

## **Digital Environment Data: Identifying Anomalies from Source to Final Databases**

**Robert F. Richbourg, Timothy M. Stone, and George E. Lukes**

**Institute for Defense Analyses**

**Alexandria, VA**

**rrichbou@ida.org, tstone@ida.org, glukes@ida.org**

### **ABSTRACT**

Digital representations of the environment are being used in a wider spectrum of applications. Military units deploy to operational areas with a digital database of that area. This is one example of the many factors that have led the National Geospatial-Intelligence Agency (NGA, formerly NIMA) to adopt a new data production strategy focused on providing digital geospatial data in addition to traditional paper map and chart products. The NGA emphasis on digital geospatial data promises new opportunities for simulation database developers and users. Simulation systems already utilize large-scale, high-fidelity, geo-specific databases to execute joint experiments including urban area operations. Expanded availability of digital source data can only increase this trend; however, the new geospatial data production model, coupled with more rigorous demands and expectations from the operational community, results in continued tension between data quantity and data quality within a crisis-response production environment. Processes developed and refined to produce traditional maps and charts are not sufficient to meet the demands for multi-purpose digital geospatial data. This paper reports on results of research into identifying the data anomalies that may arise in such an environment and describes the development of automated tools that can be applied early in the production process to detect those anomalies.

### **ABOUT THE AUTHORS**

**Robert F. Richbourg** is a member of the Research Staff in the Simulation Center at the Institute for Defense Analyses. He is a retired Army officer who earned his Ph.D. in computer science in 1987. In his last active duty assignment, he was an Academy Professor and Director of the Artificial Intelligence Center at the United States Military Academy, West Point. He has been working in the area of simulation environments for 10 years under sponsorship of DARPA, DMSO, and STRICOM.

**Timothy M. Stone** is a Research Programmer/Analyst in the Simulation Center at the Institute for Defense Analyses. His research interests include: Computational Geometry, Triangulated Irregular Networks (TINs), Simulation, Algorithm Design, Line-of-Sight (LOS), and Computer Graphics. He obtained his Master's degree in Computer Science from New Mexico State University in 1996. He is a co-developer of the SEE-IT application.

**George E. Lukes** is a member of the Research Staff in the Simulation Center at the Institute for Defense Analyses. From 1994 to 2000, he served as a Program Manager at the Defense Advanced Research Projects Agency where his responsibilities included the Synthetic Environment Program for the Synthetic Theater of War. Previously, he led research and development efforts at the U.S. Army Topographic Engineering Center in terrain database generation for advanced distributed simulation (e.g., SIMNET, Project ODIN, Battle of 73 Easting, STOW-Europe, Bosnia).

## **Digital Environment Data: Identifying Anomalies from Source to Final Databases**

**Robert F. Richbourg, Timothy M. Stone, and George E. Lukes**  
**Institute for Defense Analyses**  
**Alexandria, VA**  
**rrichbou@ida.org, tstone@ida.org, glukes@ida.org**

### **INTRODUCTION**

Have you ever used a virtual, platform-level simulator with computer generated forces and seen a tank stand on end, suddenly “teleport” to a new location, float or submerge above or below the terrain, refuse to move, or simply incorporate a few steps from the tango into its cross-country movement? Unusual entity behaviors occur and are not always rare events. Root causes could stem from errors in behavior control software, in the construction of the environment (terrain) over which the entities travel, or in some unexpected interaction between control software and the environment. This paper focuses on the synthetic environment as a source of error in results derived from use of that environment.

The fact that environmental constructions or representations can lead to abnormal results is not a new observation. Analysts and designers have often devoted substantial pre-exercise effort moving simulation entities back and forth across the terrain to identify areas where anomalous behaviors might occur (“crawling the database”). The goal for this kind of analysis is to identify areas within the environment that need to be repaired or to restrict simulation scenarios to portions of the environment database where simulation entities will operate as expected. At best, this practice offers only a partial solution to ensuring correct entity behaviors.

First, it is very doubtful that every possible entity will interact with each terrain data object. Thus, this approach only provides an approximation to identifying areas where “everything will go smoothly.” That is, an observation that entity E1 interacts correctly with terrain data object O1 does not guarantee that E1 will also interact correctly with the topological neighbor of O1. Nor does it guarantee that other entities, E2, will interact correctly with O1 when E1 and E2 use different behavior control software.

A second problem with this approach relates to time requirements. The time required for an entity to traverse a database, as illustrated in Figure 1, is a function of spacing between adjacent traversal passes, vehicle speed and database size. Digital environment databases are continually expanding in feature content and expanse. Thus, time requirements alone dictate that this kind of database inspection becomes less and less of a useful approximation of quality.

