# Data Mining for CEOP Data

## October 9, 2002

### Shin-ichi Sobue/NASDA

# Purpose of Data Mining

"The extraction of hidden predictive information from large databases."

Hidden!? Data

# Example - Medical Research

July 18, 2002

Healthy Lifestyle May Help Cut

Alzheimer's Risk, Scientists Find

Study of Alzheimer's - 1,449 people for 21 years.

High cholesterol and high blood pressure raise the risk of developing Alzheimer's more than carrying ApoE-e4, a gene variation considered the most important genetic risk factor for the disease.

High cholesterol and high blood pressure!

On Thursdays

Diapers -->  Beer
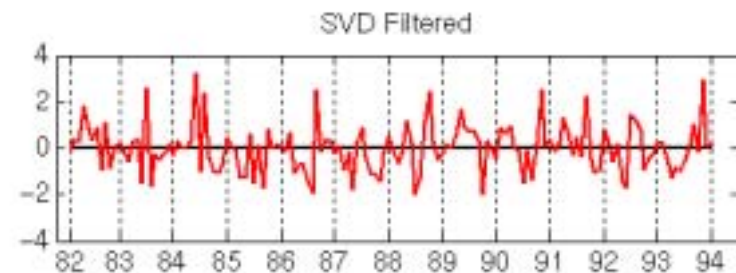
1. Displays of Beer and Diapers close together (near the registers).
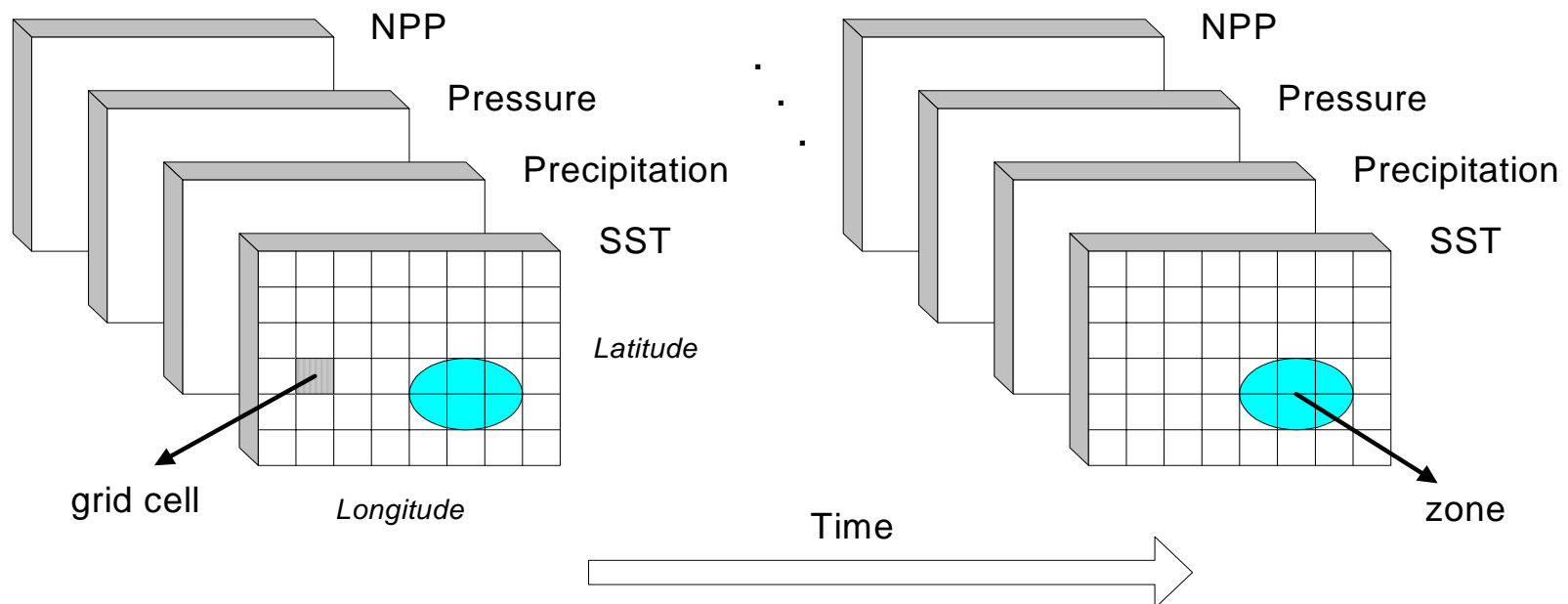
2. ?

