**Figure 7**. The geometrical query of CC-SIS

on can be logically separated into geometrical parts, image and thematic. In this paper the focus is put on the geometrical part of the database. Our pilot applications show that V3D is a suitable structure for the representation of 3-D objects, images and thematic data. It is possible to answer most of the questions about topology, position and shape of objects by means of geometric or SQL queries.

Further research will be concerned with the following themes:
(1) More extensive investigation of the possibilities of the implementation of V3D
(2) Further analysis of topology and queries. Special attention will be paid to the relationships between the different geometrical elements and object types.

## REFERENCES

[1] Breunig, M., 1996. "Integration of spatial information for geo-information systems", Springer Cop., Berlin.

[2] Fritsch, D., Pfannenstein, A., 1992. "Integration of DTM data structures into GIS data models", Int. Archives of Photogrammetry and Remote Sensing, Vol. XXIX, B3, pp. 497-503, Washington.

[3] Gruen, A., Wang, X., 1998. "CC-Modeler: A topology Generator for 3-D City Models", Int. Journal of ISPRS, Vol.53, No.5, pp.286-295.

[4] Molenaar, M., 1992. "A topology for 3D vector maps", ITC Journal, No. 1, pp. 25-33.

[5] Rikkers, R., Molenaar, M., 1994. "A query oriented implementation of a topologic data structure for 3-dimensional vector maps", International Journal of Geographical Information Systems, Vol. 8, No. 3, pp. 243-260.

---

# Effect of Compression on the Accuracy of DTM

Zhilin Li, W. K. Lam and Chengming Li

Dept. of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic University
Kowloon, Hong Kong

**Abstract**

In digital terrain modelling, if the source data (e.g. acquired by image matching) is too dense, a lossy compression procedure, e.g. the VIP (very important points) in ARC/TIN, is then applied to remove those less important data points (or select those important points) from the original data set for efficient data storage and process. This paper describes some experimental tests on the effect of such data compression on the accuracy of resultant digital terrain models (DTM). Results from two areas show that (a) the number of data points retained after compression is linearly proportional to the increase in threshold; (b) a compression of 30% doesn't have too much effect on DTM accuracy. An empirical model is established to predict the accuracy loss due to such compression.

## I. INTRODUCTION

In the case of digital terrain modelling and any other spatial modelling, accuracy, efficiency and economy are three main factors to be considered. Accuracy of DTM is a traditional topic in photogrammetric community since photogrammetrists are usually the DTM producer. Numerous papers in this area have been published in journals particularly in *Photogrammetric Record* and *ISPRS Journal of Photogrammtry and Remote Sensing* (formerly *Photogrammetria*) as well as the proceedings of ISPRS – *International Archive of Photogrammetry and Remote Sensing*. Both theoretical (Makorovic, 1972; Kubik and Botman, 1976; Frederiksen, 1981; Li, 1993) and experimental studies (Ackerman, 1979; Ley, 1986; Torlegard et al, 1986; Balce, 1987; Carter, 1989; Li, 1992; Li, 1994; Monckton, 1994) of DTM accuracy have been conducted by researchers.

If accuracy is the only concern, one would be better off to have data sets as dense as possible. For example, in the case of image matching, up to 50,000 to 700,000 points per stereo-model can be measured with the GPM-2 (Petrie, 1990). Such very dense data sets are not always suitable or appropriate for use by application specialists such as engineers, planners, etc. Indeed, in some cases, their sheer volume constitutes a definite deterrent or drawback to their use. If the efficiency of modelling process and the cost of computation are being considered, as they should be, then a lossy compression (or filtering) procedure needs to be applied to such data sets so that only a minimum number of data points will be selected while the specified accuracy of the final digital terrain model (DTM) will still be achieved. This general principle applies equally to both regular gridded data and irregular data but this article will deal with only the regularly gridded data.