A third problem area stems from the emerging use of digital environmental data. Increasingly, actual operations are supported by reference to digital representations of the physical environment. As an example, a central concept for the Army’s Future Combat System (FCS) is that there will be a single environmental representation and that it will support all aspects of FCS needs including training, planning, analyses and operations. Clearly, the database must support error-free operations over the full extent.

A more viable course of action requires determining the source of error, making repairs, and then focusing on the desired result of using the terrain database to determine the areas that will be utilized. Further, the database analysis needs to be automated; the amount of environmental data and the inherent error-prone nature of human inspection make any other approach infeasible. The Defense Advanced Research Projects Agency (DARPA) began work on such a capability in 1997 and the Defense Modeling and Simulation Office (DMSO) has sponsored continued research and development. DMSO has made the Synthetic Environment Evaluation - Inspection Tool (SEE-IT) (Richbourg and Stone, 1998) a freely available resource for the simulation community, fostering the construction of “correct” digital representations of the environment. (See <http://tools.sedris.org>.) SEE-IT has been used in support of a variety of simulation efforts (Praeger et al, 2002).

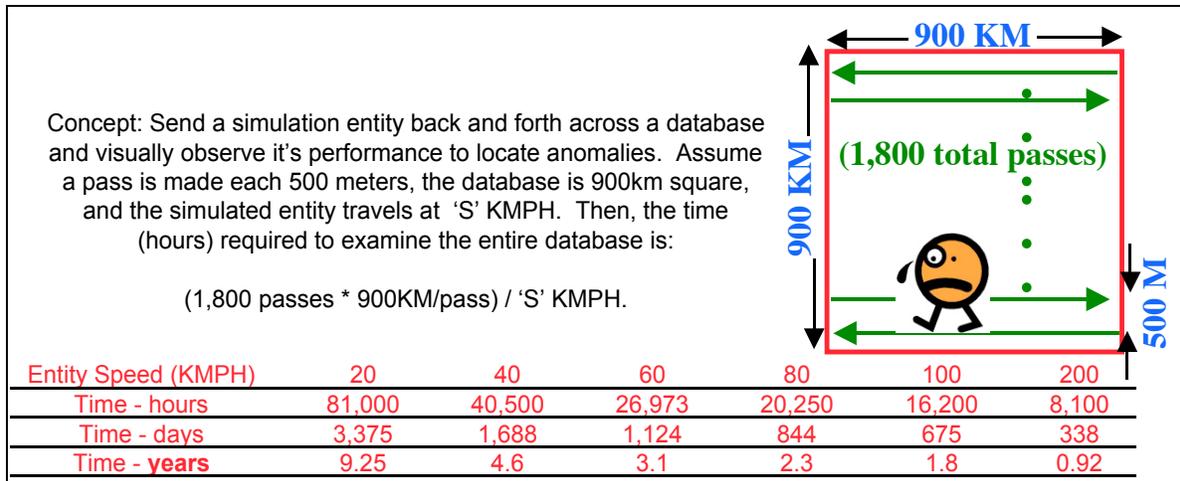


Figure 1. Time Requirements for Crawling a Database

**PROCESS**

Development of SEE-IT analytical capabilities has largely been based on research to identify environmental conditions that can lead to aberrant simulation entity behaviors. An example of the desired process occurred early in SEE-IT development efforts. Analysts at the US Army's Topographic Engineering Center (TEC) hypothesized that ModSAF entities frequently had difficulty traversing any terrain surface with a specific type of topology (T-Vertices). Subsequently, SEE-IT analysis identified all such topological constructions in the simulation environment for Camp Pendleton and provided the results to TEC. Then TEC analysts sent simulation entities across the terrain at each of the T-Vertex locations. Each time, an abnormal behavior resulted. In fact, the image in Figure 2 was captured during this experiment. As a result of the knowledge gained during this experiment, T-Vertex constructions no longer pose an obstacle for terrain movement.

This exemplifies a type of result that needs to be extended. Extensions are possible in both the numbers of conditions that require identification and in the types of data that can be analyzed. In particular, analyses need to be applied to source data (e.g., data from which simulation environments are constructed). Early identification of undesirable conditions in source data will decrease costs and required timelines for simulation database construction. Analyses at the source data level will also support applications that use that data directly (e.g., military operational or command and control systems), a trend that will increase along with increased levels of automation in military systems.

Source data for military simulation environmental database construction is typically provided by the National Geospatial – Intelligence Agency (NGA), formerly the National Imagery and Mapping Agency (NIMA). Products representing terrain, bathymetry and feature data (e.g., transportation, hydrography, cultural, vegetation, and soils) have been produced in a variety of formats. Various sources of errors are possible, particularly as products are applied to more demanding applications than they were initially designed to support.

