# CTACT: Consistent Terrain Analysis through Computational Topology

This document tries to answer Yehuda Avniel's questions and define the thrust for the Triangle:

- What are the problems that will be addressed? (Introduction & Goal)
- What are the roles of the individuals in the team (Team)
- How are the problems going to be addressed? (Tasks)
- What will the midterms and final exams be? (Roadmap)
- Who is the user/integrator going to be?  (Transition Path)

## Introduction

Continual improvement in sensor technology promises to gather massive amounts of detailed spatial data with less cost and less risk to equipment and personnel. The benefit by this new technology can only be realized, however, if the information implied in this data can be quickly extracted, succinctly represented, and correctly communicated at the appropriate level of abstraction to those who can use it most effectively. The bottleneck in the digital battlefield is not lack of data, but too much data and not enough information. In fact, the flood of data from new sensors can delay processing, obscure the essentials, and hinder the tasks that it is intended to support.

In military applications terrain analysis is performed by various echelons at different scales, often using overlay operations. Current systems use a number of different data representations, each tailored to a specific application and mapping technology.  Often there is a loss of information and consistency while transforming data between different representations, different levels of detail, and different types of discretization. Conflation – bringing together various layers of data – reveals topological inconsistencies that can create serious problems in military planning.  Processing and comparing topography (for example to detect changes) using heterogeneous data gathered by various mapping techniques at different times and scales can lead to misleading results and artifacts, ultimately affecting the decision making process in a negative (and sometimes dangerous) way – e.g.
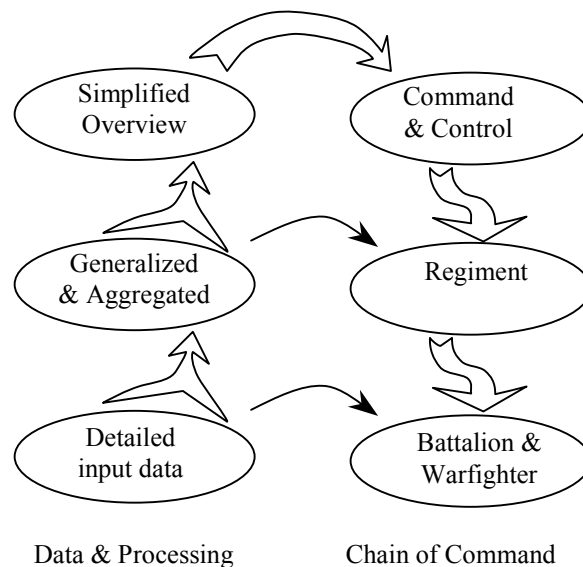


A smoke grenade marks position of $2^{nd}$ Battalion $502^{nd}$ Infantry Regiment for Apache helicopters in Kosovo, September, 2001.  (U.S. Army Photo by Sgt.1st Class Martin J. Cervantez)

removing an important terrain feature or creating artificial features or patterns that are not present.

Moreover, current techniques are not scalable to the large amounts of data that are now becoming available. They are susceptible to noise and uncertainty, and do not consistently handle many levels of representation. New efficient techniques for terrain classification, for extracting user defined features, and for draingage, trafficability and visibility analysis that use hierarchical representations of data and combine multiple data layers (e.g., LIDAR and hyper spectral data) are needed in military systems. These systems must be able to maintain features of the terrain and their interrelationships in a robust and globally consistent way across all scales.

Many teams are working on aspects of these problems in the representations embedded in commercial off the shelf (COTS) systems. We believe that new representations and new algorithms will be needed to achieve a complete solution to the following issues of representation and processing:

- Increasing data volumes: The shuttle radar topography mission (SRTM) is a good example of the rates at which data can be gathered and information processed. NIMA reports that "After 30 years of collection NIMA only had DTED for about 70 percent of the Earth's surface with measurements taken every 90 meters (DTED Level 1) and 5 percent with measurements taken every 30 meters (DTED Level 2). In only 9 days and 18 hours, SRTM collected elevation data for 80 percent of the world's landmass to enable the production of DTED Level 2. Areas the size of Alaska were mapped in 15 minutes and Florida in 90 seconds." But the project had 18 months of data processing before this data could become available.

