

# EXTRACTING STRUCTURAL LAND COVER COMPONENTS USING SMALL-FOOTPRINT WAVEFORM LIDAR DATA

*J. McGlinchy<sup>a</sup>, J. Van Aardt<sup>a</sup>, H. Rhody<sup>a</sup>, J. Kerekes<sup>a</sup>, E. Ientiluci<sup>a</sup>, G.P. Asner<sup>b</sup>, D. Knapp<sup>b</sup>, R. Mathieu<sup>c</sup>, T. Kennedy-Bowdoin<sup>b</sup>, B.F.N. Erasmus<sup>d</sup>, K. Wessels<sup>c</sup>, I.P.J. Smit<sup>e</sup>, J. Wu<sup>a</sup>, D. Sarrazin<sup>a</sup>*

<sup>a</sup>Rochester Institute of Technology, Rochester, NY, USA

<sup>b</sup>Carnegie Institution for Science, Stanford, CA, USA

<sup>c</sup>Council for Scientific and Industrial Research, Pretoria, South Africa

<sup>d</sup>School of Animal, Plant and Environmental Science, University of the Witwatersrand, Johannesburg, South Africa

<sup>e</sup>Kruger National Park Scientific Services, Skukuza, South Africa

## ABSTRACT

Previous work has shown the ability of waveform LiDAR sensors to accurately describe various land cover types [1] and biomass estimates made in the field [2]. What is lacking, however, is a way to describe the different structural components that are embedded in the digitized backscattered energy from the LiDAR pulse. This study aims to extract structural components from waveform LiDAR data in terms of woody, herbaceous, and bare ground components from data collected over a savanna environment in and around Kruger National Park (KNP), South Africa. These components are comprised of metrics extracted from the waveforms and validated using biomass measurements made in field plots. Different size windows around plot centers, 3x3 pixels and 9x9 pixels (resulting in 1.5m and 4.5 m footprint, respectively), were used to examine scale effects of larger footprints. It was found that composite waveforms resembling plot sizes (9x9) most often are able to describe more than 80% of the woody biomass variability across the entire study site, and individually for two of the three land uses within the area. However, the herbaceous component of the waveform did not correlate well with the field measurements, while the bare ground component was verified visually in a side-by-side comparison with optical imagery.

**Index Terms**— LiDAR, waveform, modeling, structure

## 1. INTRODUCTION

Remote sensing using light detection and ranging (LiDAR) technology has seen considerable advancement with the advent of full waveform digitizing sensors. LiDAR remote sensing systems operate by transmitting a monochromatic light pulse and measuring the reflection of this light pulse off of a scattering surface. The intensity of the laser pulse is recorded as a function of the time it takes for the energy to leave the emitter, interact with the surface, and return to the sensor. Waveform LiDAR sensors have the advantage of being able to record the backscattered energy at a very high sampling rate, typically on the order of nanoseconds. The combination of high temporal resolution detection and full backscatter digitization enable the extraction of structural information that is embedded within the waveform [3]. Various studies have shown that signal metrics, calculated from large footprint LiDAR waveforms (on the order of 10s of meters), can be used to assess vegetation structure in forested environments

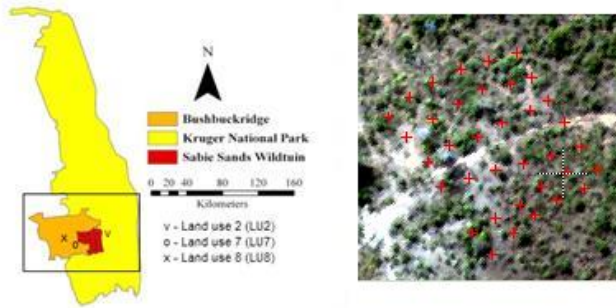
[e.g., 4, 5], while small-footprint LiDAR waveforms can be used to accurately classify various land cover types [1]. Measures such as tree height, crown volume, and biomass have been accurately predicted and modeled, resulting in good correlation between waveform-derived metrics and available field data [e.g., 2]. However, two specific challenges remain in terms of land cover assessment: (i) most previous work has dealt with large-footprint systems, which results in the measured field data typically being an order of magnitude smaller in actual ground area than the footprint size and (ii) a detailed breakdown of woody, herbaceous, and bare ground structural components along the laser trajectory, similar to the "end member" concept in an imaging spectroscopy context [6, 7], is still lacking. This latter aspect has bearing on our ability to map land cover types in the structural (3D) domain, as opposed to the traditional spectral approaches.

