

# GEOBIA SYSTEMS FOR MASSIVE DATA PROCESSING

N. Ahles<sup>a</sup>, S. MacFaden<sup>a</sup>, J. O'Neil-Dunne<sup>a</sup>, A. Royar<sup>a</sup>, T. Engel<sup>a\*</sup>

<sup>a</sup> University of Vermont, Spatial Analysis Laboratory, Burlington, Vermont USA - (nahles, smacfade, jonieldu, aroyar, tengel)@uvm.edu

**KEY WORDS:** geographic object-based image analysis (GEOBIA); eCognition; LiDAR; multispectral imagery; massive data processing; ISPRS

## ABSTRACT:

Large portions of the Earth's landscape are now captured by high-resolution remotely-sensed datasets and turned into corresponding thematic data. Despite these advancements the number of comprehensive, high-resolution land-cover maps is surprisingly low. The value of high-resolution land-cover data in landscapes that are increasingly fragmented and heterogeneous is great, but so are the challenges associated with turning these disparate datasets into information. We argue that effective geographic object-based image analysis (GEOBIA) system design, while rarely discussed in the literature, is perhaps the most important factor in determining the success of projects whose focus is on broad-area mapping. Human resources, data, hardware, and software must be tightly integrated to make the system efficient and effective. At the same time, the object-based approaches used by such systems for land-cover mapping must try to replicate the human cognitive process as much as possible, using stable, context-based approaches to feature extraction that leverage the strengths of the various input datasets while minimizing their weaknesses. Drawing on our experience deriving 12 terabytes of high-resolution land cover for more than 232,000 km<sup>2</sup> in the United States, we describe the design considerations for GEOBIA systems that are capable of processing huge volumes of data. In addition, we provide examples of the techniques and approaches deployed within these systems that overcome the challenges associated with mapping land cover from massive, disparate datasets.

## 1. INTRODUCTION

The limiting factor of spatial analytics is no longer the lack of up-to-date remote-sensing data but our ability to rapidly turn those data into information. High-resolution multispectral satellite and aerial imagery has greatly proliferated in recent years, and some acquisition programs (e.g., Planet Labs) are now approaching single-day temporal resolution for the entire globe. LiDAR, while not as universally accessible as imagery, has also grown dramatically in availability and quality, and linear LiDAR systems continue to provide the market with datasets covering much of the developed world. LiDAR coverage and revisit rates will improve further with Single Photon and Geiger-Mode LiDAR systems, which are capable of dwarfing linear LiDAR acquisition rates. Much of the world has already turned these imagery and LiDAR datasets into thematic information in some capacity, most often in the form of vector GIS layers. Developed countries in particular have highly detailed datasets funded through governmental initiatives. Many developing countries have not yet invested in such GIS datasets but nonetheless have access to road networks, building footprints, and other vector features through crowd-sourced initiatives such as OpenStreetMap. However, very few areas of the globe have comprehensive, high-resolution land-cover maps. Such maps are crucial for the effective management and understanding of gray and green infrastructure (Benz, 2004).

A considerable amount of work in the GEOBIA field has focused on the development of individual algorithms, comparative analysis, and case studies. Relatively little has been published on what makes an effective GEOBIA system, particularly one capable of capitalizing on this modern area of vast quantities of geospatial data.

In this paper, we discuss the framework we have developed to integrate data preparation, GEOBIA, and manual editing into an iterative system that streamlines the production of land-cover datasets covering large geographic extents. High-resolution image analysis at county, state and regional scales necessitates efficient system design and processing at multiple steps, and personnel, hardware, software, and input datasets all must be coordinated in this effort. Specialized analysts and technicians are trained to manage isolated aspects of production, relying on multi-core workstations to maximize processing efficiency, standard operating procedures that codify workflows, videos and other training materials that expedite information transfer, extensive script templates and libraries that encapsulate pre-existing modeling routines, and data servers that allow cross-network processing and editing.