- Data representations: Representations for commercial off-the-shelf (COTS) systems have dominated the research as well as the practice in Geographic Information Systems, despite their low level of abstraction. E.g., the raster vs. vector debate is a debate over mathematical representations that is carried on at the level of simple data structures. Others try to recast this debate at a higher level as the fields vs. objects debate. New sensors, such as LIDAR, show that more sophisticated representations are needed. The choice of representation also depends on how the terrain will be used for information display and further computation (image rectification, change detection or visibility/trafficability/drainage evaluation).

- Error: Spatial error in data is not represented directly, which makes it difficult to correct erroneous data by latter evidence.

- Batch processing: Commercial systems are oriented towards batch processing of spatial data. Instead, data could be interpreted as giving partial information samples about dynamic geometry.

- Generalization and overlay: Detailed data is often generalized and simplified before it is analyzed. Data from various sources is fused by digital overlay, a process inherited from manual cartography, which would overlay map layers on acetate. This data conflation operation brings problems of consistency between data from different sources and processing.

- Geological scale: Although seamless, multiscale maps are considered a Holy Grail in GIS, it is actually very important to represent different phenomena at appropriate map scales. This exacerbates the data conflation problems in generalization and overlay.

- Hierarchy of use: The hierarchy of military command makes the conflation problem more serious. The higher the level of command, the more generalized the pictures, so as to consider the overall situation. As decisions move back down the chain of command, it would be good if the relevant portions of the data could be restored in the generalization, so that at the bottom levels present the detailed data that can assist a war-fighter in the field. At present, data is generalized to prepare overall views for the top of the command hierarchy, but due to topological inconsistencies and large data volumes, the details on specific regions are not passed back down the



Data & Processing          Chain of Command

command hierarchy. With new, topologically consistent representations and I/O efficient processing, we will be able to add cross-links to solve this *de-conflation* problem.

# Project Goal

In this project we will provide:

- Rapid solutions to data conflation problems based on consistent topological representations, approximation techniques, and memory-aware computation,
- Solutions to de-conflation problems, enabling the delivery of consistent detail information to the levels at which it is may be effectively used, and
- Improved topographic analysis: change detection, visibility, trafficability, and information display.

As described below, we propose to focus initially on LIDAR, as an example of a sensor that produces large, detailed, focused data sets that must be rapidly processed to update representations. We incorporate other types of data as the project proceeds. The new representations and algorithms produced in this project will help users to rapidly incorporate up-to-date situational data into their planning processes and their communications at all levels of the command hierarchy.

# Project Team

By assembling an interdisciplinary team with expertise in mathematics, computer science, engineering, environment science, and geography, and by keeping close contact with the users (military personnel), we believe that this project will infuse GIS with new mathematical vigor, which is currently lacking, and will narrow the gap between GIS and mathematical and computational techniques.

The close interaction between computational and application researchers will provide a vertical integration of GIS – from developing richer representations of spatial data to designing and implementing algorithms that work with multi-level representation of data and that adapt themselves to the underlying computing resources. Our approach will use topological and geometric information explicitly and concurrently instead of the traditional approach in which topological information is derived implicitly and after the fact from the geometric information. This will clearly help to alleviate many of the problems that arise in topographic analysis due to noise, uncertainty, and inconsistency in data.

John Harer, lead-PI; Duke Mathematics
> Consistent topological representations, computational topology

Lars Arge; Duke Computer Science
> I/O and cache efficient algorithms, and their application in GIS

Helena Mitasova; NCSU Marine, Earth, and Atmospheric Sciences
> Spatial interpolation, feature detection, data fusion, open source GIS.

Jack Snoeyink; UNC Chapel Hill Computer Science
> Computational geometry in GIS, integration

Pankaj Agarwal; Duke Computer Science
> Geometric algorithms in GIS; approximation, dynamic algorithms

Pat Halpin; Duke School of the Environment
> Defining user needs and system integration

Lawrence Band; UNC Chapel Hill Geography
> Hydrography and its impact on terrain classification, representation, and analysis

Herbert Edelsbrunner; Duke Computer Science
> Computational topology.

