# Modelling and Optimising Flows Using Parallel Spatial Interaction Models

Ian Turton and Stan Openshaw

Centre for Computational Geography, School of Geography, University of Leeds,
Leeds, UK, LS2 9JT
email:{ian,stan}@geog.leeds.ac.uk , URL:http://www.geog.leeds.ac.uk/

**Abstract.** The paper demonstrates some of the benefits that high performance computing has to offer geographers. It reports the results of porting and then using very large spatial interaction models on the Cray T3D parallel supercomputer in modelling and spatial optimisation applications that would otherwise have been judged computationally infeasible.

## 1  Background

The entropy-maximising spatial interaction model is used for modelling many different types of flow data in both research and applied contexts. Examples of flow data are: journey to work trips, airline traffic, retail behaviour, world trade and telephone traffic in which flows of people, money, information etc., connect origin zones to destination zones. These zones have a geographical expression on maps whilst the flow data provides a useful summary of the very complex processes and behaviour patterns that generated them; for example, journey to work flows summarise the spatial dimensions of labour markets covering home-workplace interactions whilst in a retail context credit card purchases connect peoples home addresses to where they shop. Not surprisingly building computer models of flow data has many commercial applications. The purpose of this paper is to describe a few applications as a means of illustrating some of the benefits that high performance computing (HPC) can offer computationally minded geographers interested in this area of human systems modelling.

## 2  Flow data and spatial interaction models

The idea of developing a parallel spatial interaction model is not new. Harris (1985), discussed how (in theory) it could be ported onto a parallel machine. However, it seems that not until the early 1990s was the parallel implementation of the spatial interaction model developed, see Birkin et al (1995). The reasons for wanting to parallelise spatial interaction models are: to allow much bigger datasets to be processed more speedily than previously, to improve the underlying science both by being able to use finer resolution data, and to seek an improvement in the quality of the results by using more flexible methods of

parameter estimation and model optimisation: In the spatial interaction model the compute load is an approximate power function of the number of zones, small numbers of zones run well enough on a PC but large numbers (i.e. few thousand) require high performance computing hardware in the form of vector or parallel supercomputers. An ability to model large flow data matrices is now important because the availability and volume of flow data has rapidly increased as a result of developments in IT, viz. data warehousing and there is a prospect of UK flow data tables with up to 1.6 million origins and destinations becoming available soon (viz. credit card transactions) or 32 million origin and destinations (if telephone traffic were to be modelled at the subscriber level). However, currently in the UK the biggest public domain flow matrices are those from the special workplace statistics of the 1991 census (Openshaw, 1995). This is a 10,764 by 10,764 matrix of journey to work flows between all wards in Britain.

A simple origin constrained flow model developed using entropy-maximising methods is specified as follows:

$$T_{ij} = O_i D_j A_i e^{-\hat{\beta} C_{ij}} \tag{1}$$

$$A_i = 1/\sum_j D_j e^{-\hat{\beta} C_{ij}} \tag{2}$$

where $T_{ij}$ is the predicted number of trips from an origin place $i$ to a destination place $j$, $O_i$ is the "size" of $i$, $D_j$ is the "size" of $j$, and $C_{ij}$ is a measure of the distance or travel cost between $i$ and $j$. The parameter $\hat{\beta}$ is estimated to optimise the fit of the model using maximum likelihood or nonlinear least squares methods. Note that equation (2) ensures that the predicted flows satisfy the following accounting constraint:

$$\sum_j T_{ij} = O_i \tag{3}$$

A doubly constrained model is slightly more complex:

$$T_{ij} = O_i D_j A_i B_j e^{-\hat{\beta} C_{ij}} \tag{4}$$

$$A_i = 1/\sum_j D_j B_j e^{-\hat{\beta} C_{ij}} \tag{5}$$

$$B_j = 1/\sum_i O_i A_i e^{-\hat{\beta} C_{ij}} \tag{6}$$

# 3  Porting and scalability experiments

Fortran code was written for the two models so that they could cope with the 10,764 origin and destination zones of the 1991 ward level journey to work data. The two matrices used in the model to hold observed trips and costs ($T_{ij}$ observed, $C_{ij}$), could not be stored in the memory of a standard UNIX workstation, since they would require over one Gbyte of memory. A solution to this is to store

the trip matrix as a singly linked list which reduces the storage needed to about five Mbyte but this requires that the cost matrix can be recomputed when needed rather than stored. Unlike the trip data, the cost matrix is dense and cannot be stored in sparse format. Repeatedly recalculating the $C_{ij}$ values as they are needed dramatically reduces the memory requirements, however, this is achieved at the expense of a large increase in the amount of computation. On a Sunsparc 10/41 workstation the singly constrained one parameter model took 17.6 hours to run. The non linear optimisation routine performed 29 model evaluations in the search for an optimal parameter. The doubly constrained version ran for 264 hours. It is clearly difficult to make much use of a model that runs for 11 days before a single result is obtained. The question now is how much faster would a parallel version run and what new applications a dramatic speed-up in model performance might provide.