This framework moves GEOBIA from the realm of localized feature extraction to regional mapping efforts covering entire drainage basins and political administrative units (e.g., individual American states). We are currently mapping some of the largest sub-meter resolution land-cover projects ever attempted in the United States, including the Chicago Region (14,530 km<sup>2</sup>), the Chesapeake Bay Watershed (165,760 km<sup>2</sup>), and the Delaware River Basin (35,066 km<sup>2</sup>). Overall, our workflows have mapped land cover for more than 1 trillion pixels of data across North America. And size is not the only factor requiring efficiency; many of our large-extent mapping projects must be completed on very short timelines, often less than a year. The maps produced for these regions will document baseline conditions and serve as inputs for future land-cover and change-detection analyses.

---

\* Corresponding author

Many projects also focus on improving or updating existing datasets rather than mapping landscapes from scratch. GEOBIA is uniquely positioned to perform this type of mapping; it can ingest existing land-cover datasets and vector GIS layers (e.g., building footprints, road polygons), extract new features from updated imagery or LiDAR, and detect changes using contextual information. This process mimics human cognition (Blaschke, 2010). When a high density of roads and buildings is observed, it can be assumed with a high degree of probability that parking lots and other impervious surfaces are nearby. Similarly, extensive concentrations of large, tall buildings are much more likely to occur in an urban zone than an agricultural landscape. This understanding of contextual relationships, and the ability to incorporate them into GEOBIA rule sets, is essential to effective land-cover mapping across large, heterogeneous geographic areas (O'Neil-Dunne, 2011).

## 2. METHODOLOGY

### 2.1 Framework

A GEOBIA system is ultimately a collection of data, hardware, software and people. Development of accurate high-resolution, comprehensive land cover maps requires the effective synergy of these technological and human resources.

#### 2.1.1 Data

We have found that the most efficient and effective approach to high-resolution land cover mapping is to leverage all existing remotely-sensed and thematic datasets. There is no point in extracting features that have already been mapped. Moreover, any land cover mapping should insure consistency, to the extent possible, with existing mapped information. Flawed datasets are not excluded in our workflow, rather they are evaluated and understood so that they can be harnessed to minimize their limitations and maximize their valuable information. We draw from datasets that exist in raster, vector, and point cloud formats.

#### 2.1.2 Software

Although the term “image” is still part of the GEOBIA definition we are of the opinion that the true power of GEOBIA lies in its ability to serve as a data fusion platform in which the object breaks down the barriers that exist between traditional geospatial formats. As such, one of our key requirements is that the software be capable of handling point cloud, raster, and vector datasets in their native format without the need for conversion. Furthermore, the system must be capable of applying vector, raster, and point cloud algorithms. eCognition (Trimble) serves as the foundation of our GEOBIA system. eCognition fulfills our key requirements in terms of data fusion while at the same time supporting enterprise-level, distributed processing of massive datasets. In addition, we make use of a number of domain-specific software packages for data preparation and manual corrections. Key data preparation tasks include LiDAR point cloud classification, generating raster mosaics, and editing vector data layers.

#### 2.1.3 Hardware

A given regional high-resolution land cover project may entail working with hundreds of gigabytes or even terabytes of input datasets. In a GEOBIA processing object information is stored in computer memory, making large amounts of RAM a key requirement of our hardware systems. Multiple, high-performance CPUs also form a key foundation of our GEOBIA-

specific hardware. Much less stringent hardware requirements exist for our computer workstation used for manually editing. Online storage with fast networking and regular backups are also a key component of our GEOBIA system.

### 2.1.4 People

The human resource component of the GEOBIA system is the one we consider to be the most important. There exist two facets of the people component of our GEOBIA system, the first is the individual skills that each person possess, the second is the ability to function as a team. Personnel on our team are assigned to one of three sections: data preparation, GEOBIA feature extraction, and manual corrections. The data preparation section is responsible for assembling input datasets, including such tasks as mosaicking raster data, classifying LiDAR point clouds, and checking the quality and consistency of vector data. The GEOBIA section develops automated approaches to feature extraction using GEOBIA software. The manual corrections section has the responsibility of reviewing the output from the automated process, determining if it meets the standard, and making manual corrections to the data.