# Project Tasks

We believe that we can achieve our goals of handling data conflation and de-conflation problems through new representations and algorithms. These will be developed through performance of several intertwined tasks, which we illustrate using LIDAR (Light Detection And Ranging) data. These tasks describe the representation,

management, analysis and management of massive terrain data sets. For each aspect they address issues of hierarchical structure, memory aware implementation, approximation and uncertainty in the data.
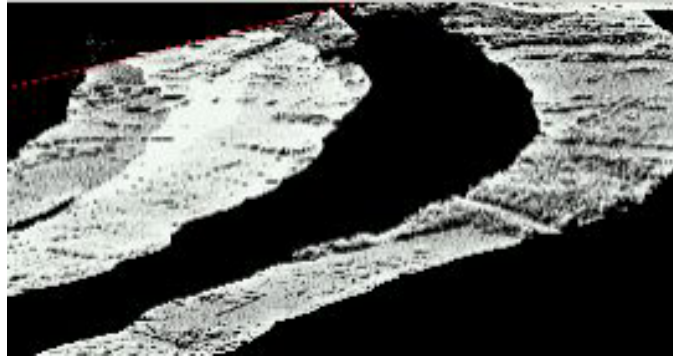
We list these tasks in roughly chronological order of their beginning. Because of the spiral nature of the development in this project, most of the tasks will be revisited with richer representations and larger data sets throughout the project. A little more detail on scheduling of the tasks is given in the Project Roadmap. Although many of these tasks are interdisciplinary and will be carried out in close collaboration within the project, we list the PIs who will be principally in charge of each.

## (i) Establishing close communication with interested users (Halpin, Arge)

Most of the PIs have participated in university/industry or university/government partnerships, and know the value of close communication with interested users to make sure that the right problems are being solved. Dr. Patrick Halpin is identified as the lead in this task because he served in the Special Forces before joining the faculty at Duke. As we interact with users, the specific data that we process and products that we compute will evolve, but we expect that we will continue to focus on representations and algorithms for massive geometric data sets that arise in GIS, and to face many of the issues that are illustrated below for LIDAR data.

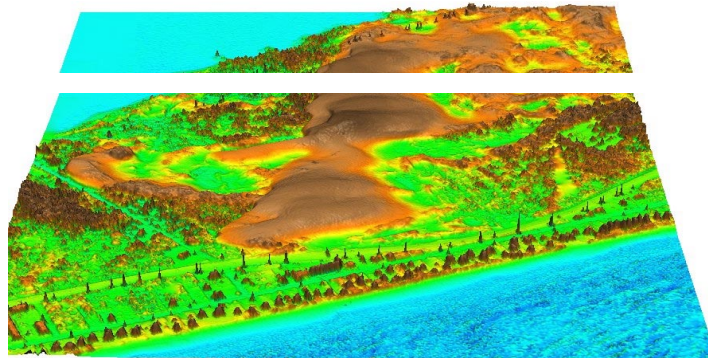## (ii) Handling of LIDAR data (Snoeyink, Mitasova, Arge)

Although general techniques will be developed in this project, we focus first on LIDAR data, and use it to illustrate the issues that arise when remote sensing advances lead to fast, automated acquisition of detailed (and massive) data sets. North Carolina is the first state to be fully mapped using LIDAR technology, providing us with a unique opportunity to develop and test the methodology using a large scale, real world data set, and to experiment with novel applications. Moreover, the coast of North Carolina has been mapped by LIDAR regularly providing time series of LIDAR data that will allow us to test the methods and algorithms for detection of subtle changes in topography.

Lidar point samples in swaths as the airplane turns

There are also regions which have high accuracy elevation data (such as 2ft contours) obtained by more traditional methods which will be used to develop techniques for combining LIDAR data with base maps to identify changes and update terrain information. We can map local areas that are experiencing rapid changes (due to construction, coastal erosion or other human or nature induced processes) using Real Time Kinematic GPS on ground. This can provide us with additional, detailed data for development of methods for detecting terrain changes and their impact on visibility, mobility, and landscape processes.

Although visual output is not a primary focus, it is a necessary tool in our research. The particular form of user output will depend on user requirements, but the representations and algorithms that we produce will allow us to adapt to displays ranging from display walls in command centers to hand-held devices. We will leverage from our previous work and from other projects in communication networks and display technology.

