → Knowledge Base → Performance Element

**CIS 732: Machine Learning and Pattern Recognition**

# Case Study:
# Fraud Detection
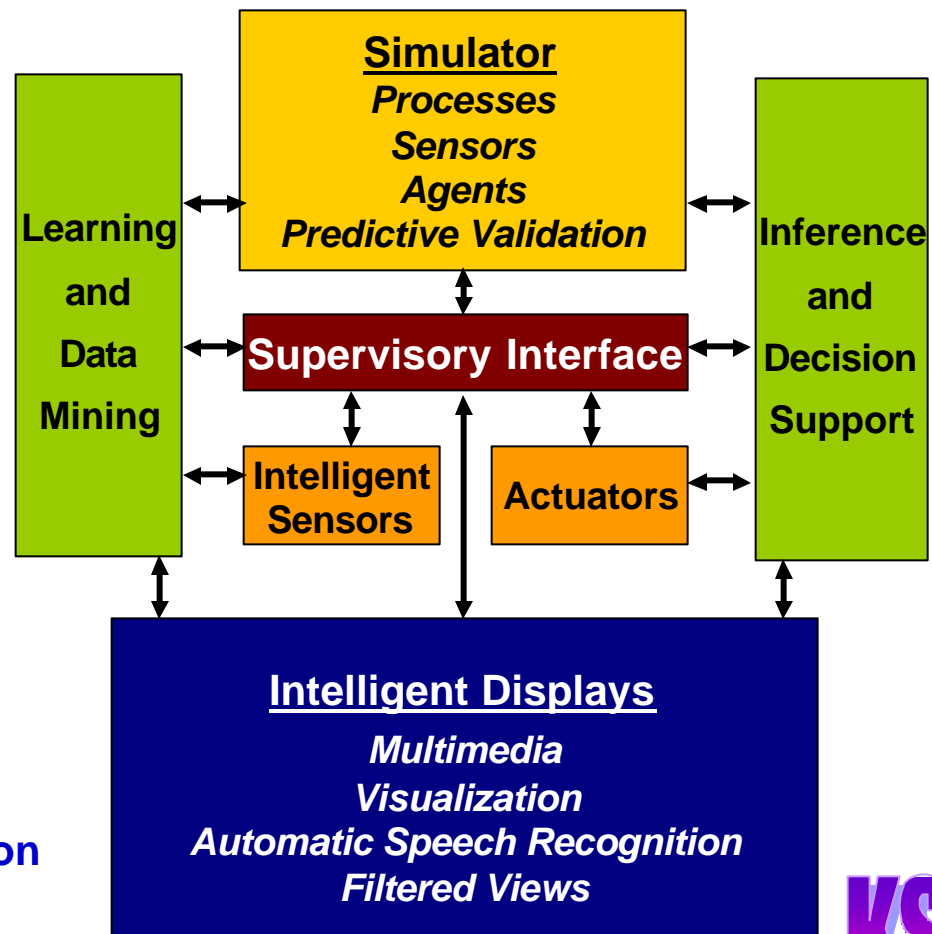
**CIS 732: Machine Learning and Pattern Recognition**

# Case Study: Prognostic Monitoring

- **Control Interfaces**
  - Actuators: fire/smoke suppression, electrical isolation, counterflooding
  - Intelligent sensors
- **Simulation Module**
  - Process/agent simulation
  - Automation simulation
  - Predictive validation for sensors
- **Learning Modules**
  - Time series learning
  - Control knowledge acquisition
- **Intelligent Reasoning Modules**
  - Crisis recognition
  - Casualty response
- **Intelligent Displays Module**
  - Interactive design and visualization
  - Supervisory interface

**Learning and Data Mining**

**Simulator**
*Processes*
*Sensors*
*Agents*
*Predictive Validation*

**Supervisory Interface**

**Intelligent Sensors**

**Actuators**

**Inference and Decision Support**

**Intelligent Displays**
*Multimedia*
*Visualization*
*Automatic Speech Recognition*
*Filtered Views*

**CIS 732: Machine Learning and Pattern Recognition**

# Terminology

- **Data Mining**
  - <u>Operational definition</u>: automatically extracting *valid*, *useful*, *novel*, *comprehensible* information from large databases and *using it to make decisions*
  - <u>Constructive definition</u>: expressed in stages of data mining

- **Databases and Data Mining**
  - <u>D</u>ata <u>B</u>ase <u>M</u>anagement <u>S</u>ystem (<u>DBMS</u>): data *organization, retrieval, processing*
  - <u>Data warehouse</u>: repository of integrated information for queries, analysis
  - <u>O</u>nline <u>A</u>nalytical <u>P</u>rocessing (<u>OLAP</u>): storage/CPU-efficient manipulation of data for summarization (descriptive statistics), inductive learning and inference

- **Stages of Data Mining**
  - <u>Data selection</u> (*aka* <u>filtering</u>): sampling original (<u>raw</u>) data
  - <u>Data preprocessing</u>: sorting, segmenting, aggregating
  - <u>Data transformation</u>: change of representation; feature construction, selection, extraction; <u>quantization</u> (<u>scalar</u>, e.g., <u>histogramming</u>, <u>vector</u>, *aka* <u>clustering</u>)
  - <u>Machine learning</u>: unsupervised, supervised, reinforcement for model building
  - <u>Inference</u>: application of performance element (pattern recognition, *etc.*); evaluation, assimilation of results

**KSU**

**Kansas State University**
**Department of Computing and Information Sciences**

# Summary Points

- **Knowledge Discovery in Databases (KDD) and Data Mining**
  - <u>Stages</u>: selection (filtering), processing, transformation, learning, inference
  - Design and implementation issues
- **Role of Machine Learning and Inference in Data Mining**
  - Roles of unsupervised, supervised learning in KDD
  - Decision support (information retrieval, prediction, policy optimization)
- **Case Studies**
  - Risk analysis, transaction monitoring (filtering), prognostic monitoring
  - Applications: business decision support (pricing, fraud detection), automation
- **Resources Online**
  - Microsoft DMX Group (Fayyad): http://research.microsoft.com/research/DMX/
  - KSU KDD Lab (Hsu): http://ringil.cis.ksu.edu/KDD/
  - CMU KDD Lab (Mitchell): http://www.cs.cmu.edu/~cald
  - KD Nuggets (Piatetsky-Shapiro): http://www.kdnuggets.com
  - NCSA Automated Learning Group (Welge)
    - ALG home page: http://www.ncsa.uiuc.edu/STI/ALG
    - NCSA *D2K*: http://chili.ncsa.uiuc.edu

**CIS 732: Machine Learning and Pattern Recognition**