The objectives of this study are to (1) establish a method by which to extract structural components, e.g., woody, herbaceous, and bare ground from small-footprint LiDAR waveforms, (2) assess how these components and their extraction vary across different footprint sizes, and (3) establish how these structural components can be mapped across the landscape. We will accomplish this by using plot-level waveforms, generated by compositing small-footprint waveform LiDAR (0.56 m footprint) returns, and extracting waveform-derived metrics to identify unique structural components and map woody and herbaceous biomass for both a conserved and communal savanna land use area. This scalable approach will increase our understanding of the interaction between waveform footprint and land cover object sizes and aid in the development of improved relationships between structural waveform metrics and measured field data.

## 2. STUDY AREA AND DATA

The study area (Figure 1) is located in and around the Kruger National Park (KNP) in South Africa. The area is bounded by (22°8'00"S; 30°34'52"E) and (25°32'48"S; 32°2'50"E). Field and remote sensing data were collected for structural assessment of land degradation across a land use gradient that includes the KNP and an adjacent subsistence farming, communal area; this layout effectively juxtaposes a "protected" and "communal" area (Figure 1, left). An example of a degraded communal savanna site from the study area is shown in Figure 1 (right), which also shows placement of plot-level field data on 10 m grid spacings. The field

data are based on 4-5 sites per land use type for a total of 9 sites.



**Figure 1** Left: Study area for this research. Our study focuses on the protected Kruger National Park and degraded Bushbuckridge (communal) areas. Right: An example of a degraded communal land use site in Bushbuckridge, South Africa (LU7). Field plots in red.

Each site, in turn, consists of 36 plot-level measurements of herbaceous biomass, tree height and diameter, species, and a qualitative assessment of cover (crusting, bare soil, herbaceous, and woody cover).

The field data were collected during May 2008 in association with an airborne data collection campaign, and are summarized in Table 1; woody biomass calculations were derived from allometry equations. Waveform LiDAR data were collected by the Carnegie Airborne Observatory (CAO), using an Optech waveform digitizer, at 0.56 m footprint size and 1 ns temporal (vertical) resolution [8].

### 3. METHODS

#### 3.1. Bare ground component

Bare ground is one structural component that is present in every waveform. Removing and quantifying this component potentially will allow us to determine how much bare ground is present in a waveform in the 2D domain, i.e., without vertical interactions.

Three measurements were taken from the waveform once this component was removed. They are defined as “Road Ratio”, “Leftover”, and “Ratio Removed”. “Road Ratio” is measured as the ratio of an amplitude scaled dirt road sample to an original dirt road waveform sample extracted from LU8. “Leftover” is measured as the ratio of the sum of what remains in the ground pulse to the sum of these same points in the original waveform. “Ratio Removed” is measured simply as the area under the road-removed waveform divided by the area of the original waveform (See Figure 2).

A combination of these measures is hypothesized to provide a means to identify the amount of bare ground present in a waveform. These measures range in value from 0 to 1, can be performed on a per-waveform basis, and are independent of any amplitude variations observed waveform to waveform, given the normalization approach.

#### 3.2. Waveform features

Once the bare ground was removed, the composite waveforms were formed by summing values in identical height bins of waveforms within the pixel window. From these composite

**Table 1** Land use biomass summaries (‘w’ denotes woody, and ‘h’ denotes herbaceous). Land use defined by LUx, where ‘x’ is 2, 7, or 8 (see Figure 1).

LUx (sample size)		kg (w)	kg/ha (w)	g (h)
LU2 (82)	median	20.24	7.16E+03	34.00
	mean	59.62	2.26E+04	39.84
	std	104.00	4.08E+04	19.52
LU7 (67)	median	4.57	1.23E+04	20.00
	mean	174.05	5.25E+05	27.37
	std	1085.33	3.46E+06	21.80
LU8 (32)	median	4.83	4.98E+03	26.00
	mean	21.67	5.00E+04	30.24
	std	64.77	1.87E+05	16.19

waveforms we extracted numerous measurements related to plot structure. Metrics such as height of median energy (HOME), canopy energy, and ground energy developed by the United States Geological Society (USGS) [9] and Neunswander *et. al.* [1] were extracted along with numerous other measures.

Firstly, two height related metrics were extracted. One is similar to the canopy ratio (CRR) defined by the USGS, but adapted to focus on the energy of the first backscattered return in the waveform as opposed to the signal immediately above the ground pulse. This is defined as ‘aCRR’. The other is simply the volume of the plot as measured by taking the height of the first interaction in each waveform in a plot and multiplying by the spatial resolution of the data (0.56 m). This is defined simply as VOL.