Relaxed spline fit to LIDAR data of Jockey Ridge, NC

## (iii) Mathematical representations (Harer, Snoeyink, Edelsbrunner)

New sensors, such as LIDAR, show that more sophisticated representations are needed – terrain is not a collection of square parcels with the same elevation value, nor a set of elevation contours, nor even described simply by a mathematical function $f(x,y)$. LIDAR collects dense measurements of the ground and vegetation in a narrow swath; terrain is then a set of functions defined over partial domains that are densely, but irregularly, sampled (with characteristic errors), and which may be overlaid with coarser data (such as DTED Level 2). The choice of representation depends on how the terrain will be used for information display and further computation (image rectification, change detection, or drainage, visibility, or trafficability evaluation). We recognize that phenomena change at different scales, and will seek for topologically persistent features that can be used to tie these phenomena into a consistent representation.

In creating and using hierarchical representations for terrain data, especially when we merge that data from several sources, we need to simplify by progressively suppressing information and to refine by processing more samples.

In defining simplification procedures, we must first assess the importance of particular features and set a threshold beyond which they are considered unimportant. Then we have to proceed with the simplification by eliminating those features below the threshold – treating them as noise in the data. To do this while maintaining topological consistency requires special attention to the global pattern of the data. For example, a small variation in the height of a terrain may suggest noise and require smoothing, but which ridge-line we follow to do this smoothing cannot be determined locally and instead must be determined by an appropriate global structure. The Morse Complex is one tool that we have used successfully to define procedures that smooth noise in certain data sets. Other global structures are possible, however, and will be developed as needed. We anticipate results not only for static terrain, but also for dynamic data where the structures will be 3-dimensional and, thus, more complex.

In defining refinement procedures, we aim to support the "hierarchy of use" diagram by rapidly identifying the relevant data and inserting it into the simplified picture while maintaining consistency. Keeping terrain data topologically consistent across varying geographic scales requires the maintenance of break lines and other structures. These structures form the basis for a less uniform representation that traditional grids. This will be crucial in dealing with enormous data sets like LIDAR. Mathematical characterization of these structures will be a fundamental task of our early work.

## (iv) Memory-aware computation (Arge, Agarwal)

Current algorithms do not scale to massive data sets, primarily because they do not take the hierarchical nature of modern memory systems into account. Different platforms (hand-held devices, laptops, desktops, mainframes) may have very different memory sizes and configurations, but they all use a hierarchy of larger and slower memory, further and further away from the processor, to give the impression of one large fast memory. To amortize the access time of slow memory, data is transferred between memory levels in large, contiguous blocks. Thus, it is important to design algorithms with a high degree of locality in their memory access patterns; that is, data accessed close in time should also be stored close in memory. To efficiently handle of massive datasets on a variety of platforms, we will develop "memory-aware" algorithms and data structures, designed to improve access locality and automatically adapt to the hierarchy of a given platform. Since the gap between the access times of main memory and disks is particular large, we especially consider I/O-efficient algorithms for processing disk-based LIDAR data on a desktop platform.

## (v) Feature detection (Snoeyink, Harer)

We plan to develop methods for automated extraction of break-lines and user-defined topographic features from over-sampled, noisy and/or heterogeneous elevation data (such as LIDAR). This will require an adaptive set of quantitative measures that adequately characterize a feature, allow the user to define the feature by setting those measures and enable the preservation of the structure of the features under refinements. Techniques that extract user defined features will be inherently geometric (local) and topological (global) in nature. Techniques from function theory, initially computational Morse Theory, but later involving more elaborate aspects of the theory of smooth functions, will be applied to LIDAR and other data in experiments designed to determine best-possible representations of features. Of course, these representations will need to be hierarchical and consistent with existing techniques of filtering noise of various sorts. We will also explore possible primitives that could be used to represent features by looking for recurring patterns, especially in regions with similar terrain.

## (vi) Change detection (Harer, Agarwal)

We will develop methods to combine and compare baseline elevation data given by contours or lower resolution DEM with new, mission specific data from modern mapping technologies, such as LIDAR, which are often high resolution, over-sampled, and noisy. This will require the development of methods to distinguish between the changes in topography and noise/artifacts produced by mapping (or cartographers fantasy). Temporal aspects of data are fundamental, since the variance of data from one survey to another is often the very thing that defines the features of interest. Key issues are how changes are detected and how they affect visibility, drainage, and trafficability. Tracking changes will be especially challenging when the data is not complete and when the mapping was performed by different techniques leading to different distribution and accuracy of data points.



A change detection scenario from Iraq.

The coast of North Carolina has been mapped by LIDAR regularly, providing time series data, which will allow us to test methods and algorithms for detection of subtle changes in topography.