Initially the code was ported to the Kendall Square Research KSR1 parallel super computer at Manchester University (Openshaw and Sumner, 1995). This is a virtual shared memory MIMD with 64 processing units. Each processor is a 20 Mhz super-scalar RISC chip with a peak 64 bit floating point performance of 40 Mflop/s and 32 Mbyte of local memory.

The code was also ported onto the Cray T3D parallel supercomputer at Edinburgh. The Cray T3D version runs on 256 DEC alpha processors each rated at 150 Mflop/s with 64 Mbytes of memory, providing a peak theoretical speed of 38.4 Gflop/s for the whole machine. The same data parallel strategy was used for the KSR. In both cases the parallelism was at the outer DO LOOP level, so that in effect the computational load of the model was split into 10764 concurrent pieces, one for each i value in equation 1. This is a fairly fine grained level of within model parallelisation but was the highest level relevant to this model in a parameter estimation context. Here the model subroutine was going to be called 50 to 100 times by a nonlinear optimiser as it searches for an optimal parameter values and for this application there is no benefit in seeking to parallelise the nonlinear optimisation algorithm itself.

Figure 1 shows the performance for both models with comparisons being made to the KSR1. As can be seen performance scales extremely well and whilst hardly surprising this is very reassuring.

It is interesting to note that one doubly constrained model evaluation on a SunSparc 10/41 workstation took 91 hours of compute time compared to 171 seconds on the 256 processor T3D. It is now possible for the first time to easily model journey to work patterns for the whole of Britain at the finest available level of geographic data resolution. Furthermore, the compute times are now sufficiently small to allow newer and more complex types of models to be investigated which are more compute intensive and more advanced forms of parameter estimation procedures to be applied; (Diplock and Openshaw 1996).
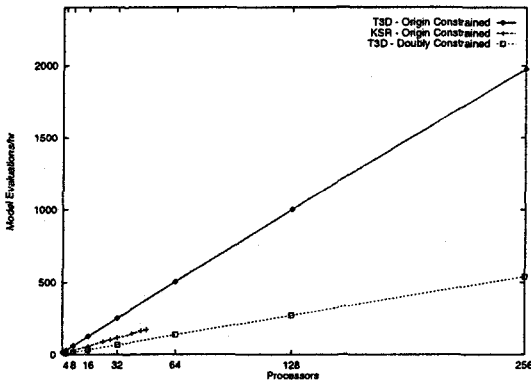
**Fig. 1.** Speed-up of models

# 4 Retail network optimisation

A major applied use for these models is to embed them in a spatial optimisation framework (Birkin, et al., 1995). For example in a retail context or hospital planning situation, it is often of interest to determine which set of $D_j$ values (in equation 1) will yield maximum profits or optimise some other performance indicator of global network benefits. This involves solving a discrete nonlinear combinatorial programming problem using various computationally intensive heuristics; e.g. genetic algorithms, tabu search and simulated annealing. The quality of the results now depends heavily on, put crudely, the number of times a model such as that in equation 1 can be evaluated in a fixed time period on a realistically sized spatial interaction dataset.

Consider the problem of optimising a national retail network using data for 2755 consumer origin zones and 822 destinations. The aim is to solve a massive combinatorial optimisation problem that involves finding the best 60 sites from 822 alternative locations that maximise profits assuming consumers behave according to the spatial interaction model. There are 822!/60! possible solutions. This spatial network optimisation model can be written as:

$$\text{maximise} \sum_{i=1}^{2755} \sum_{j=1}^{822} T_{ij}(\hat{D}_j^1) \tag{7}$$

$$T_{ij} = O_i D_j^1 A_i e^{-\hat{\beta} C_{ij}} \tag{8}$$

$$A_i = 1 / \sum_{k=1}^{2} \sum_{j}^{822} D_j^k e^{-\hat{\beta} C_{ij}} \tag{9}$$

The aim is to define a set of 0, 1 values for variable $D_j^1$ such that the objective function is optimised and $\sum_{j=1}^{822} D_j^1 = 60$. The $D_j^2$ values represent existing competitor sites and are fixed. Obviously different sets of $D_j^1$ will alter the total

market share attracted to these locations. Other constraints on the distribution of $D_j^1$ values may also be imposed; for example, minimum distances between nearest non-zero $D_j^1$ values should exceed a distance threshold.