In general, the process of devising analyses applicable to geospatial source data is quite similar to that used to detect problems with simulation databases. It requires forming a hypothesis, conducting experiments to judge hypothesis validity, and then implementing analysis and identification capabilities when the hypothesis proves valid. This is, of course, a very standard approach. Differences here arise in the hypothesis formulation stage and in the effort to implement analyses

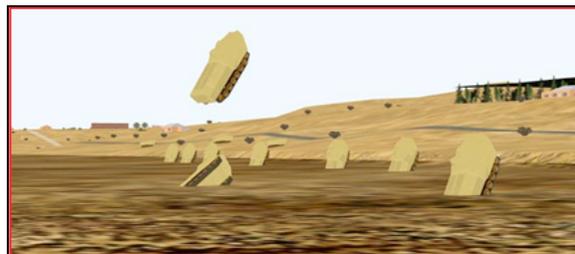


Figure 2. Simulation Landing Craft React to Terrain Topology at Locations Identified by SEE-IT

based on valid hypotheses. Regarding the former, we need to consider the methods used to produce the NGA data as an approach to forming hypotheses on “what could go wrong.” With regard to the latter, automating the detection of errors in the data involves solving computational geometry problems (Preparata and Shamos, 1985) over a very large data space, so efficiency will be a paramount concern. As an example, a recent analysis effort started with the delivery of a 15 GB data set covering a single country. This paper focuses on the problem of cataloging what could go wrong.

## DATA SOURCES AND PRODUCTION METHODS

Traditional cartographic practice produces topographic maps where the terrain surface is depicted by contour lines of equal elevation while natural and man-made features are represented by a variety of symbols, colors and patterns. Early work in digital cartography focused on production of paper maps and charts requiring development of specialized hardware<sup>1</sup> and software. Later efforts focused on scanning finished cartographic products to field the electronic “images” of maps and charts<sup>2</sup> that proliferate today.

The Digital Radar Landmass Simulation (DRLMS) program initiated in the late 1960’s was a landmark event for both the military mapping community and the training and simulation community. DRLMS was a special purpose digital mapping program conceived to provide data for real-time flight simulators, initially the navigation radar display for B-52 aircrew training, and was the origin of the now well-known Digital Feature Analysis Data (DFAD) and Digital Terrain Elevation Data (DTED) products.

Much has changed in more than forty years of experimentation and practice. While the fundamental cartographic issues of abstraction, generalization, metric accuracy and timeliness endure, the characteristics of generated data, including sources of error, continue to evolve based on available source materials and production methods.

### Feature Data

DFAD is a classic *vector* format product in which cartographic objects are represented by attributed *point*,

*linear* and *areal* features. In the 1980’s and the 1990’s, the Defense Mapping Agency (DMA, a precursor organization to NGA) developed additional vector products<sup>3</sup> to support terrain analysis and a series of general-purpose vector map products.<sup>4</sup> For the most part, vector feature data has been and continues to be captured interactively by human operators. The major changes have been in available data sources and the evolution of interactive digitizing workstations.

Initial efforts involved operators capturing cartographic features from two-dimensional (2D) manuscripts, orthophotographs and/or finished maps. In time, far more sophisticated workstations<sup>5</sup> were developed to extract feature data directly from stereoscopic aerial photography (film-based or “hardcopy” systems). Today, PC-based “softcopy” workstations<sup>6</sup> are in widespread use operating in 2D from digital raster maps or orthoimages<sup>7</sup> or in 3D from stereo imagery. What have not changed are reliance on human operators and the potential for errors to be introduced as part of the interactive digitizing process.

Today, NGA is migrating from a strategy based on production of standard map and digital geospatial products to implementation of a multi-purpose Geographic Information System (GIS) that supports direct query by operational users as well as production of standard products. To field a Geospatial-Intelligence Feature Database, NGA Centers, contractors and coproducers are working to integrate previously collected geospatial data (e.g., VMAP, VITD) in a very large GIS. The concept calls for operators to use current imagery with the existing geospatial data to (a) resolve conflicts, (b) perform metric refinement, (c) intensify feature density, and (d) update the database over time. This ambitious undertaking presents both great promise for greater responsiveness as well as new opportunities for incorporating errors into the geospatial data repository.

---

<sup>1</sup> Large-format scanners, interactive digitizing tables, graphic displays and large-format plotters.

<sup>2</sup> For example, NGA Compressed ARC Digitized Raster Graphics (CADRG).

---

<sup>3</sup> Interim Terrain Data (ITD) succeeded by Vector Interim Terrain Data (VITD), each consisting of separate files representing transportation, vegetation, drainage, soils, obstacles, slope and bridges.

<sup>4</sup> VMAP products content similar to traditional paper maps and charts (e.g., ONC, JOG, TLM).