## (vii) Terrain generalization (Snoeyink, Agarwal)

Most of the work on terrain simplification has focused on geometric simplification. Although these methods are suitable for visualization, other applications of terrain require the preservation of different cartographic, geological, or physical features. We propose to develop terrain simplification methods that consistently and accurately approximate application-specific properties at all levels, based on our representations of terrain. In particular, we will simplify to preserve computed properties, such as visibility or drainage.

## (viii) Terrain Characterization (Band, Mitasova)

We can incorporate into terrain simplification knowledge about the physical processes that shape terrain - developing methods that combine physical principles with geometry. We will try to decide if there are a few basic primitives, perhaps derived from the basic types of regions (coastal, piedmont, mountains) that define the features of interest.

In fluvially eroded terrain, the topographic skeleton formed by the complementary networks of ridges and drainage lines contain much of the geometric and topological information describing the landscape. Representation of these networks as formal graph models has been pursued for some time and we would implement similar methods for extraction, representation and generalization. Inclusion of break lines representing (geological) structural features and erosion/deposition transitions can extend this representation. Elementary hillslopes, for drainage analysis, are formed as surfaces subject to boundary conditions on these linear features, with specific characteristics that may vary on the basis of climate and dominant geomorphic processes. Terrain generalization the respects drainage can take the form of graph pruning, and reinterpolation or remapping of hillslope forms.

In non-fluvially-eroded terrain, a modified or augmented set of landscape primitives may be used. We have previously developed methods to extract and represent glacial features in alpine environments from combinations of terrain, spectral imagery, and knowledge-based methods including cirques, U-shaped valleys, and aretes. Similar modifications would be required for the expected forms of coastal or karst dominated terrain.

## (ix) Topographic analysis: drainage, trafficability, & visibility (Harer et al.)

Algorithms for topographic analysis (e.g. visibility, watershed regions, trafficability) are needed that can work with hierarchical representations of data. User can define the level of accuracy and details. The algorithms should produce consistent results at all scales. Since we are working with massive datasets, we cannot hope for exact solutions. In many applications, algorithms have to optimize conflicting goals. For example, to analyze the trafficability on a terrain, one has to consider the length of the path, its visibility, and the topography of the terrain along the path. In such cases, we need an approach that can interpolate between various parameters, depending on the application. Instead of reconstructing a function on the point samples (e.g. constructing a TIN on LIDAR data) and doing classical analysis of



Photo by Master Sgt. Michael J. Haggerty, US Air Force

the function, an alternative approach would be to transform the operations of classical analysis of the function to an algorithm that works directly on the samples. Efficiency is gained in this way by avoiding the reconstruction step provided we can work with sparse point samples, but the main benefit is that the essential information in the samples is used for the computation.

## (x) Data conflation (Halpin, Arge, Snoeyink, Mitasova)

This project will attack data conflation problems of increasing complexity. As can be seen in item (ii), algorithms for LIDAR data must deal with gaps and overlaps. Gaps may be handled by infusing the LIDAR data with other low-resolution data. Our new representations aim to support the replacement of data removed in simplification, and the addition of new data. We aim to rapidly update the results of topographic analyses with new data, a process that will become more sophisticated as the project proceeds. And new types of data, such as IFSAR, satellite imagery, GPS survey points, etc. will be included in our analyses, again as the project proceeds.

The selection of data types and regions to study will, we hope, be determined by user needs. We note, however, that we do have access to interesting data sets for areas in North Carolina and other locations where remotely sensed data sources have been backed up with extensive ground truth. In addition to LIDAR, there are regions with high accuracy elevation data (such as 2ft contours) obtained by more traditional methods, which will be used to develop techniques for combining LIDAR data with base maps to identify changes and update terrain information. Local areas that are experiencing rapid changes (e.g. due to construction, coastal erosion or other human or nature induced processes) are being mapped on ground using Real Time Kinematic GPS providing us with additional, detailed data for development of methods for detecting terrain changes and their impact on visibility, mobility as well as landscape processes. Optical and hyperspectral satellite imagery is also available and will be incorporated.

There have been significant efforts for using statistical techniques to classify hyperspectral data. We have some expertise in these techniques, but may wish to add additional expertise into this project, at a later stage, to classify based on combining hyperspectral data with the geometric information that we represent.

# Project Roadmap

This project begins with a focus on the efficient handling of LIDAR data, and spirals outward to incorporate other types of data, richer representations, and more complex tasks.