In the case of regular grid sampling for DTM data acquisition using photogrammetric method, the usual practice is to employ a sampling interval which is suitable for area with roughest terrain in order to ensure that the accuracy of the final DTM will meet the requirements set by users. However, the problem is that, in flatter areas, too many data points will have been measured. One may argue that the progressive sampling suggested by Makarovic (1973) can be used to avoid such a problem. However, as pointed out by many specialists (see Toomey, 1984) that there are many fundamental problems with this techniques. At a later stage, following the line of thought opposite to progressive sampling, Makarovic (1977) himself proposed a method called "Regressive Rejection", which, as the name implies, rejects less important data points regressively. The working principle of this method is similar to the "quadtree" data structure. Another method which has been in use is the VIP (very important points) procedure (Chen and Guevara, 1987) in ARC/TIN. However, the effect of such procedures on the accuracy of resultant DTM has rarely been evaluated. Indeed, such an evaluation is very desirable in order to help users understand the potential consequence of using such a procedure.

This project aims to experimentally investigate into the relationship between the level (or degree) of data compression and the loss of accuracy for resultant DTM. In this study, a modified VIP procedure is used for experimental testing, with data sets of two testing areas.

Following this introduction is a section reviewing the VIP procedure and describing the modified VIP procedure. Then experimental testing is conducted to the effect of the VIP selection on the accuracy of resulting DTM using the modified VIP procedure. After that, an analysis of results is made and a mathematical model is established to predict the accuracy loss due to the compression.

## II. COMPRESSION METHODS USED: A MODIFIED VIP PROCEDURE

Before any experimental test can be reported, the procedure used for the testing should be described. As the procedure is modified from the VIP procedure, the VIP procedure should be described in some detail.

### A Review of the VIP Procedure

The VIP procedure was proposed by Chen and Guevara (1987). The basic principle is that "only those points whose significance is greater than the given threshold shall be selected" or "any points whose significance is smaller than a threshold shall be deleted". Then the most important element in this procedure is to assign a significance value for each point in the set.

The basic idea of obtaining significance values is illustrated in Figure 1. The technique used by Chen and Guevara (1987) is the so-called "spatial differential operator", which is a measure of the changing behaviour of a point from its neighbours. The working principle in one-dimension case is as follows:

Suppose, the height (H) of a point along a profile with equal interval is a function of its position (X) as follows:

$$H = f(X) \qquad (1)$$

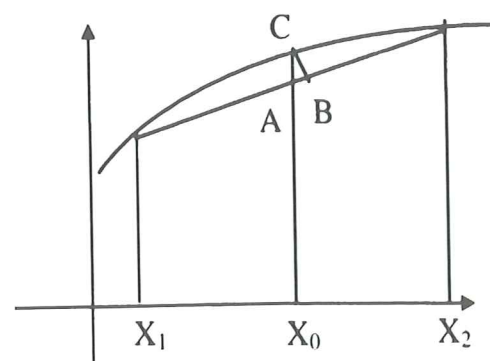then, its second order differntial value at point $X_0$ is



**Figure 1.** Significance value of Point C equals to the second differential of the point -- AC

$$\frac{d^2 H}{d^2 X} = f''(x_0) = 2 \times (f(x_0) - \frac{f(x_1) - f(x_2)}{2}) \qquad (2)$$

The distance AC in the Figure 1 is the second order differential value at point $X_0$  In the VIP procedure proposed by Chen and Guevara (1987), the value of BC was used instead of AC. They also consider four spatial directions (up-down, left-right, upper/left-lower/right, and lower/left to upper/right). For every point, the second differential values for all four directions are added together to represent the degree of the significance of this point. That is, the significance value for point $X_0$ is the sum of BC values in four directions or the average of these four sums.

Actually, the line of thought is very similar to that used by Makarovic (1984), who uses a Laplacian operator as follows:

$$L = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix} \qquad (3)$$

After each point is assigned a significance value, the points can be sorted out according to such a value. When the number of points, say N, required for a project is specified, the procedure will automatically select N points with largest significance values.

### A Modification to the VIP Procedure for Compression

As has been discussed in the previous sub-section, the threshold used for point selection in the VIP procedure is the number of points required. However, in practice, one can know only the required accuracy of the resulting DTM, instead of the required number of points. Therefore, the it is very desirable to find a threshold which can be directly related to the accuracy of resulting DTM.

To find out such a threshold for data selection, a close examination of Figure 1 must first be undertaken. It can be found that AC in Figure 1 is the error at $X = X_0$, if $X_0$ is removed and the profile is then linearly constructed from points at $X_1$ and $X_2$. In the procedure used for this test, the AC is selected as value of the so-called degree of significance for point at $X_0$. As a result, the accuracy loss due to such data compression can be directly related to threshold because the AC itself represents the DTM error resulted from the removal of these non-VIPs (not very important points).