- **Remove seasonality to see events (anomalies) of interest.**
  - Monthly Z Score
    - Subtract monthly mean and divide by monthly standard deviation
  - 12 month moving average
  - Discrete Fourier Transform
  - Singular Value Decomposition

# Time Series Data

- Global snapshots of values for a number of variables on land surfaces or water.

- Monthly values over a range of 10 to 50 years.

- Gridded values (NPP $0.5° \times 0.5°$,  SST $1° \times 1°$)



NPP

Pressure

Precipitation

SST

Latitude

grid cell

Longitude

NPP

Pressure

Precipitation

SST

zone

Time

# Data Clustering

- Cluster Formation
  - Find regions of the land or ocean which have similar time series behavior.

- Teleconnections
  - Teleconnections - time series behavior is (very) similar over widely separated points on the Earth.

- Interested in relationships between regions, not "points."

- For ocean, clustering based on SST (Sea Surface Temperature), SLP (Sea Level Pressure), etc.

- For land, clustering based on NPP or other variables, e.g., precipitation, temperature.
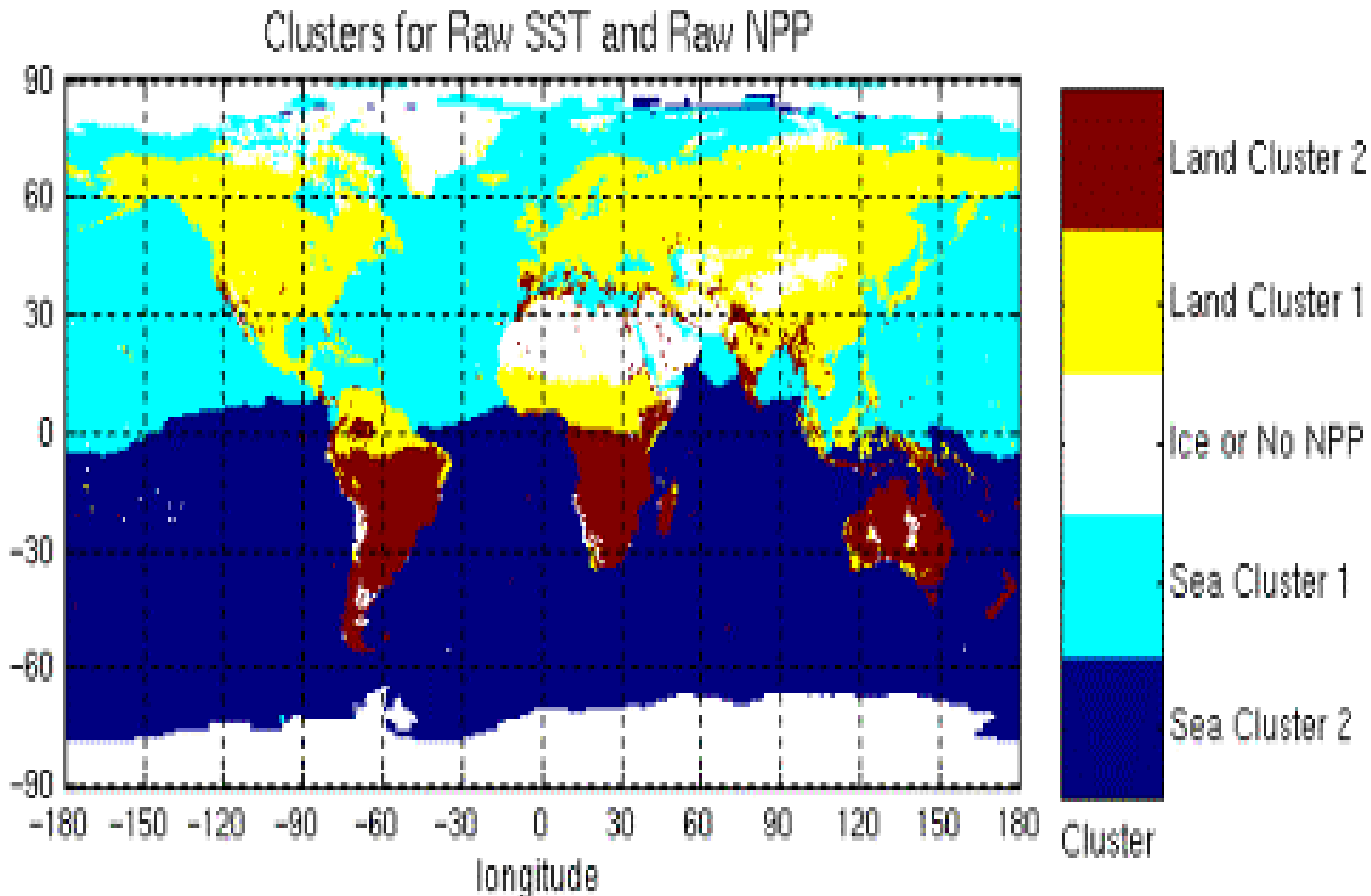
- The K-means algorithm
  - "K" is number of (nonoverlapping) clusters
  - "means" - center of the cluster is the mean or median of the "nearness"of the points in its cluster, where "nearness" is defined by a similarity function (Pearson's correlation coefficient).
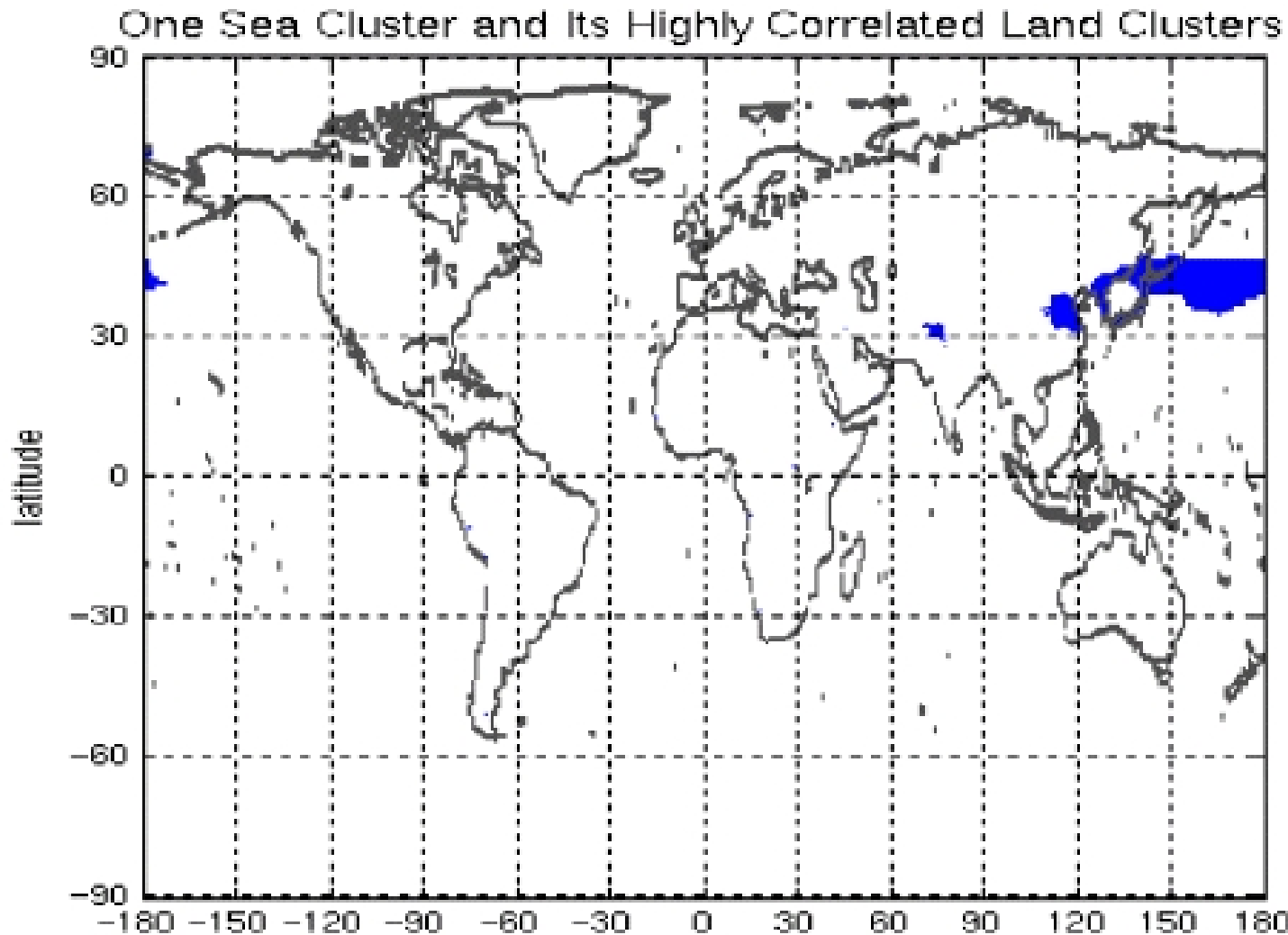
# K-Means Clustering of Raw NPP and Raw SST



