

# Parallel Data Mining on a Beowulf Cluster

Peter Strazdins, Peter Christen, Ole M. Nielsen and Markus Hegland

<http://cs.anu.edu.au/~Peter.Strazdins> (/seminars)

Data Mining Group

Australian National University, Canberra

<http://csl.anu.edu.au/ml/dm/>

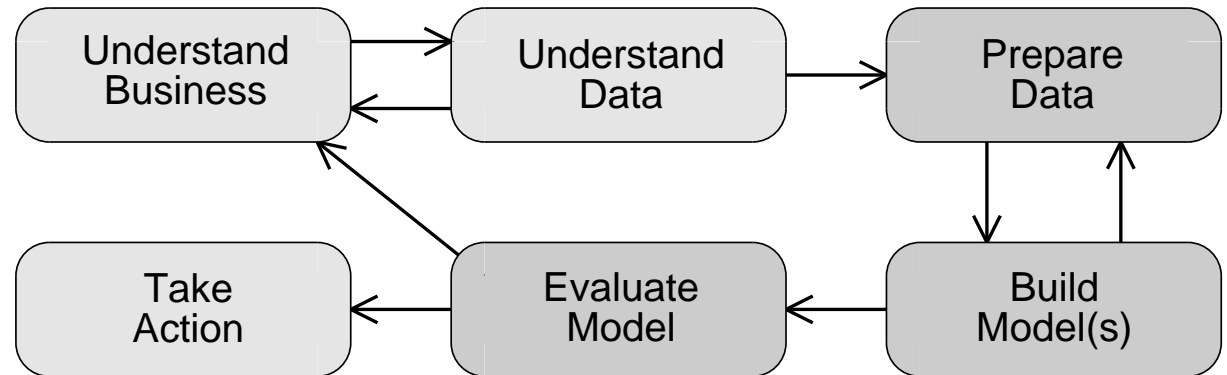
*5th Int'l Conference on High-Performance Computing in the Asia-Pacific Region*

*25 September 2001*

# Talk Outline

- The Data Mining Process
  - Data Mining Issues
  - What we do at the ANU
- Predictive Modelling
- Parallel Implementation
  - Matrix Assembly and Linear System Solve
  - The ANU Beowulf *Bunyip*
  - Performance Results
- Outlook

# *The Data Mining Process*



- Analysis of large and complex data sets
- Find previously unknown relationships and patterns that are useful
- Data Mining is iterative and interactive
- Challenges: data size and data complexity

# *Data Mining Issues*

- Large data size and data complexity require
  - more computational power
  - higher memory and I/O bandwidth
  - more secondary storage
- Iterative and interactive data mining process requires rapid prototype development
- Data security and privacy (personal data)
- Heterogeneity and distributed data collection

# *"What we do at the ANU"*

- Development of algorithms that are scalable both wrt. data size and number of processors
- Apply numerical techniques for predictive modelling  
(including thin plate splines, finite elements, wavelets and additive models)
- Data mining toolbox (DMtools)  
(for data exploration, analysis and preprocessing)
- Consultancies