## III. EXPERIMENTAL TESTING

The problem arising now is how much will the loss of

accuracy be in terms of standard deviation or RMSE if all the data points with the so-called degree of significance smaller than a specific value (the threshold) are removed. In other words, the relationship between the accuracy loss and the threshold needs to be investigated. In order to find such a relationship, some experimental tests have been carried out and the results will be reported in this section.

### Design of the Experiments

The steps of experimental tests are as follows: (a) The value of the so-called degree of significance of every point is computed; (b) Different threshold values are used for VIP selection to generate data sets with different degree of compression; and (c) The accuracy of the final DTM derived from each data set (after the removal of non-VIPs) is assessed using check points.

In the assessment of DTM accuracy, the standard procedure used by other researchers (Torlegard et al, 1986; Li, 1992) is followed. The steps are as follows: (a) to measure a set of ground points (check points) in the testing area with much higher accuracy than the resulting DTM; (b) to interpolate the height values from the resulting DTM for all points at the same locations (i.e. same X and Y coordinates) as the check points; (c) to use these check points as ground truth to compute the height differences between ground truth and interpolated values; and (d) to estimate the RMSE (root mean squre error) and SD (standard deviation) of these differences, which are two of the possible measures of DTM accuracy.

### Testing Data Sets

The test areas used in this study are two of the six (6) areas used for the ISPRS DTM test which was conducted by Working Group 3 of Commission III (Torlegard et al, 1986). A description of these two areas is given in Table 1 and a graphic view is given in Figure 2.

A set of square-gridded data for each area was meas-

ured on analytical photogrammetric instruments. The data sets were kindly made available to the authors through the courtesy of Prof. H. Ebner and Dr. W. Reinhardt of the Technical University of Munich (Germany). A description of these data sets is given in Table 2.

To make the testing results reliable, three paramters for check points need to be considered (Li, 1991), i.e. the accuracy, sample size (number) and distribution. In this test, the check points were measured on analytical photogrammetric instruments from aerial photographs with much larger scale than those used for the measurement of the source (raw) data for modelling so than the accuracy of check points is much higher than the accuracy of source data. Also a large number of points were measured to make the testing more reliable. The points are well distributed in the whole testing area. Detailed information about these check points is given in Table 3. These check points were used for the ISPRS testing and were kindly made available to the author through the courtesy of Prof. K. Torlegard and Dr. M. Li of the Swedish Royal Institute of Technology.

### Testing Results

For the first test area (Uppland), a set of threshold values is 0.4m, 0.5m, 0.6m, 0.7m, 0.8m and 0.9m, resulting six (6) new data sets. The numbers and percentages of points retained after the data compression are listed in the 2nd and 3rd columns of Table 4.

A triangulation-based program was used for the genaration of DTM and linear interpolation was used to interpolate heights from the generated DTM surface. The RMSE and SD values of the DTM errors at check points for different filtering level are also listed in Table 4. Diagrammatic presentations for RMSE, SD and the mean are also given in Figure 3.

For the second test area (Sohnstetten), the set of threshold values is 0.2m, 0.3m, 0.4m, 0.5m, 0.6m, and 0.7m, again, resulting six (6) new data sets. The num-

**Table 1.** Description of test area

| Test area | Description | Height Range | Typical slope |
|---|---|---|---|
| Uppland | Farmland and forest | 7 m - 53 m | 6° |
| Sohnstetten | Hills of moderate height | 242 m -538 m | 25° |