Secondly, the time (in nanoseconds) it takes to complete the 10-90% integration range of the entire integrated waveform energy was extracted. This is defined as the duration of the waveform, where a longer duration is proportional to the presence of complex above ground structure.

Lastly, summary statistics (mean, mode, median, standard deviation, variance, and range) of the first and second derivative were taken, including the zero- and nonzero-values of the waveforms. Contrary to the aforementioned features, the derivative measures were performed on the composite waveform prior to ground removal. Also, these derivatives were peak normalized because of the amplitude dependency of the discrete derivative process.

#### 3.3. Regression and component mapping

These measurements were taken from plot waveforms that comprise synthetic footprints of 5 meters and 1.5 meters, i.e., 9x9 and 3x3 pixel windows, respectively. The waveform features were incorporated into forward regression models ( $\alpha=0.1$  for entry) to predict woody and herbaceous biomass measured at the plot level. Features that best explain the variance in these measurements were measured site-wide to show how the structural components could be mapped.

Past work [10] has shown nonlinear relationships between waveform LiDAR metrics and field measured biomass. For this reason, various transformations of the field data and waveform features (e.g., squared/square root value) were tested to determine ideal fit. This was performed across all land uses (generic model), across each individual land use (site specific models), and across

**Table 1** Regression results for modeling woody and herbaceous biomass from waveform features. Analysis performed across all land uses (LU278), land uses 7 & 8 (LU78), and each land use (LUx). “var” is the biomass transformation commensurate with R<sup>2</sup> value to its left.

Woody best R <sup>2</sup> at SL $\alpha = 0.1$										
	LU2 (w)	var (w)	LU7 (w)	var (w)	LU8 (w)	var (w)	LU78 (w)	var (w)	LU278 (w)	var (w)
best	0.4819	ln wha	0.8946	w ha	0.9917	w ha	0.8088	sqrt wha	0.8049	w ha
z3x3	0.4054	ln wha	0.8156	w ha	0.963	w kg	0.8088	w ha	0.5013	sqrt ha
nz3x3	0.3714	ln wha	0.8688	w kg	0.8796	sqrt wha	0.4428	sqrt wha	0.2253	sqrt ha
z9x9	0.4819	ln wha	0.8946	w ha	0.9917	w ha	0.8048	w ha	0.8049	w ha
nz9x9	0.4579	ln wha	0.7054	w ha	0.9331	sqrt wha	0.6304	w ha	0.4827	w ha
Herbaceous best R <sup>2</sup> at SL $\alpha = 0.1$										
	LU2 (h)	var (h)	LU7 (h)	var (h)	LU8 (h)	var (h)	LU78 (h)	var (h)	LU278 (h)	var (h)
best	0.1932	sqrt hg	0.2577	ln hg	0.538	hg	0.2604	ln hg	0.2556	sqrt hg
z3x3	0.1357	sq hg	0.1478	ln hg	0.5099	sq hg	0.1188	ln hg	0.1994	hg
nz3x3	0.1932	ln hg	0.1492	ln hg	0.3753	hg	0.1367	ln hg	0.1593	hg
z9x9	0.1461	hg	0.2577	ln hg	0.538	hg	0.2586	ln hg	0.2556	ln hg
nz9x9	0.1271	hg	0.2577	ln hg	0.4848	hg	0.2604	ln hg	0.216	ln hg

land uses 7 and 8 as a single data set (conservation vs. communal models). Linear representations of the data were also modeled.

The bare ground component was evaluated by comparing the three ground-removal components to the visual estimate made by field researchers. It should be noted this is a purely subjective measurement being evaluated with objective data.

#### 4. RESULTS AND DISCUSSION

The results from this analysis (summarized in Table 2) showed that the features extracted from the waveform can be applied in modeling woody biomass estimates in a savanna environment. Importantly, all woody biomass estimates were adequately modeled in a linear fashion (adjusted R<sup>2</sup>=0.80) when scaled from per-meter to per-hectare measurements. The equation that describes this general model is (scale in subscript, land use in superscript)

$$\begin{aligned}
 w_{ha}^{278} = & 3.79 \times 10^8 + 2.19 \times 10^8 (\text{leftover}) \\
 & - 1.26 \times 10^8 \left( \text{leftover}^{\frac{1}{2}} \right) - 1.25 \times 10^8 (\text{leftover}^2) \\
 & + 1.17 \times 10^9 (\text{ratioOut}) - 1.25 \times 10^8 \left( \text{ratioOut}^{\frac{1}{2}} \right) \\
 & - 2.83 \times 10^8 (\text{ratioOut}^2) + 3.35 \times 10^5 (\text{HOME}) \\
 & - 3.79 \times 10^5 (\text{CRR}) - 6.42 \times 10^5 (\text{aCRR}) \\
 & + 3.36 \times 10^9 (\text{mode}_1) - 8.33 \times 10^8 (\text{mode}_2).
 \end{aligned} \quad (1)$$

with <sub>1</sub> and <sub>2</sub> subscripts denoting first and second derivative.