The ability to recognize features from a broad array of types of remotely sensed data is the foundational skill that we consider to be paramount to all team members regardless of their role. Image interpretation is a qualitative skill, one that is developed over time through the application of the elements of image interpretation. Our GEOBIA feature extraction analysts, who are the most seasoned members of our team spent many thousands of hours manually extracting information from remotely sensed data. We see this development as a crucial preparatory phase that enables them to harness the elements of image interpretation for automated feature extraction.

Communication is essential. We use a cloud-based team collaboration software to insure that all project personnel must understand project goals, the prescribed workflow, and their specific responsibilities in the workflow. By developing standard but adaptable protocols, it is possible to build a team of technicians and analysts (sometimes including 40 or more contributors) that works time- and cost-efficiently toward final products that satisfy end-user needs. A second but no less important goal is to document and store each final product for historical use: what were its inputs, how was it processed, and where is it located?

### 2.2 Approach

Our approach to broad-area, wall-to-wall land cover mapping with GEOBIA includes three primary phases: data preparation, GEOBIA rule-set development and processing, and manual editing (Figure 1). These tasks require different levels of technical expertise, ranging from undergraduate students who have completed GIS and remote-sensing coursework to seasoned research analysts with considerable experience in GEOBIA theory and practice. They also require unique considerations with respect to the geospatial hardware and software employed.

#### 2.2.1 Data Preparation

Data preparation includes acquisition and evaluation of all existing imagery, LiDAR, and vector datasets that could be useful in the feature extraction process. It also includes initial data processing and creation of derivative products that aid object segmentation and classification (e.g., rasterized LiDAR surface models). Our typical projects will incorporate anywhere from a

few to over twenty existing vector layers. These vector layers, representing individual features (e.g. buildings) or geographies (e.g. property parcel boundaries) are valuable, but most likely outdated when compared to the source imagery and LiDAR. A crucial part of the data preparation phase is the evaluation of these datasets, to understand their strengths, weaknesses, and their value in the GEOBIA feature extraction process. The review of these datasets is also crucial for determining whether or not

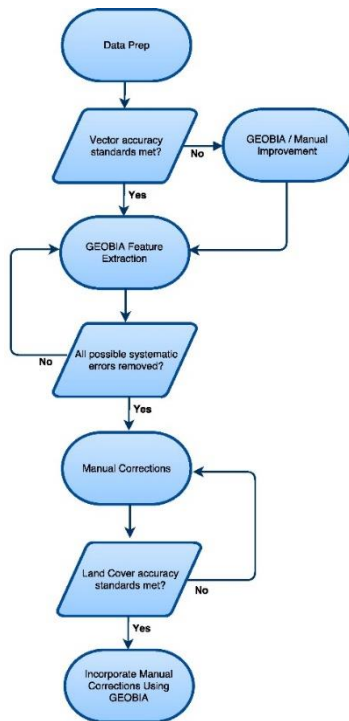


Figure 1. Flowchart depicting the GEOBIA process.

manual corrections should be made to the data prior to incorporating them into the GEOBIA system. An example of this is an incomplete building dataset that is serving as an input to a high-resolution land cover mapping project. If the dataset contains a relatively low number of missing or changed buildings manual editing would be a more efficient approach than an automated GEOBIA workflow. A larger number would make a manual editing process too costly. Determinations such as this are judgement calls, requiring input from both the GEOBIA feature extraction analysts and manual editing technicians.

### 2.2.2 GEOBIA Rule Set Development and Processing

In the GEOBIA step, rule-set development is the primary task, requiring iterative modification and testing in eCognition. Individual rule sets are generally produced for separate LiDAR collections, which in the United States is often by county but can sometimes include larger portions of states. This focus on LiDAR collections helps maximize data quality and temporal consistency, especially in counties or states that have coordinated programs for acquiring or updating imagery, LiDAR, and planimetric vector layers. However, the rule sets are usually structured with a generic modeling flow (e.g., preliminary identification of tall features followed by discrimination of tree canopy from buildings) that can be readily modified for use with other LiDAR collections. By building a library of rule sets, GEOBIA analysts can adapt modeling routines from project to project, re-using well-tested

approaches while adjusting specific segmentation and classification parameters as necessary; we do not start from scratch unless project goals or input datasets necessitate it. It is also possible to apply templates to specific modeling tasks depending on data availability. For example, if no planimetric building footprints exist for an individual study area, we can implement a previously-created routine that uses a combination of imagery and LiDAR to model footprints. This routine would replace a step that simply incorporates existing buildings into the draft classification. Similarly, if no roads polygons exist for an area of interest but road centerlines are available, the centerlines can be used to estimate road surfaces with a pre-existing routine designed for this purpose.