**Table 2.** Description of source data for testing

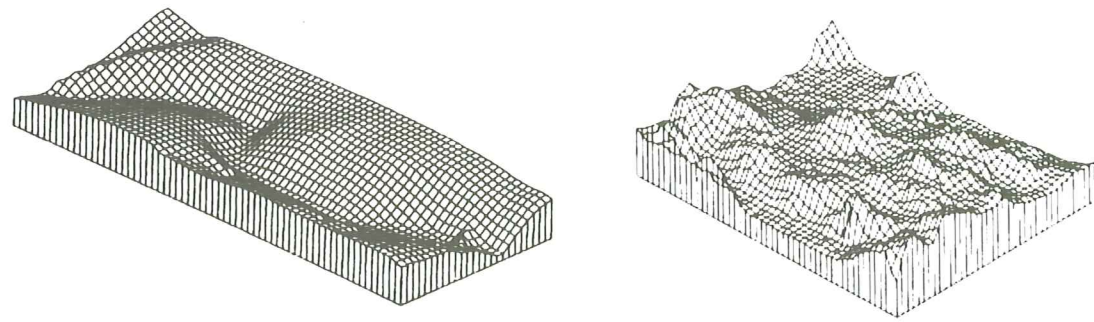| Test areas | Flying Height | Scale of photo | Grid interval | Height accuracy |
|---|---|---|---|---|
| Uppland | 4500m | 1:30 000 | 40m | ±0.67m |
| Sohnstetten | 1500m | 1:10 000 | 20m | ±0.16m |

**Figure 2.** Two test areas: Sohnstetten (left) and Uppland (right) (not scaled)

**Table 3.** Description of check points (ground truth)

| Test areas | Flying Height | Scale of photo | Number of pts | Height accuracy |
|---|---|---|---|---|
| Uppland | 900m | 1:6 000 | 2314 | ±0.09m |
| Sohnstetten | 750m | 1:5 000 | 1892 | ±0.054m |

**Table 4.** Variation of the accuracy of resulting DTM with threshold used for data selection for the Uppland area

| Threshold Value | (Number of Data Points | Percentage of Data Points | RMSE | Standard Deviation ($\sigma$) | Mean |
|---|---|---|---|---|---|
| 0.0 m | 1,862 | 100.0% | 0.770 m | 0.764 m | 0.102 m |
| 0.4 m | 1,489 | 80.0% | 0.778 m | 0.772 m | 0.098 m |
| 0.5 m | 1,392 | 74.8% | 0.779 m | 0.775 m | 0.086 m |
| 0.6 m | 1,301 | 69.9% | 0.783 m | 0.777 m | 0.084 m |
| 0.7 m | 1,206 | 64.8% | 0.801 m | 0.795 m | 0.102 m |
| 0.8 m | 1,123 | 60.3% | 0.810 m | 0.803 m | 0.104 m |
| 0.9 m | 1,044 | 56.1% | 0.825 m | 0.818 m | 0.108 m |



(a) Threshold and RMSE      (b) Percentate of pointsand RMSE

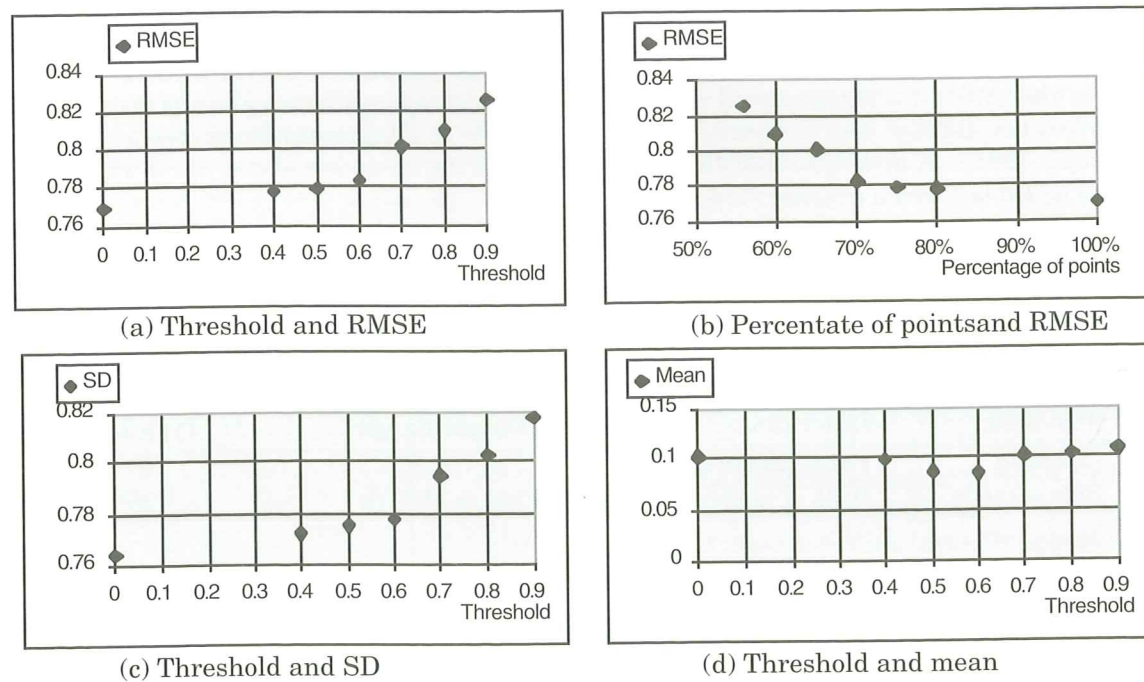(c) Threshold and SD      (d) Threshold and mean