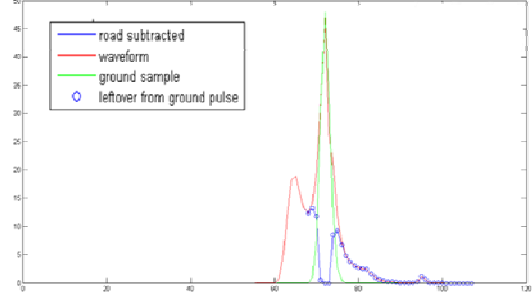
When land use 7 and 8 were considered together, the waveform features were able to explain 80.48% (adjusted R<sup>2</sup>=0.80) of the variation in woody biomass estimates. The equation that describes this model is

$$\begin{aligned}
 w_{ha}^{78} = & 1.22 \times 10^6 - 1.89 \times 10^6 (\text{CRR}) \\
 & + 6.09 \times 10^9 (\text{mode}_1) - 2.76 \times 10^9 (\text{mode}_2).
 \end{aligned} \quad (2)$$

When each land use was considered separately, the waveform features were able to explain 48.19% (adjusted R<sup>2</sup>=0.38), 89.46% (adjusted R<sup>2</sup>=0.89), and 99.17% (adjusted R<sup>2</sup>=0.99), of the variation in woody biomass estimates for land use 2, 7, and 8, respectively. The equations that describe these models are

$$\begin{aligned}
 \ln(w_{ha})^{LU2} = & 13.45 - 4.42(\text{ratioOut}^2) + 0.05(\text{duration}) \\
 & - 5.33(\sigma_2) - 16.52(\mu_2),
 \end{aligned} \quad (3)$$

$$\begin{aligned}
 w_{ha}^{LU7} = & 1.58 \times 10^6 - 2.67 \times 10^6 (\text{CRR}) + 2.39 \times 10^9 (\text{mode}_1) \\
 & - 3.98 \times 10^{25} (\text{range}_2),
 \end{aligned} \quad (4)$$



**Figure 2** Result of removing the dirt road sample from a multiple interaction waveform. The dirt road sample is aligned with the ground pulse of the waveform and removed, leaving the points marked in circles from the ground pulse of the original waveform.

and

$$\begin{aligned}
 w_{ha}^{LU8} = & -4.38 \times 10^5 + 3.92 \times 10^5 (\text{ratioOut}^2) \\
 & + 8.64 \times 10^3 (\text{VOL}) + 3.80 \times 10^4 (\text{aCRR}) \\
 & + 6.69 \times 10^8 (\text{mode}_1) - 3.51 \times 10^8 (\text{mode}_2) \\
 & + 2.45 \times 10^5 (\sigma_2).
 \end{aligned} \quad (5)$$

The best model chosen to model woody biomass estimates for land use 2 unfortunately did not adequately explain the variance in estimates (<50%). This was attributed to the narrow range of biomass, and influence of outliers.

Disappointingly, none of the waveform features were able to create a model that efficiently described the herbaceous biomass estimates at the plot level. This was not wholly unexpected, as the herbaceous backscatter of the waveform is embedded in the ground pulse. Since the “bare ground” component was removed from the ground pulse, some information related to the herbaceous content may also have been removed. It would appear that previous work investigating multiple scattering of bi-modal waveforms generated more accurate regression models in estimating herbaceous biomass at the plot level (e.g., Wu *et al.* [10]).

The regression models for the woody biomass for land uses 2, 7, and 8 did not all contain the same features. Also, the model for woody biomass across the entire study site contained features that were not present in any of the individual land use models. Of note is the number of features that were included from removing the bare ground component from the waveform.