Clusters for Raw SST and Raw NPP

Land Cluster 2
Land Cluster 1
Ice or No NPP
Sea Cluster 1
Sea Cluster 2

Cluster

longitude

# Discovering new Teleconnections

- Find land and sea clusters that are highly correlated, to identify new potential teleconnection patterns.

  – Produced 100 clusters for the land (NPP) and 100 clusters for the sea (SST).

  – Calculate correlations between sea and land clusters.

  – Next page shows diagram of sea cluster 19 and land clusters 56 and 58 (correlations of 0.56 and 0.50).

- Sea cluster 19 is highly correlated (-0.77) with the Pacific Decadal Oscillation

# One Sea Cluster and Associated Land Clusters



One Sea Cluster and Its Highly Correlated Land Clusters

12

# Discovery of Ocean Climate Indices

- 1. SOI ( Southern Oscillation Index)
- 2. NAO (North Atlantic Oscillation)
- 3. AO (Artic Oscillation)
- 4. PDO (Pacific Decadel Oscillation)
- 5. QBO (Quasi-Biennial Oscillation Index )
- 6 .CTI (Cold Tongue Index)
- 7. WP (Western Pacific)
- 8. ANOM12  (Normalized version of NINO12)
- 9. ANOM3   (Normalized version of ANOM3)
- 10. ANOM4 (Normalized version of NINO4)
- 11. ANOM34 (Normalized version NINO34)
- (Note: 1, 6, and 8-11 are El Nino related indices)