## Year 1:

In the first year of the project, the team will focus on efficient handling of the massive amounts of LIDAR terrain data now available. This requires the development of new algorithms that are scalable and able to deal with gaps and overlaps. Mathematically based concepts for feature representation as well as for detection of feature and feature change will be developed, together with the prototype software to test these concepts against the LIDAR data sets. An example would be a new representation for topographical structure for terrain that is annotated with cover parameters, such as height and density, which will later be used to evaluate trafficability and visibility. Especially challenging will be algorithms that adapt to the memory-hierarchy of modern machines. At the same time integrating the data with existing data and processing in an efficient way will create a variety of challenges. Our deliverable objective is a suite of tools to provide fast processing of LIDAR data, based on memory-aware algorithms.

We also have two internal objectives for the first year: First, to develop a prototype representation for the data and processing stages of the "hierarchy of use" diagram in the introduction. Second, to build the foundation of an effective relationship with Military user groups and to understand more carefully the particular challenges that they face. This will ensure that the team's research in years 2 and beyond addresses the right problems.

## Years 2 and 3:

In years 2and 3 the team will test the models developed in year 1 against a variety of data sets. This will be accompanied by the development of algorithms for computing watershed, drainage, visibility, and trafficability information from LIDAR data. Additionally the team will begin developing and implementing adaptive algorithms that work with hierarchical representations. The work on Morse complexes will be incorporated into topological simplification and denoising algorithms and uncertainty models will be introduced. These are all steps in the development of the "Chain of Command" aspects of our "hierarchy of use" diagram, and will be carried out in the context of the feedback obtained from the users identified in year 1. *The first major test of the team's success, therefore, will be the successful development of tools that address all elements of the Hierarchy of Use diagram.*

## Years 4 and 5:

The team will also revisit the prototype developed in the first three years, incorporating more types of data, accelerated interaction, and improved solutions based on the experience gained from the previous implementations and interaction with the users.

In the final years, the team will develop robust software libraries that can be incorporated into existing and new systems. The open source community has a well developed infrastructure for this and GRASS is already benefiting by using the PROJ library to support over 120 cartographic projections or GDAL library to import various raster data formats and fftl (fast fourier transform library) for image processing. Our product will be a library with demonstration modules which will illustrate how to use the library to create a module, e.g. for visibility analysis or other tasks, in GRASS, ArcGIS, or Manifold (or whatever the preferred GIS will be).

# Tangible Benefits

There are clear tangible benefits to the military from this project, including:

- Faster and more flexible processing of large elevation data sets.

- The presentation of information based on terrain data in the form, time, and detail needed for a particular task.

- A framework that supports generalized data for a high-level overview, and detailed data for those who need it for their missions.

- Faster and more reliable detection of changes in topography.

- Tools to effectively use of terrain information on a wide range of devices and computational environments.

- Topographic analysis (visibility, mobility, optimizing locations of camps, etc.) optimized for a given situation.

# Transition Path

Software application programmer interfaces (APIs), plug-ins, OpenGIS, object-oriented libraries, and GML (the Geographic Markup Language, an XML extension) will give the first transition paths for some aspects of this work. They cannot incorporate all aspects, because they focus on the low-level point, polygon and raster representations that are the mainstay of present GIS.

It is instructive, when considering a transition path, to look at how COTS systems incorporate new representations and algorithms by revolution, when they do incorporate them. ESRI's ArcInfo was built on 70s technology from Harvard University. The success of the company restricted innovation – although ESRI has incorporated some revolutions in spatial processing (e.g., the Spatial Engine, and adoption of COM technology), when Jan van Roessel wrote more accurate and faster overlay processing, it was incorporated as a (non-default) option, since changing the 1970's code would have changed the results for too many users. Manifold uses 80s technology to produce a powerful and more extensible COTS system, available at considerably less cost. They must still compete against ESRI's install base. GRASS, produced by the Army, is the first open system; Dr. Mitasova has had a long history of work on GRASS.

This project will also leverage from DARPA research into sensor technology, feature detection, data storage, and communication networks. It supports NIMA's goal of maintaining Foundation Data for the globe and fielding Mission Specific Data for areas of potential or actual involvement. The results of this project will help these reach their full potential to extract consistent information from data and provide it to the level at which it may be used effectively.