<sup>5</sup> The CAPIR prototype was developed at the US Army Engineer Topographic Laboratories in the late-1970s. Later, the Feature Extraction (FE) Workstation was fielded by the Defense Mapping Agency as part of the Mark 85 Modernization Program.

<sup>6</sup> Popular products include ArcGIS (ESRI), GeoMedia Pro (Intergraph), and Microstation (Bentley).

<sup>7</sup> CADRG and NGA Controlled Image Base (CIB).

## Elevation Data

The NGA DTED product is a classic example of a Digital Elevation Model (DEM) where the terrain surface is represented by discrete elevation values for each node of a grid. A variety of methods have been developed to produce DEMs – in all cases, quality of the end product varies depending on source material and the production process.

Traditional map compilation relied on interactive extraction of contour lines from stereoscopic aerial photography by human operators using optical-mechanical stereo plotters. Early efforts in digital cartography involved tracing and later scanning contour line manuscripts to yield contour line vectors, followed by vector-to-raster conversion processes to yield a DEM. This approach, now implemented with improved scanners and more sophisticated vector-to-raster conversion algorithms, remains in widespread use -- when applied to high-quality contour sheets, it continues to produce some of the highest quality DEMs.

Pioneering work in automatically extracting digital elevation data from stereo imagery started in the 1960's (Bertram, S., 1969) and continues today in a variety of commercial computer-assisted mapping systems.<sup>8</sup> Human operators are tasked with initializing and supervising the elevation extraction process that can fail in areas of low spatial frequency content, water bodies, dark shadows and repetitive terrestrial patterns.

Radar mapping systems have the attractive property of operating in most weather conditions. Novel techniques for extracting terrain elevation data with Interferometric Synthetic Aperture Radar (IFSAR) were pioneered in the 1990's with aircraft mapping systems,<sup>9</sup> and then applied to space radar systems. The Shuttle Radar Topographic Mission (SRTM) flown in February, 2000 acquired data to produce DTED Level 2 products covering over 80% the earth's landmass (60 degrees north latitude to 54 degrees south latitude). An anomaly often present in this data is repeated void areas introduced by radar shadows and water bodies. Extensive post-processing is underway to bring the data into correspondence with national map accuracy standards.

The newest and most rapidly growing source of digital elevation data is derived from airborne Light Detection

And Ranging (LIDAR) systems. LIDAR systems use a laser scanner to illuminate the terrain surface and sophisticated on-board signal processing to compute range to the terrestrial surface based on the time difference between transmitting a coherent laser signal and receiving the reflected return signal. Although constrained by modest operating altitudes and scan swath widths, contemporary LIDAR systems are now producing the highest resolution data generally available over urban and floodplain areas.

## FEATURE DATA ANOMALIES

Feature data anomalies that we have identified to date come from one of two main processes. First, conflation is the process by which two or more digital data sets describing the same region are combined into a single database of that region. As an example, the data producer might have access to digital data used to create two different scale maps of the same region. Typically, these kinds of data sources will intersect each other as measured by area of coverage or feature content, but one data set will not be a proper subset of the other in either measure. Thus, the possibility of feature duplication is high. Duplication can occur in many different ways, as explained below.

The second process that contributes to the inclusion of error in digital data sets is manual data capture, also known as digitization. In this process, a human operator uses a digitizing tablet (or other device) to capture digital feature data from a map or image source. Manual data capture frequently produces errors that are too small to be detected by the human eye. These are inconsequential when the final product is to be an analog product produced from digital data (e.g., a printed map or chart). As an example, manual data capture could introduce small gaps between road network segments that would not be discernable to a human reading a map, but would create network breaks when using a computer for route planning. As is the case with conflation, manual data capture leads to several types of specific errors, as explained in the data capture section below.

### Conflation Errors

The potential conflation error that seems most obvious is the possibility of introducing duplicate features. Here, two or more of the data sources being merged would contain an identical feature (as measured by feature attribute set and values) at the same location. The merged data would contain both features, a condition that is near impossible to detect by visual inspection. Imagine an electronic map display that contains

---

<sup>8</sup> Contemporary products include ImageStation (Intergraph), SOcET SET (BAE), SoftPlotter (Boeing IIS) and Stereo Analyst (Leica)

<sup>9</sup> Intermap Technologies and EarthData currently operate commercial IFSAR mapping systems.



**Figure 3.** Intertwined Roads

“duplicate” point features. Visual examination of the map would never reveal the fact that a single point feature had been drawn two or more times. Human inspection of the data is also unlikely to detect duplication errors simply due to data volume. In a recent examination of a country data set, automated analysis detected 41 duplicate features within a population of 976,580 point, areal, and linear features. This is a very small percentage of the data and their identification is tantamount to “finding needles in a haystack;” however, they are not so insignificant as to be ignored. Duplicate features are often clustered in small areas that could impact performance in simulation systems that rely on load modules during run-time processing.