# Discovery of Ocean Climate Indices

- **SST clusters are potential Ocean Climate Indices (OCIs).**

- **Determine if the clusters match known OCIs.**

- **Evaluate the influence of clusters (potential OCIs) on land points.**

- **For clusters that don't match know OCIs, consider them to be potential OCIs and conduct further analysis to see if they are useful and interesting.**
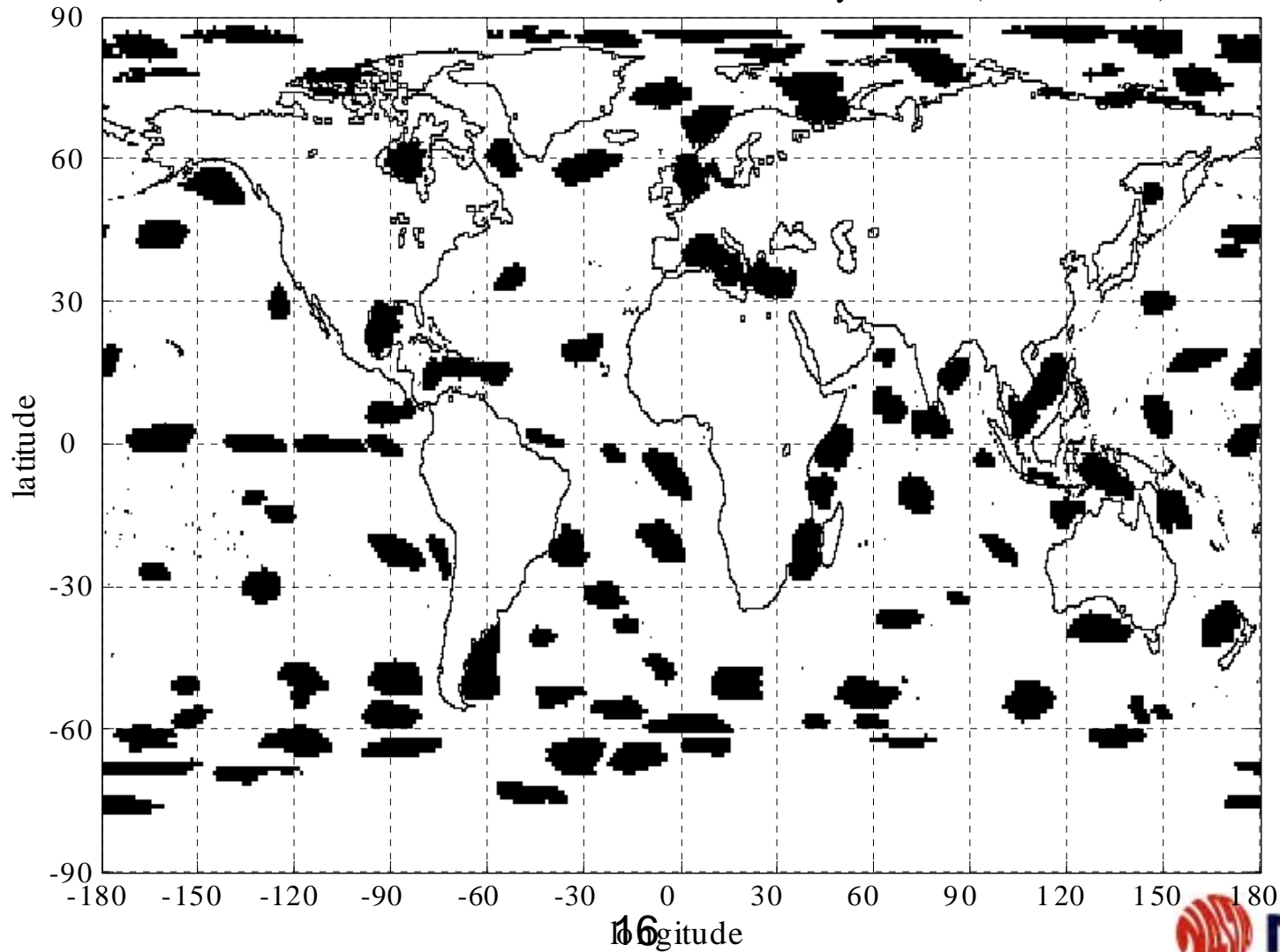
- **Find the nearest neighbors of each data point.**
  - **In this case data points are time series.**
  - **Examine the similarity between pairs of points in terms of how many nearest neighbors two points share.**
- **Eliminate noise, which are points with low density.**
- **Build clusters around the core points, which are points with high density.**
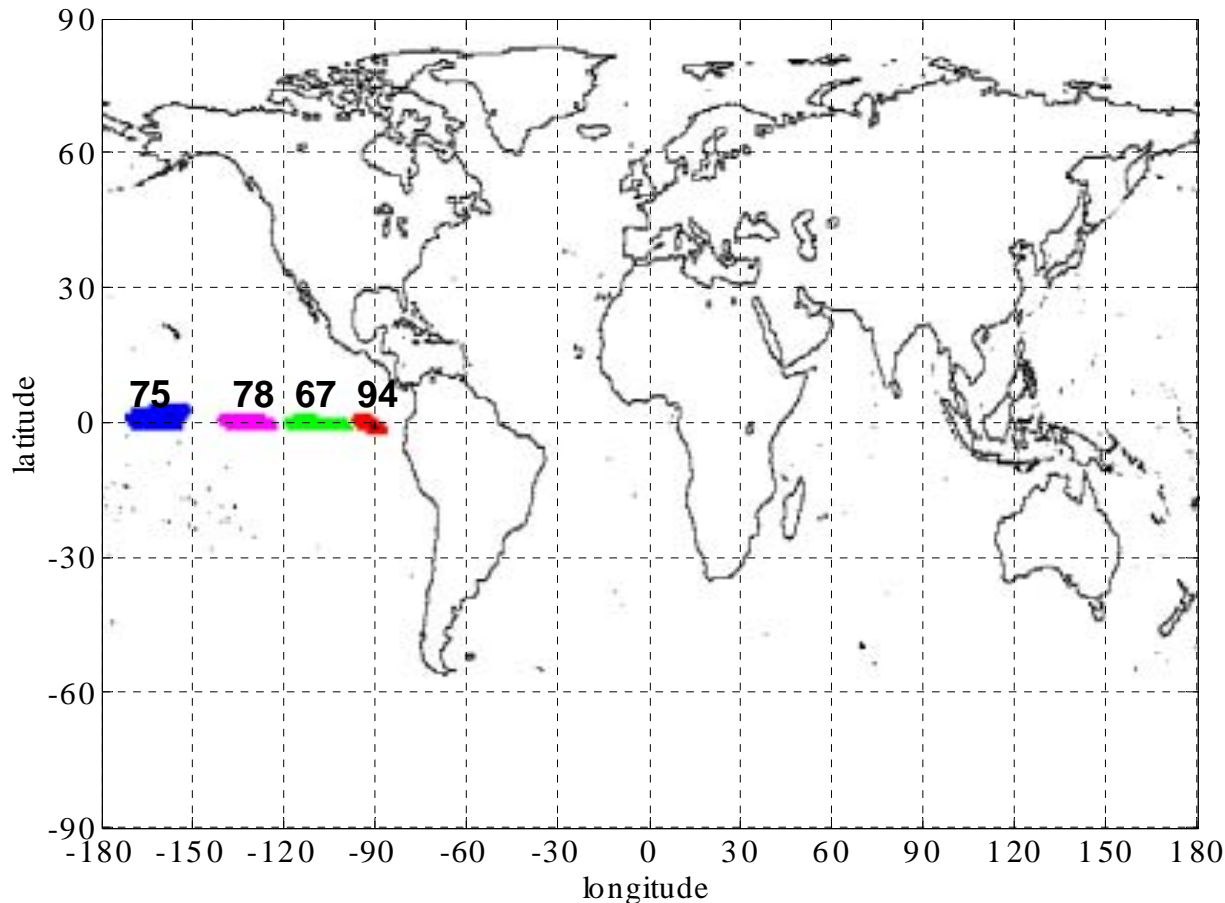
# SST Clusters