### 2.2.3 Manual Review

Once a rule set has been sufficiently honed by iterative testing, a point typically reached when all systematic errors have been addressed, a draft land-cover dataset is produced for manual review and editing. The team responsible for this effort usually consists of well-trained undergraduate students led by one or more experienced research technicians. Each member reviews individual tiles of the draft map, comparing it to the imagery and LiDAR used in initial modelling and identifying obvious errors of omission and commission (MacFaden, 2004). Polygons are then drawn around the misclassified features and labelled according to the specific changes required (e.g., a powerline erroneously classified as tree canopy must be re-assigned to the unclassified category). These polygons are then incorporated into a final classification in a second eCognition run (Figure 2). The manual corrections team follows a project-specific image interpretation key, and the lead research technicians review edits and provide feedback on the quality and volume of edits. The tiling scheme used to distribute sections for editing provides an effective system for gauging and documenting progress.

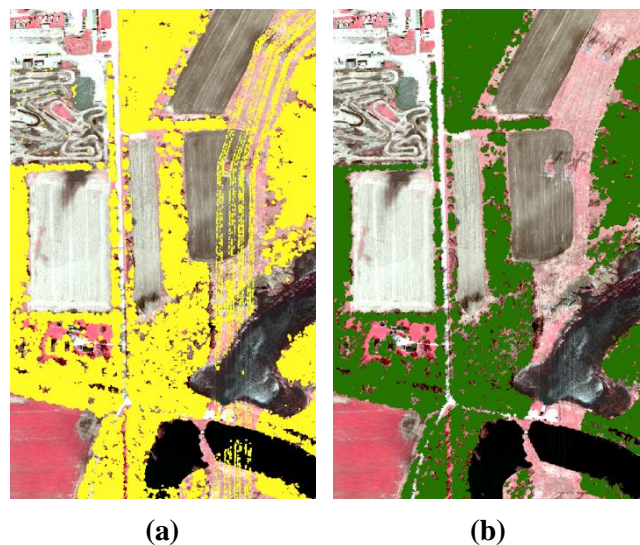


Figure 2. Tree canopy after automated feature extraction (a). Powerlines and other false tree canopy removed during manual review (b).

### 2.2.4 Data Management

Data management and organization are the pillars on which the GEOBIA workflow is built. Projects often require dozens of input layers possessing heterogeneous acquisition dates, data types, projections, resolutions, and quality (O'Neil-Dunne,

2013). As such, it is essential to know when and how each was acquired or developed.

### 2.2.5 Workflow Examples

GEOBIA is an iterative process, not only within eCognition-based modeling but also between automation and manual corrections. One way of mimicking human cognition in a rule set is by applying contextual metrics using distance and density maps derived from building footprints and road centerlines. By using these maps in tandem, it is possible to depict approximate areas of urbanization (O'Neil-Dunne, 2011). However, these contextual analyses are only as good as their input datasets. For example, when classifying buildings using a combination of current imagery and an older LiDAR dataset, new buildings constructed after LiDAR acquisition will inevitably be classified as non-building impervious surfaces due to their low Normalized Difference Vegetation Index (NDVI) value, used to identify healthy vegetation, and ground level elevation. This is an error that is not easily remedied using an automated approach, but a round of manual corrections focused on reclassifying misidentified buildings quickly fixes the problem. Once these errors are fixed, the new building dataset will improve the building-distance and density maps, which in turn improves the contextual analysis of agricultural fields, impervious surfaces, and sidewalks.