Analysis of the same data set also detected “near duplicates”, features that had exactly the same location, but different attribution. Analysis of the same 976,580 features detected a total of 442 “geometric duplicates” of this type. In this case, these duplicates were clustered together, and could impact on run-time performance.

Both of these conditions have readily apparent repair strategies. Complete duplicate features can simply be deleted. Geometric duplicates can also be repaired by deletion, as long as some mechanism exists to ensure the correctly attributed feature remains after repairs. In general, the requirement for access to corroborating information applies to almost all repair endeavors

(given the exception of complete duplicates). The underlying assumption is that the data is “geo-specific;” it accurately describes a specific location in the physical world. In the special case of “geo-typical” or fictitious) terrain, repair decisions can be made more readily.

Another error stemming from conflation of similar data results in a merged data set that has very similar features placed very closely together. The old saying that “no two database descriptions of the same object ever agree” plays a central role here. As an example, suppose that two data sets describe the same road linear feature, but provide slightly different attribution and geometry. The different attribution could prevent detection of the duplication prior to conflation. After conflation, the result could include multiple roads that seem intertwined. Figure 3 illustrates this occurrence. The zoom insert (right side) shows that there are two roads present. The distance between these two, at the widest point, is only 0.8 meters. Again, the human perception system is a poor detector of this condition. Repair of this type of condition requires informed deletion.

Here again, we have an example of a condition that is insignificant when an analog product is the desired end state. The human eye will gloss over such inaccuracies and the duplicate roads would probably not be noticed; however, an automated system that, for example, calculates miles of roadway in specific areas, cannot ignore data such as this, and will produce an incorrect result.

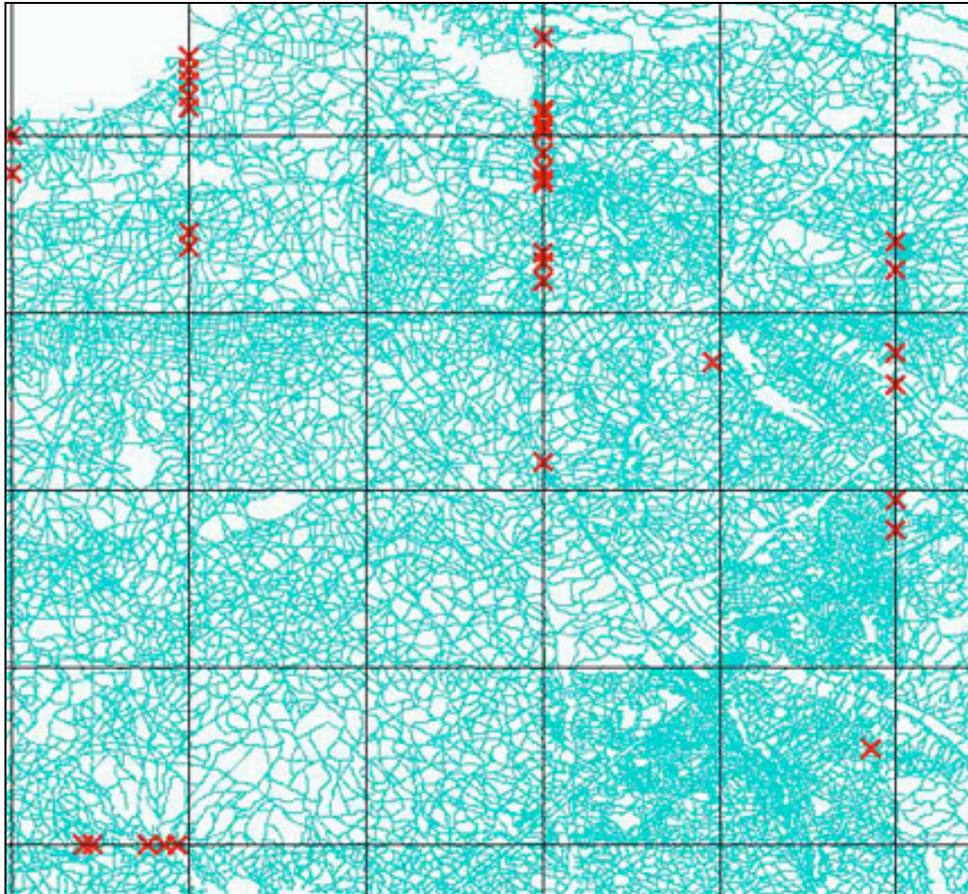
Point and areal features are subject to this same conflation problem. The areal features typically intersect each other while the point features are placed in very close proximity. Automated methods can be used to detect these relationships between feature geometries; however, it is usually impossible to determine the best repair strategy without consulting another source of information for corroboration.