**Figure 3.** Relationship between threshold and various statistical estimates for Uppland

bers and percentages of points retained after the data compression are listed in the 2nd and 3rd columns of Table 5. The RMSE and SD values of the DTM errors at check points for different compression level are also listed in Table 5. Diagramatic presentations for RMSE, SD and the mean are also given in Figure 4.

## IV. ANALYSIS OF TESTING RESULTS

The testing results have been reported in the previous section. An analysis of these results will be made in this section.

### Relationship between Threshold Value and Percentage of Points Retained

The percentage of points retained after data compression with different thresholds are listed in Tables 4 and 5. A diagrammatic presentation is given in Figure 5. From this diagram, it can be found surprisingly that the relationship between the percentage of points retained and the threshold value used for the filtering is quite linear.

### Relationship between Threshold Value and Resulting DTM Accuracy

From Figures 3 and 4, it can be found that the mean values for both testing areas are quite stable and there is no particular variation with the threshold values used for data compression. The results also demonstrate clearly that, with an increase in the threshold value, the accuracy of the final DTM becomes lower, i.e. the value of RMSE or standard deviation are greater. There is a dramatic increase in RMSE and SD values when the percentage of points retained becomes less than 70%. This is the case for both test-

**Table 5.** Variation of the accuracy of resulting DTM with threshold used for data selection for the Sohnstetten area

| Threshold Value | (Number of Data Points | Percentage of Data Points | RMSE | Standard Deviation ($\sigma$) | Mean |
|---|---|---|---|---|---|
| 0.0 m | 1,716 | 100.0% | 0.572 m | 0.561 m | -0.112 m |
| 0.2 m | 1,442 | 84.0% | 0.579 m | 0.567 m | -0.114 m |
| 0.3 m | 1,321 | 77.0% | 0.581 m | 0.571m | -0.112 m |
| 0.4 m | 1,220 | 71.1% | 0.580 m | 0.571 m | -0.102 m |
| 0.5 m | 1,110 | 64.7% | 0.586 m | 0.576 m | -0.106 m |
| 0.6 m | 1,017 | 59.3% | 0.596 m | 0.588 m | -0.095 m |
| 0.7 m | 942 | 54.9% | 0.613 m | 0.606 m | -0.087 m |