107 SNN Clusters for Detrended Monthly Z SST (1958-1998)

# SST Clusters that Correspond to El Nino Climate Indices
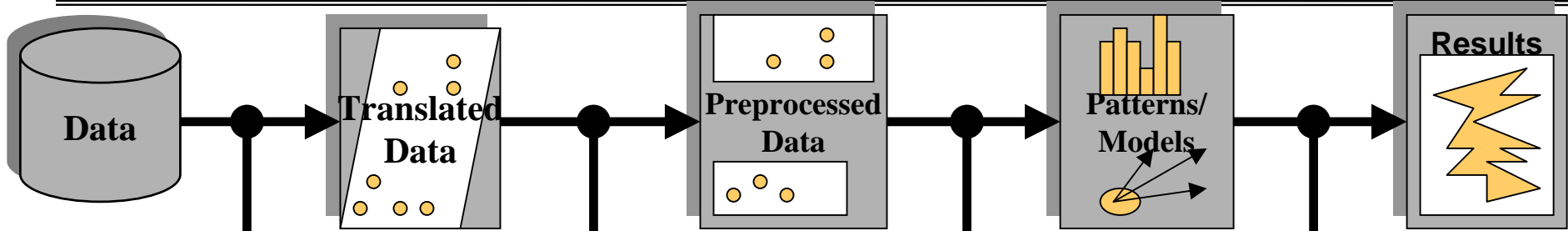
## EL Nino Related SST Clusters



| Niño Region | Range Longitude | Range Latitude |
|---|---|---|
| 1+2   (94) | 90°W-80°W | 10°S-0° |
| 3      (67) | 150°W-90°W | 5°S-5°N |
| 3.4    (78) | 170°W-120°W | 5°S-5°N |
| 4      (75) | 160°E-150°W | 5°S-5°N |

El Nino Regions Defined by Earth Scientists

SNN clusters of SST that are highly correlated with El Nino indices, ~ 0.93 correlation.