The conflation process can also introduce errors when data describing different layers (e.g., hydrology, transportation, vegetation, etc.) are integrated into a single database. As an example, analyses have detected

buildings from the cultural layer positioned inside perennial lakes from the hydrology layer, or intersected by roads from the transportation layer. These conditions arise because each layer of thematic data is created as a stand-alone product. Combining layers of data makes the positional conflicts evident. Again, proper correction requires some information source that can determine which layer provides a more accurate portrayal of the physical world.

These kinds of errors are significant. Planning an assault on a building completely surrounded by water is quite different from planning an attack against a building along a shoreline. Similarly, depicting a road passing through a building implies a tunnel or other passage that could be operationally significant, but that may not exist. Humans can reference maps, overhead imagery, or other descriptions of the area to make appropriate decisions in these cases; however, automated reasoning systems (including simulations) must depend on the digital data provided to them, and they can be misled by data containing these kinds of errors.

A related error condition could be classified as either conflation related or data capture related. This condition results in discontinuities of features as they span whole degree (geodetic coordinate system) lines. Typically, the data creation process features a spatially derived division of effort. That is, regional data sets are often created as collections of 1-degree cells. When these cells are juxtaposed to form the data set for the entire country, small errors at the cell boundaries become evident as feature discontinuities. As an example, a linear road that runs across degree lines would be digitized in parts (within each one degree cell). Frequently, such features include disconnects or gaps as they cross over a degree boundary, a reflection of the data being captured in different projects or by different operators within their own (spatial) areas of responsibility. Figure 4 provides an example of this phenomenon by illustrating road network breaks (as red ‘X’ marks on top of blue road networks) and their positional relationship to degree boundaries (black grid lines).



**Figure 4.** Road Network Breaks and Degree Boundaries

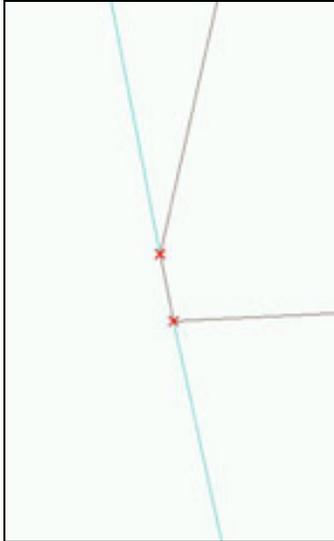
### Data Capture Errors

Data capture errors are often the result of capture device positional inaccuracies as features are digitized. Again, these types of conditions are not evident when a printed map is the final product; however, direct automated use of the digital data files makes these errors important. A prime example is the road network break as discussed above. Such network breaks will defeat shortest-path and other graph search algorithms when applied to road network routing problems. Network breaks of this type can be described as either “overshoots” or “undershoots.” Overshoots describes lines that extend beyond their intended end point. Undershoots describes lines that fall short of their intended end point. Typical repair strategies for network breaks will “snap” one disconnected vertex to the other. Selection of the best “snap-to” point requires a corroborating information source.

An interesting subclass of the “undershoot” case occurs when the gap between the linear features is actually bridged by another type of feature. Figure 5 illustrates

a potential instance of this phenomenon. Here, roads are depicted in blue and cart tracks are drawn in brown. The red ‘X’ marks denote the end points of the road linears, which fall short of meeting each other by 1.5 meters; however, the cart track feature bridges this gap. It is easy to understand how such an error might occur during the digitization process.

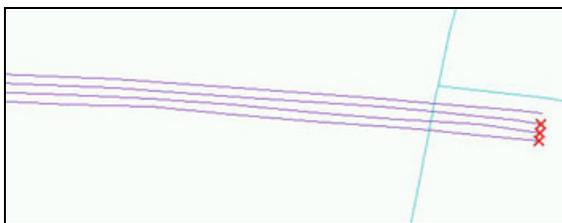
Detection of graph breaks requires solving “closest point” problems. That is, for each linear feature end node, find the closest end node of a second linear feature or node of an areal feature. Whenever this distance is greater than zero, a potential graph break exists. While this can be an expensive calculation if not managed properly, the real difficulty lies in separating the disconnected end nodes that accurately reflect the physical world from those that are the result of error. A heuristic that minimizes false positives while not discarding the most difficult to find errors is based on distance; if the closest node to the disconnected end node is more than 5 meters away, then we treat the situation as correct. Of course, this is only a heuristic



**Figure 5.** Road Gap Bridged by a Cart Track

and can provide incorrect results. Figure 6 provides an example of a typical false positive.