*Data mining is one of 13 APAC Expertise Programs*

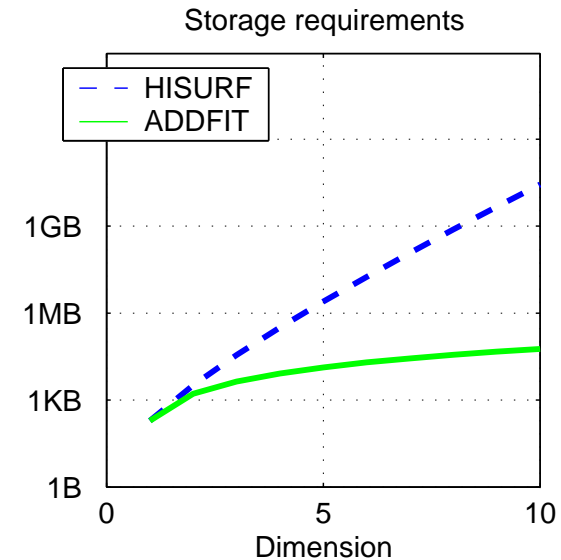
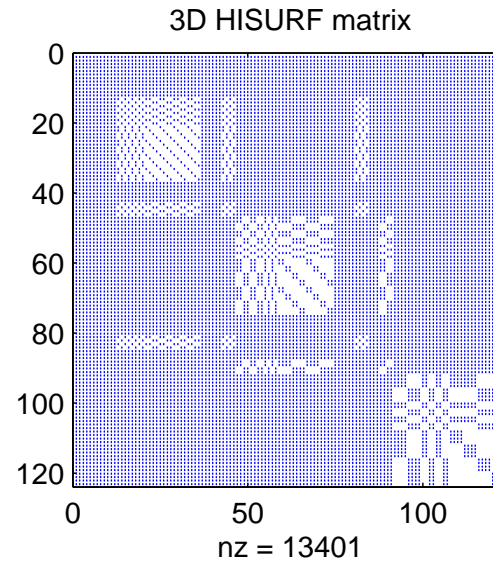
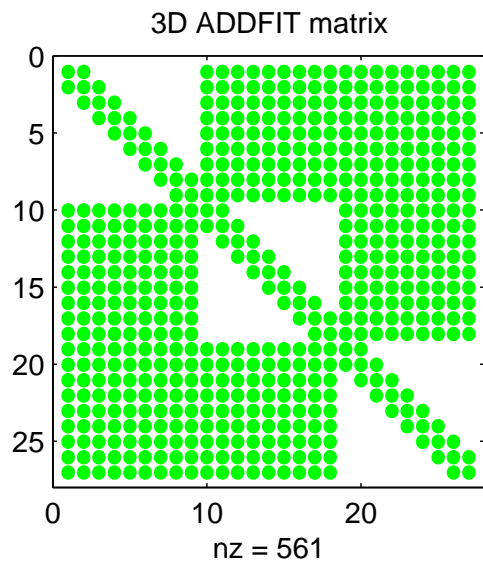
# *Three Predictive Modelling Methods*

- 1. ADDFIT: Additive models**  
Lowest computational costs, but coarsest approximations, suitable for high-dimensional problems
- 2. HISURF: High dimensional surface smoothing**  
Uses a hierarchical interpolatory wavelet basis, suitable for three to seven dimensional problems
- 3. TPSFEM: Thin plate splines finite element method**  
Piecewise multilinear finite elements, most accurate approximation at highest computational costs, suitable for two and three dimensional problems

# *Advantages of these Methods*

- All three methods have two steps
  - 1) Read data and assemble a linear system
  - 2) Add constraints and solve the linear system
- The first step is easy and efficiently to parallelise (as long as matrix fits into main memory)
- Data has to be read from disk only once
- Size of the symmetric linear system is independent of the number of data records
  - ADDFIT and HISURF result in dense matrix
  - TPSFEM in larger sparse matrix(only ADDFIT and HISURF are presented here)

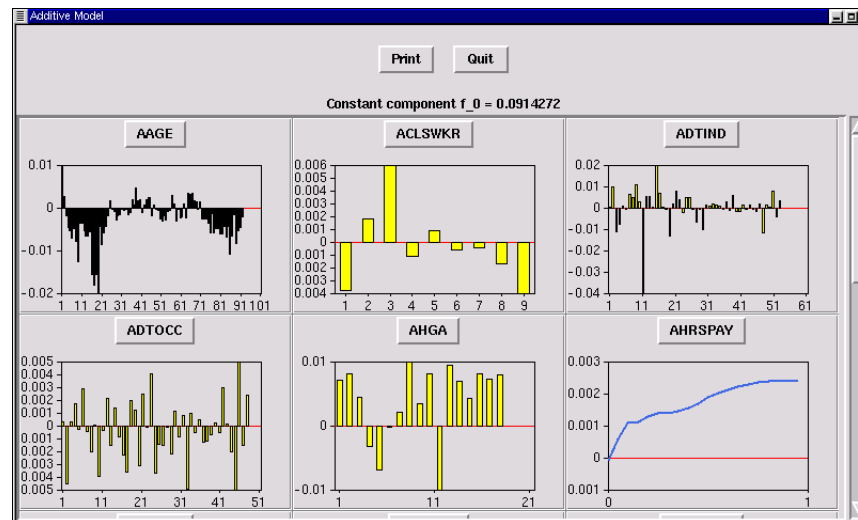
# Matrix Structures and Complexity



- A 3-dimensional examples
- Resolution is 9 grid points in each dimension
- HISURF matrices becomes much larger with increasing number of dimensions