When the serial code for this model was parallelised at the $i$ loop level the performance no longer scaled with the number of processors. This problem was smaller than previously and the code was also heavily optimised with the result that the amount of work being performed in the parallel regions rapidly became less than the communication overheads. As the number of processors increased, the performance peaked (at 16 processors) and then got worse. The problem can be temporarily fixed for the T3D by increasing the amount of work being done (i.e. removing some of the optimisation) but this defeats the objective of maximising the number of model evaluations per hour. Another solution was to make the entire model a parallel activity and not just the major loop. This was possible because of the nature of the optimisation task in particular the simulated annealing process used to optimise equation (7) was based on a single move heuristic that used a local neighbourhood search. The search neighbourhood could be of any size and by setting this to some integer multiple of the number of processors being used then high levels of parallel efficiency and linear scalability were achieved; see Table 1. This leads to a peak performance of 28.2 MFlops/PE which is good for a Cray T3D.

It is interesting to note that tuning the serial code for the T3D resulted in a single processor speed–up of 1097 times. Running on the T3D with 256 processors resulted in a combined speed–up of 2.8 million times compared to the original serial code. As well as a speed up the new algorithm also produced nearly twice as good a result as the serial version; see Figures 2a and 2b.


# 5  Conclusions


The paper describes the development of a parallel spatial interaction model and some applications involving two of the largest flow matrix yet modelled. The results are interesting both in their own right and also as a means of increasing awareness of what high performance computing (HPC) in general, and parallel super computing in particular, can offer. Increasingly geographical modellers are being freed from the computational constraints that have existed ever since the early years of the quantitative revolution. As Openshaw (1994b) points out, this could be the dawn of a major new era in geography, the beginning of what can be termed 'computational geography'; defined as new computationally intensive ways of doing geography. HPC is not just a means of making old models run two or three orders of magnitude faster but much more significantly it creates an opportunity to solve many previously unsolvable problems, to develop new computationally based technologies and to create entirely new ways of thinking about and doing all kinds of geography not all of which were previously quantitative. Hopefully the results reported here will help stimulate awareness and assist in this process of exploiting HPC in geography and the social sciences.
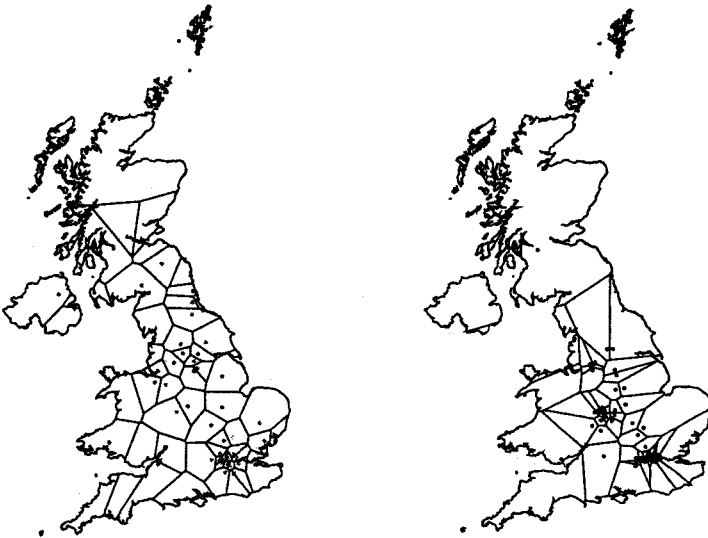
**Fig. 2.** a) The optimal dealer network found with the serial code, b) The optimal dealer network found using the improved algorithm on the Cray T3D

# References

Birkin, M., Clarke, M., and George F.: The use of parallel computers to solve nonlinear spatial optimisation problems: an application to network planning. Environment and Planning A, **27**,(1995) 1049–1068

Diplock, G., Openshaw, S.: Using simple Genetic Algorithms to calibrate Spatial Interaction Models. Geographical Analysis, (forthcoming)

Harris, B.: Some notes on parallel computing with special reference to transportation and land use modelling. Environment and Planning A, **17**, (1985), 1275–1278

Openshaw, S.:Some geographical applications of high performance computing. Working Paper 94/18, School of Geography, Leeds University, (1994a)

Openshaw, S.: Computational Human Geography: towards a research agenda. Environment and Planning A, **26**, (1994b), 499–505

Openshaw, S.: Census Users Handbook GeoInformation International, Cambridge, (1995)

Openshaw, S. and Sumner, R.: Parallel spatial interaction modelling on the KSR1-64 supercomputer, Working Paper 95/15, School of Geography, University of Leeds, (1995)