(a) Threshold and RMSE      (b) Percentate of pointsand RMSE

(c) Threshold and SD      (d) Threshold and mean

**Figure 4.** Relationship between threshold and various statistical estimates for Sohnstetten

(a) For Uppland area
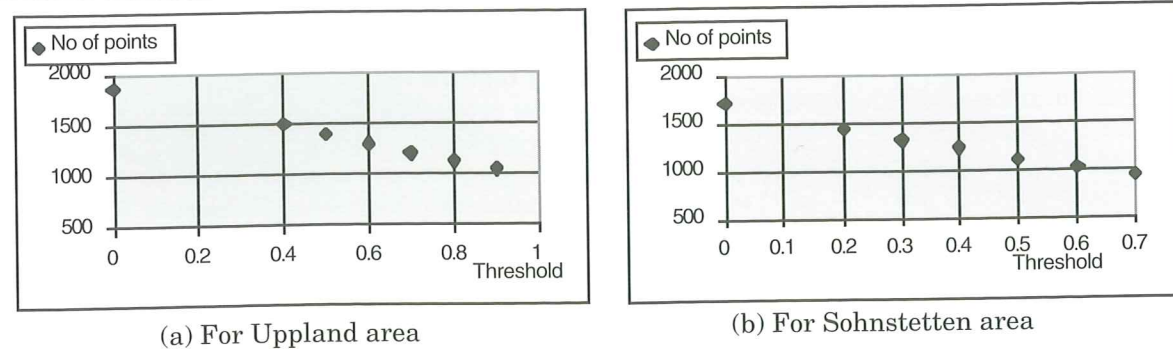


(b) For Sohnstetten area

**Figure 5.** Relationship between percentage of points retained and the value of threshold used for data compression

ing areas.

However, it is obvious that if there is a large degree of redundancy in the data set, then one may find that a 40% or 50% of reduction has little effect on the final result. Therefore, it might be dangerous to consider the 70% as a critical value.

**Relationship between threshold value and resulting DTM accuracy loss**

The accuracy estimates such as RMSE and SD are the results of errors at source data and errors introduced by the approximation of DTM surface to the real terrain surface. The mathematical expression for the standard deviations is as shown in Equation (4).

$$\sigma^2_{DTM} = \sigma^2_{Source} + \sigma^2_{Approx} \tag{4}$$

The expression for RMSE is similar to Eq.(4). The Accuracy of the resulting new DTM after data compression will be poorer than that of the old DTM. The difference is the accuracy loss due to data compression. The mathematical expression of such a loss can be given in Equation (5).

$$\sigma^2_{loss} = \sigma^2_{DTM/new} - \sigma^2_{DTM/old} \tag{5}$$

It is clear from Equations (4) and (5) that this accu-

racy loss is due to the poorer approximation of terrain surface by the new data sets after compression. In other words, the accuracy loss due to data compression is a function of the threshold used for data compression and it should take a form like Equation (6).

$$\sigma_{loss} = \frac{Threshold}{k} \tag{6}$$

If the distribution of these errors are known, then the k value can easily be set out. For example, if the errors belong to normal distribution, then the k value for Equation (6) is 3 according to error theory.

However, as has been discovered by researchers (Torlegard et al, 1986; Li, 1992), DTM errors don't obey the normal distribution. Therefore, the k value is not necessarily 3. However, in any case, with 96% confidence, the value of k will be smaller than 5, as predicted by the Chebyshev's Theorem: *the probability is at least as large as $1-1/k^2$ that an observation of a random variable X will be within the range $|\mu - k\sigma|$:*

$$P(|X - \mu| \le k\sigma) = 1 - \frac{1}{k^2} \tag{7}$$

For example, it predicts that the probability of DTM errors located within 2s, 3s, 4s and 5s are at least 75%, 89%, 94%, and 96%.

**Table 6.** Relationship between threshold and accuracy loss

| Threshold | Uppland Area | | | | Sohnstetten Area | | | |
|---|---|---|---|---|---|---|---|---|
| | $R_{loss}$ | $k_{RMSE}$ | $\sigma_{loss}$ | $k_\sigma$ | $R_{loss}$ | $k_{RMSE}$ | $\sigma_{loss}$ | $k_\sigma$ |
| 0.2 | — — | | — — | — — | 0.090 | 2.22 | 0.082 | 2.43 |
| 0.3 | — — | | — — | — — | 0.102 | 2.94 | 0.011 | 2.65 |
| 0.4 | 0.111 | 3.60 | 0.111 | 3.60 | 0.100 | 4.00 | 0.011 | 3.64 |
| 0.5 | 0.118 | 4.24 | 0.130 | 3.84 | 0.127 | 3.94 | 0.131 | 3.61 |
| 0.6 | 0.142 | 4.23 | 0.142 | 4.24 | 0.167 | 3.59 | 0.176 | 3.41 |
| 0.7 | 0.221 | 3.17 | 0.220 | 3.18 | 0.220 | 3.14 | 0.229 | 3.05 |
| 0.8 | 0.251 | 3.19 | 0.247 | 3.10 | — — | | — — | — — |
| 0.9 | 0.296 | 3.04 | 0.292 | 3.08 | — — | | — — | — — |

The testing results summaried in Table 6 reveals that the k value ranges from 2.43 to 4.24 with an average of 3.32. If the RMSE is considered, then the k value ranges from 2.22 to 4.24 with an average of 3.44.