# Parallel Cluster Implementation

- First (sequential) prototypes in *Matlab* and *Python*
- C/MPI code for higher performance
- Python/Tkl wrapper code to facilitate user interface and present graphical output



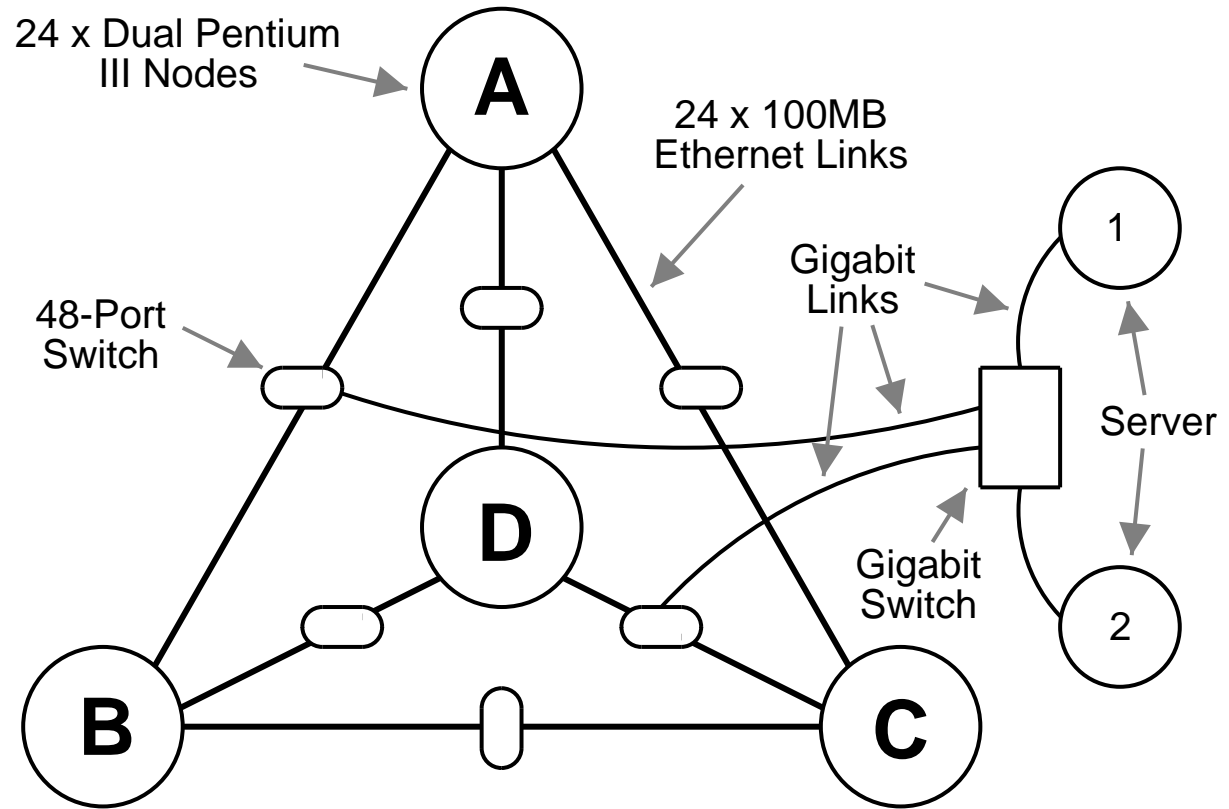
# *Parallel Assembly Step*

- Data set is initially copied to (local) disks of each processor (one-off cost)
- Each node reads a fraction  $n/p$  of the whole data set and assembles a local linear system
- Each data record adds some non-zero elements into the matrix (at data dependent locations)
- The complete matrix data structure is needed on each node
- The local linear systems are collected and summed after the assembly
- The final linear system can be solved sequentially or in parallel