Similarly, existing planimetric vector datasets depicting buildings, roads, and other impervious surfaces often need to be improved before they can be incorporated into a land-cover rule set. If available imagery and LiDAR datasets are more recent than an existing vector layer, improvements can be performed using an automated approach. For example, an old set of building footprints can be augmented to remove non-building features such as ground level porches by segmenting with new LiDAR and then evaluating new aboveground features with a combination of simple height thresholds and the size of the resulting image objects (Figure 3). This approach maximizes the resources already devoted to the buildings layer and minimizes the errors of commission in unwanted features. A similar combination of LiDAR metrics and spectral values can be used to identify newly built buildings since the creation of the original dataset.

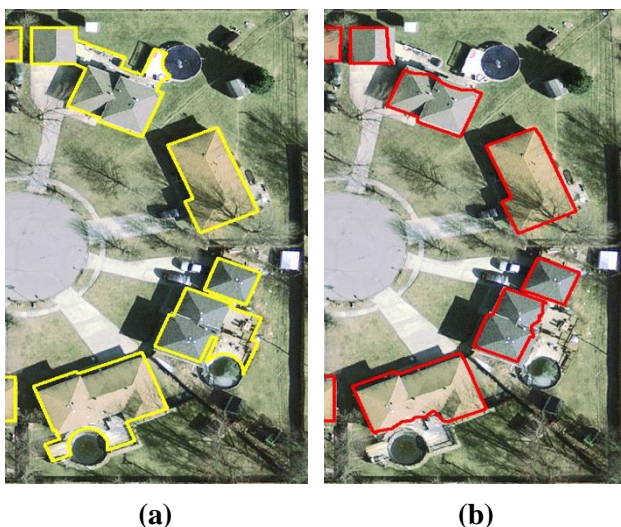


Figure 3. Old buildings vector dataset with porches (a). Porches automatically removed using OBIA (b).

### 3. CONCLUSIONS

GEOBIA is perhaps the only tool for converting the abundance of available remote-sensing datasets into timely, GIS-ready information. Much thought has been given to individual GEOBIA algorithms and applications, but little has been written on building effective GEOBIA systems, and the hardware, software, data, and human resources that are required to successfully apply this technology. Our experiences show that by implementing a framework that combines GEOBIA with standardized data processing and thorough manual review, it is possible to leverage previous investments in imagery and GIS data while providing essential land-cover data to policy makers and urban planners who need reliable, comprehensive maps for establishing baseline conditions and formulating tangible green-infrastructure goals. Additionally, this information is equally important to researchers who are working toward a better understanding of complex landscape phenomena and their social and environmental dependencies, including tree-canopy change, global carbon stocks, urban land use, water and air pollution, and human health. The time is now for GEOBIA to inform and shape landscape analysis and planning at broad regional scales.

### REFERENCES

- Benz, U.C.; Hofmann, P.; Willhauck, G.; Lingenfelder, I.; Heynen, M. Multi-resolution, object-oriented fuzzy analysis of remote sensing data for GIS-ready information. *ISPRS J. Photogramm. Remote Sens.* 2004, 58, 239–258.
- Blaschke, T. Object based image analysis for remote sensing. *ISPRS J. Photogramm. Remote Sens.* 2010, 65, 2–16.
- MacFaden, S.W.; O'Neil-Dunne, J.P.M.; Royar, A.R.; Lu, J.W. T.; Rundle, A.G. High-resolution tree canopy mapping for New York City using LiDAR and object-based image analysis. *J. Appl. Remote Sens.* 2012, 6, doi:10.1117/1.JRS.6.063567.
- O'Neil-Dunne, J.P.M.; MacFaden, S.W.; Pelletier, K.C. Incorporating contextual information into object-based image analysis workflows. In *Proceedings of ASPRS 2011 Annual Conference, Milwaukee, WI, USA, 1–5 May 2011*.
- O'Neil-Dunne, J.P.M.; MacFaden, S.W.; Royar, A.R.; Pelletier, K.C. An object-based system for LiDAR data fusion and feature extraction. *Geocarto. Int.* 2013, 28, 227–242.