The 5-meter tolerance is an analysis system default value. It is possible to dynamically alter this value (and all other tolerances we discuss in the following) because it has become very clear that a single value will not be appropriate for all situations in which an analysis might be executed. In fact, specification of the “best” value for each situation is a difficult problem in and of itself. There is always a tension between avoiding false positives and including all potential errors. Our approach has been to use previous experience to provide a reasonable default value that minimizes false positives but does not miss the errors that are difficult to locate. The default value can be changed at any time to be more or less inclusive. Further, whenever possible, we sort the error notifications according to their deviation from the tolerance. This permits a prioritized look at the potential error conditions; the “worst” conditions are always presented first.

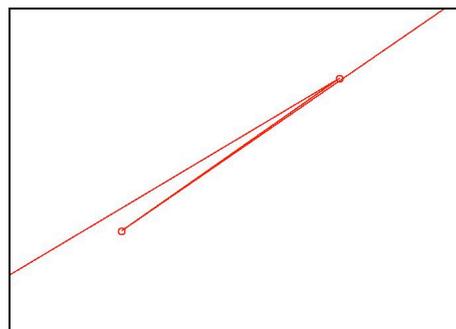


**Figure 6.** Telephone Line Features (purple) and a Series of False Positive Network Breaks (red ‘x’)

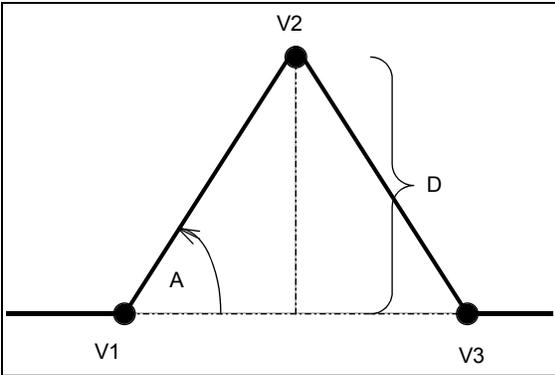
Another type of data capture error condition impacts both linear features and areal feature perimeters. Here, a point has been inadvertently included in the feature, so that a segment of the areal or linear appears to turn back on itself or preceding feature segments. This class of error has been described as a “kink,” “kick-back,” or “Z.” Figure 7 provides an actual example, depicting the geometry of a portion of a linear road feature. A simple correction strategy is to remove one or more of the vertices that produce the “kink.” Again, detection and repair is desirable to ensure correct results from, as a single example, graph search algorithms.

There is no widely accepted definition of a kink. One definition used by a well-known Geographic Information System (GIS) states that a vertex is considered to be a kink when the distance from that vertex to an imaginary line drawn between the immediately preceding and following vertices is greater than a user-specified tolerance. Figure 8 illustrates this definition where vertex V1 immediately precedes the kink vertex (V2), vertex V3 follows it, and D is the distance to be compared to the specified tolerance.

There are several problems with this definition and detection criterion. Most obviously, the tolerance is “an exercise left to the reader.” Also, overly tight tolerances will admit many false positives while overly large tolerances will fail to detect small kinks. The approach we have taken relies on comparing angular measures to a threshold of 165 degrees (as in angle A, Figure 8). The 165-degree threshold is a heuristic measure that has been used to ensure error conditions are not overlooked, at the cost of admitting a potentially high number of false positives; however, unlike the distance-based strategy, it is not sensitive to scale. Also, the approach of arranging the analysis results as a sorted list helps with the management of false positive notifications.



**Figure 7.** Road Network “Kink”



**Figure 8.** Measuring a “Kink”

The threshold or tolerance value used in applying the angular measure strategy is, like most others, context-sensitive. Similar to the example provided with Figure 6, this analysis type provides clear examples of the relationship between feature type and tolerance applied. A recent analysis of a country data set detected 1,697 potential kinks involving 23 types of linear features (e.g., roads, rivers, etc.) within a population of 435,615 linear features describing 47 feature types. The number of false positives varied greatly by feature type. Linear walls included 0 false positives, while 90% of the detected potential kinks for railroad sidings were false positives. This result is intuitive as adjacent railroad sidings are usually close to parallel before they intersect while walls normally form convex enclosures. Thus, there is a requirement for more study to determine an appropriate threshold for each type of anomaly for each feature type.

### ELEVATION DATA ANOMALIES

We have been working with four types of elevation data anomalies to date. These include data dropouts, spikes and pits, feature omissions, and errors that arise during conflation. Data dropouts are easily detected as they are recognized during data production and are required to be represented by a sentinel value as the assigned elevation.

Spikes and pits can be produced during automated or semi-automated data capture. If noise exists in the source imagery, the data production system often produces outlier elevation values as a response. These can be either falsely high (spikes) or low (pits). An intuitively appealing approach to locating these errors is to simply look for all elevations either above or below a specified threshold. This approach can work reasonably well within small areas. However, when applied to larger areas where the range of elevation values is

greater, simple thresholding by elevation only produces the worst cases of spikes or pits. Our approach to locating these anomalies attempts to take into account the elevations within localized areas. We calculate the slope from each elevation point to each of its adjacent neighbors. When over half of the slopes are greater than 75 degrees, a spike or pit notification results. (Again, 75 degrees is the system default value that can be changed before any analysis.)