## V. CONCLUDING REMARKS

In order to store and process data more efficiently, a compression procedure is normally applied for removing those points which are considered as being less important, i.e. non-VIPs. In this case, a loss in accuracy will result. The magnitude of such a loss will depend on the threshold used for VIPs selection, or non-VIP removal. The results from this investigation show that

(a) the number of data retained is quite linear with an increase in percentage of compression;

(b) a dramatic increase of accuracy loss may occur when the compression is over a critical value (30% for these two test areas);

(c) mathematical relationship between the threshold value and the resulting accuracy loss can be established if the former is carefully selected. For the modified VIP procedure described in this paper, and such a relationship is rewritten as follows:

$$\sigma_{loss} = \frac{Threshold}{k} \tag{8}$$

where k is a constant, ranging from 2.2 to 4.3 and the average value is about 3.4, according to the testing results obtained in this study.

These conclusions are only based on the results obtained from this test. These results provide DTM users some hints on the accuracy loss due to compression. However, it might be dangerous to generalise the results. Indeed, more investigation is required before a generalised conclusion can be made.

## ACKNOWLEDGEMENT

## REFERENCES

[1] Ackerman, F., 1979. The accuracy of digital terrain models. *Proceedings of 37th Photogrammetric Week*, University of Stuttgart. 113-143.

[2] Balce, A., 1987. Determination of optimum sampling interval in grid digital elevation models (DEM) data acquisition. *Photogrammetric Engineering and Remote Sensing*, 53(3): 323-330.

[3] Carter, J. R., 1989. Relative errors identified in USGS gridded DEMs. *Auto-Carto 9*, 255-265.

[4] Chen, Z., and Guevara, J., 1987. Systematic selection of very important points (VIP) from digital terrain model for constructing triangular irregular networks. *Auto Carto 8*, 50-56.

[5] Frederiksen, P., 1981. Terrain analysis and accuracy prediction by means of Fourier transformation. *Photogrammetria*, 36: 145-157.

[6] Kubik, K. and Botman, A., 1976. Interpolation accuracy for topographic and geological surfaces. *ITC Journal*, 1976-2: 236-274.

[7] Ley, R., 1986. Accuracy assessment of digital terrain models. *Auto-Carto London*, 1:455-464.

[8] Li, Z., 1991. Effects of check points on the reliability of DTM accuracy estimates obtained from experimental tests. *Photogrammetric Engineering and Remote Sensing*, 57(10): 1333-1340.

[9] Li, Z., 1992. Variation of the the accuracy of digital terrain models with sampling interval. *Photogrammetric Record, 14(79): 113-128*.

[10] Li, Z., 1993. Mathematical models of the accuracy of digital terrain model surfaces linearly constructed from gridded data. *Photogrammetric Record*, 14(82): 661-674.

[11] Li, Z., 1994. A comparative study of the accuracy of digital terrain models based on various data models. *ISPRS Journal of Photogrammetry and Remote Sensing*, 49(1): 2-11.

[12] Makarovic, B., 1972. Information transfer in construction of data from sampled points. *Photogrammetria*, 28(4): 111-130.

[13] Makarovic, B., 1977. Regressive rejection - A digital data compression technique. *Proceedings of ASP/ACSM Fall Convention*. Little Rock, October 1977.

[14] Makarovic, B., 1983. A test on compression of digital terrain model data. *ITC Journal*, 1983-2: 133-138.

[15] Monckton, 1994. An investigation into the spatial structure of error in digital elevation data. In: Worboys, M. (ed), 1994. *Innovations in GIS 1*, Taylor & Francis. 201-214.

[16] Petrie, G., 1990. Analogue, analytical and digital photogrammetric systems applied to aerial mapping. In: Kennie, T., and Petrie, G. (eds), *Engineering Surveying Technology*,. Whittles Publishing, 238-288.

[17] Toomey, M. (ed), 1984. *Digital Elevation Model Workshop Proceedings*. Edmonton, Alberta. 231pp.

[18] Torlegard, K., Ostman, A. and Lindgren, R., 1986. A comparative test of photogrammetrically sampled digital elevation models. *Photogrammetria*, 41(1): 1-16.