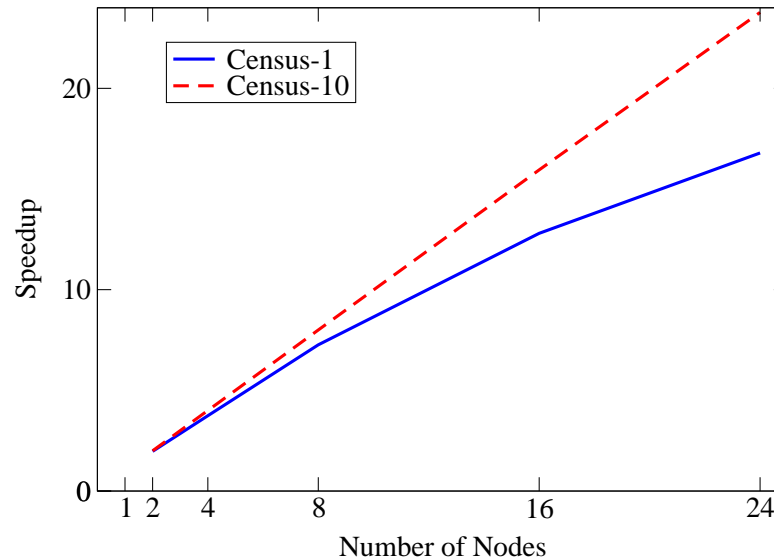
# Parallel Solve Step

- $P \times Q = p' \leq p$  processors are used with a square **block-cyclic matrix dist.**
- **HISURF**: uses a || LLT factorization / back-solve
- **AddFit**: uses a || LDLT factorization / back-solve
  - **Bounded Bunch-Kaufman** algorithm (Boeing, 1998) used: high accuracy
  - augmented with a **block search** algorithm (equally stable) to reduce || **symmetric interchange** overheads
  - other communication aspects are highly optimized (see LDLT paper in HPCAsia'01)
  - also has very high serial performance
- requires large  $\frac{N}{\sqrt{p'}}$  for good parallel efficiency

# The ANU Beowulf Cluster Bunyip

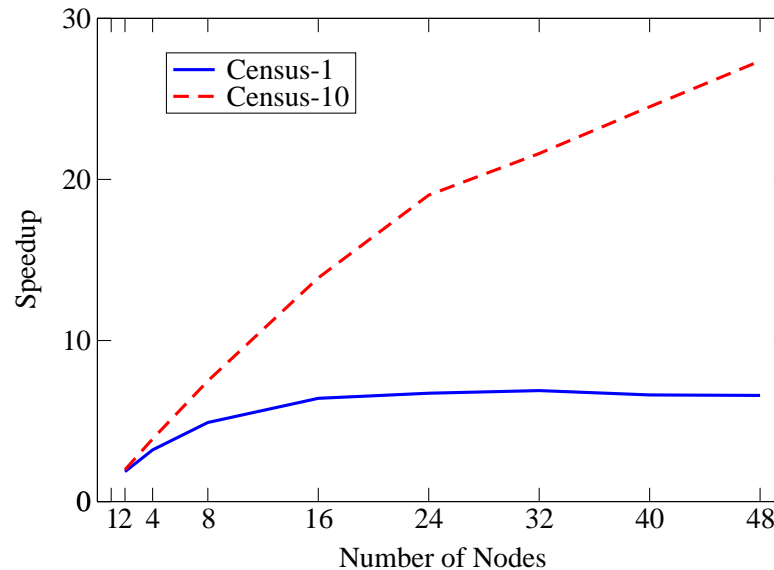


# ADDFIT Results on Bunyip



- Data set *Census* from UCI KDD archives (42 attributes,  $\sim 300,000$  records, 54 MBytes)
- 8 attributes used, matrix dimension  $100 \times 100$
- For *Census-1*, reducing and solving limits speedup

# ADDFIT Results on Bunyip (cont'd)



- Data set *Census* UCI KDD archives
- 41 attributes used, matrix dimension  $1000 \times 1000$
- Reducing (4 MBytes) and solving limits speedup (especially for *Census-1*)

# Current and Future Work

- Port to APAC National Facility (OpenMP / MPI)  
(first OpenMP prototype is working)
- Unifying ADDFIT, HISURF and TPSFEM into *adaptive sparse grids*
- Integrate into data mining toolbox *DMtools*
- Apply to other real world data
- Compare with other predictive modelling techniques (neural networks, decision trees, etc.)

Visit our web site at:

<http://csl.anu.edu.au/ml/dm/>