Data omission errors can arise from conditions in the natural environment. As an example, cloud cover could obscure terrain features in some portions of stereo imagery. Terrain that is not represented in the source imagery cannot be compiled into the resulting DEM. Our approach to isolating such anomalies is based on comparing multiple DEM that describe the same region. Non-systematic differences between the models point to potential errors of omission. Even lower resolution data generally contains sufficient spatial frequency content to identify significant inconsistencies. As an example, we have found a case where a small hill was represented in DTED Level 2 (nominally 30 meters on-ground distance between elevation posts) but was not present in higher-resolution DTED Level 3 (nominally 10 meters between elevation posts). In this particular case, cloud cover in the source imagery seems to have prevented capture of the hill feature during creation of the higher-resolution data. As it happened, this small hill was the only significant relief in an otherwise flat desert area, so it was tactically commanding terrain. This particular case is also a rare example where context provides circumstantial evidence that the hill really exists, without need for other corroborating evidence. Other feature data included a road that spiraled around from the base of the hill, winding its way to the top of the hill, finally arriving at a building on the hilltop. This arrangement of a building and road leading to it makes sense when combined with the DTED Level 2 (that includes the hill) but seems out of place when combined with the DTED Level 3, where the hill feature is absent.

DEM conflation errors can appear in ways similar to those that impact feature data. This includes cases where different elevation models are juxtaposed to describe a larger region and when feature data are combined with a DEM. The former case frequently introduces artificial “ridge lines” that coincide with individual DEM boundaries. Such instances can be automatically detected by isolating areas of steep slope (abrupt elevation change) that have a very linear presentation.

Combining feature and elevation data frequently exposes disagreement between the two layers of informa-

tion. As an example, river features might appear to flow both up and down hill, when the DEM elevations are combined with the river feature (x,y) locations. These are easily identified using slope direction calculations.

Each of these types of elevation-related errors can have implications in operational use. Again, the greatest impact will be on automated systems. Elevation values directly impact line-of-sight and mobility predictions, two of the three pillars in “move-shoot-communicate.” Elevation model errors can also be repaired if appropriate corroborating information is available.

### SUMMARY AND CONCLUSIONS

The shift in emphasis from analog geospatial products (e.g., printed maps and charts) to generation and dissemination of digital geospatial data by NGA has many ramifications for both the geospatial data producer and geospatial data consumer communities. Consumers should expect the availability of data much more appropriate for use by automated systems. The data producer is implementing a new large-scale production strategy that requires an accompanying development of new quality assurance processes. Data that meet the standards for analog product production can include many error conditions that will only become evident with automated use of that data. New standards are required that can address graph connectivity, feature duplication, combination of feature data from different layers, and positional accuracy. Moreover, these standards must be crafted so that they can provide specifications suitable for use by automated quality control utilities. Continued reliance on the human perception systems simply will not suffice given the magnitude of the problem.

In most cases, repairing the conditions described here requires access to primary source materials, typically stereoscopic imagery and/or derived orthoimagery. Such supporting information is necessary to determine the most applicable repair strategy. This information is rarely available in a form that allows complete automation of the repair process; human judgment and interpretation are vital.

Standards for automated quality control systems will have to be both precise and context sensitive. As we have seen, tolerances that can help identify errors involving one type of feature will not be suitable for use with all other features. The results described above are the initial thrusts into this line of exploration. Until such research is completed, we have taken an approach that provides a reasonable default value, allows dynamic specification of thresholds, and presents analyses results in priority order whenever possible.

Finally, we cannot blindly trust the results provided by our automated systems that consume digital geospatial data until we have some mechanism to justify confidence in the data itself. The most accurate algorithms will produce incorrect results when applied to incorrect data.

### REFERENCES

- Bertram, S., (1969). The UNAMACE and the automatic photomapper. *Photogrammetric Engineering*, 35, 576 – 596.
- Prager, D., Miller, D., and Gafford, D., (2002), “Terrain Interoperability in Large Federations: Correlation and Consistency in MC02”, Proceedings of IITSEC 2002, Orlando, FL, December 2002.
- Preparata, F. P. and Shamos, M. I., (1985), *Computational Geometry: An Introduction*. New York, Springer-Verlag.
- Richbourg, R. and Stone, T. (1998). “Automating error detection and correction in synthetic environments”. *Proceedings, Simulation Interoperability Workshop, Fall 1998*. Paper number 98F-SIW-087. <http://www.sisostds.org/confandwork.cfm> (2001, May29).