# ADaM
## Algorithm Development and Mining System

**Data** → **Translated Data** → **Preprocessed Data** → **Patterns/ Models** → **Results**

**Processing**

| Input | Preprocessing | Analysis | Output |
|---|---|---|---|
| HDF<br>HDF-EOS<br>GIF PIP-2<br>SSM/I Pathfinder<br>SSM/I TDR<br>SSM/I NESDIS Lvl 1B<br>SSM/I MSFC<br>  Brightness Temp<br>US Rain<br>Landsat<br>ASCII Grass<br>Vectors (ASCII Text)<br><br>Intergraph Raster<br>Others... | Selection and Sampling<br>  Subsetting<br>  Subsampling<br>  Select by Value<br>  Coincidence Search<br>Grid Manipulation<br>  Grid Creation<br>  Bin Aggregate<br>  Bin Select<br>  Grid Aggregate<br>  Grid Select<br>  Find Holes<br>Image Processing<br>  Cropping<br>  Inversion<br>  Thresholding<br>Others... | Clustering<br>  K Means<br>  Isodata<br>  Maximum<br>Pattern Recognition<br>  Bayes Classifier<br>  Min. Dist. Classifier<br>Image Analysis<br>  Boundary Detection<br>  Cooccurrence Matrix<br>  Dilation and Erosion<br>  Histogram<br>Operations<br>  Polygon<br>Circumscript<br>  Spatial Filtering<br>  Texture Operations<br>Genetic Algorithms<br>Neural Networks<br>Others... | GIF Images<br>HDF-EOS<br>HDF Raster Images<br>HDF SDS<br>Polygons (ASCII, DXF)<br>SSM/I MSFC<br>  Brightness Temp<br>TIFF Images<br>Others... |

# ADaM - Input

**Data Readers and Writers**

- **Binary**
- **GIF**
- **HDF**
- **HDF-EOS**
- **TIFF**
- **ASCII**

# ADaM - Preprocessing

- **Grid Operations**
- **Subsetting**
- **Image Processing**

| Collage | Invert |
| --- | --- |
| Crop | Overlay |
| Dilate | Pulse Coupled Neural Network |
| Equalize | Quantize |
| Erode | Rotate |
| Gabor | Spatial Filters |
| Filter | Statistics |

# ADaM - Analysis

## 1. Genetic Algorithms

## 2. Pattern Recognition

- Bayes Classifier
- Isodata clustering operation
- K Means
- Max/Min Operation
- Multiple Prototype Minimum Distance Classifier
- Decision Tree Classifier
- Oblique Decision Tree
- Recursively Splitting Neural Network
- Minimum Distance Classifier

# ADaM - Plan Builder Client

# Space Time Toolkit

# Space Time Toolkit

**Can visually integrate virtually any geographic data set, regardless of differences in their spatial and temporal representation.**

**Low-level sensor data**

**Gridded products**

**Vector based GIS data**

**3D model output, etc.**

- **Relationship Mining** - mining for various types of relationships includes most types of data mining.

  – **Cause or Indicator Relationships** - Find the actual causes of an "event" (e.g. determine if land clearing has disturbed hydrologic runoff and/or flood frequency).

  – **Effect Relationships** - The analysis of events having regional climate impact (e.g., volcanic eruptions, desert flash floods).

  – **Correlation Relationships** - Study of earthquakes (e.g. the space-time correlations involving triggered slips, foreshocks and aftershocks).

  – **Linkage Relationships** - An event at one time and place is related to an event at another time and place (e.g., El Nino and drought in Indonesia). Feature Transformation

  – **Prediction Relationships** - Predict storm tracts and changes in intensity.

# "Types" of Data Mining (continued)

- **Exploratory Pattern Mining** - Look for unexpected spatial/temporal patterns long-term data that cover a longer period (e.g. three to five years),

- **Complex Process Characterization** - Develop predictive statistical models that can be applied to areas such as seismic activity or the spreading of fire. Physical/mechanistic models that can be applied to areas such as earthquake fault modeling, interactions between crust/mantle and stress transfer.

1. WTF-CEOP website: list information on research in data mining of Earth Science data.

2. WTF-CEOP website: list examples of data mining.

***Example of data mining query***:

- Investigation: It has been shown that warming in the equatorial Pacific Ocean in the El Nino area affects the strength of winds at the 200 mb level in the tropical Atlantic.  During drought years in West Africa, winds increase from west to east at the 200 mb level, in the Central Atlantic.

- Query: The user wants to find the time periods where the SST anomaly in the El Nino region, is less than -1°C and winds, at 200 mb, in the Central Atlantic, are greater than   -5.00 m/s.  For those time periods the user then wants to study the precipitation and vegetation data, for the western Sahel.

3. Survey CEOP scientists for what type of data mining would be useful to them. Add these queries to item 2 above.

4. WTF-CEOP website: develop an on-line library of examples of visualization techniques for 3D/4D data and data mining.

5. Prepare a CEOP database testbed (has the same software environment as CEOP).  This database testbed will be used to adapt and test CEOS data mining algorithms before they are installed on CEOP.

6. Identify "candidate" CEOP algorithms - algorithms that have been developed by (NASA, ESA, etc.) scientists that can be used by CEOP.  NASDA will install and test the algorithms on the CEOP database testbed.

7. Survey COTS Data Mining software.