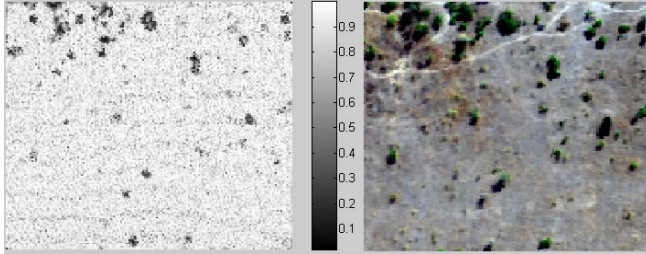
The features that best described the biomass in land uses 7 & 8 did show overlap. Both models incorporate a ratio of the canopy energy to the total energy (aCRR for LU8 and CRR for LU7). LU8’s model was more dependent on the upper canopy elements, while LU7’s model was more dependent on the entire canopy (upper and lower). LU2’s model was also dependent on CRR, but more on the most frequently occurring values (mode) of the first and second derivative. This reflects the ecology of the different land uses; the rangelands in LU8 are well-kept so the vegetation density at higher heights is more prominent than is the case with the other land uses. LU7 is heavily utilized with fewer large trees but a more uniform and dense shrubby population, so the woody components are physically more variable at lower heights. LU2 is in the protected KNP area, and the pattern and density of the biomass in this land use is driven more by wildlife than human interaction.

Finally, no model was able to adequately explain the visual bare cover field estimates using the bare ground removal

measurements. However, the model that performed best can be expressed as

$$\text{bare ground \%} = \text{ratioOut} \quad (6)$$

For example, if only 5% of the original waveform is present after subtracting the bare ground, then the waveform contains 95% bare ground. This can be verified visually in Figure 3.



**Figure 3** “Ratio Out” metric for a site located in protected land use area (KNP).

## 5. CONCLUSION

From this study we have extracted features from full waveform small-footprint LiDAR data that accurately and linearly describe the woody component in a savanna environment. By visual confirmation, the bare ground component has also been extracted from the data. Future work will focus on how the herbaceous component can be more accurately extracted from the full waveform LiDAR data and the development of a more general biomass model for this area that is not so specific to the data set used in this study.

## 6. ACKNOWLEDGEMENTS

We would like to acknowledge the Carnegie Institution for Science, the Council for Scientific and Industrial Research, and the Rochester Institute of Technology for their scientific and monetary support of this project.

## 7. REFERENCES

- [1] A. Neuenschwander, L. Magruder, and M. Tyler, “Landcover classification of small-footprint, full-waveform lidar data”, *Journal of Applied Remote Sensing*, vol. 3, 2009.
- [2] J. Drake, R. Dubayah, R. Knox, D. Clark, and J.B Blair, “Sensitivity of large-footprint lidar to canopy structure and biomass in a neotropical rainforest”, *Remote Sensing of Environment*, Vol. 81, pg. 378-392, 2002.
- [3] C. Mallet, F. Bretar, “Full-waveform topographic lidar: State-of-the-art”, *ISPRS Journal of Photogrammetry and Remote Sensing*, Vol. 64, pg. 1-16, 2009.
- [4] M. A. Lefsky, W. Cohen, D. Harding, G. Parker, S. Acker, and S. Gower, “Lidar remote sensing of above-ground biomass in three biomes”, *Global Ecology and Biogeography*, vol. 11, pg. 393–399, 2002.

- [5] J.B. Boudreau, R. Nelson, H. Margolis, A. Beaudoin, L. Guindon, and D. Kimes, “An analysis of regional aboveground forest biomass using Airborne and Spaceborne LiDAR in Québec”, *Remote Sensing of Environment*, vol. 112, pg. 3876-3890. 2008.

- [6] G. Asner, C. Borghi, and R. Ojeda, “Desertification in Central Argentina: Changes in Ecosystem Carbon and Nitrogen from Imaging Spectroscopy”, *Ecological Applications*, Vol. 13, 629-648, 2003.

- [7] N. Keshava and J. Mustard, “Spectral Unmixing”, *IEEE Signal Processing Magazine*, Vol. 19, pg. 44-57, January 2002.

- [8] G.P. Asner, D.E. Knapp, T. Kennedy-Bowdoin, M.O. Jones, and R.E. Martin. “Carnegie Airborne Observatory: In-flight fusion of hyperspectral imaging and waveform light detection and ranging (wLiDAR) for three-dimensional studies of ecosystems.” *Journal of Applied Remote Sensing* 1:1-21. 2007.

- [9] USGS DSP: Lidar-based Vegetation Metrics, USGS, 2009. [http://ngom.usgs.gov/dsp/mapping/lidar\\_vegetation\\_metrics.html](http://ngom.usgs.gov/dsp/mapping/lidar_vegetation_metrics.html)

- [10] J. Wu, J. van Aardt, G.P. Asner, T. Kennedy-Bowdoin, D. Knapp, B.F.N. Erasmus, R. Mathieu, K. Wessels, and I.P.J. Smit. “Lidar waveform-based woody and foliar biomass estimation in savanna environments”. *Rochester Institute of Technology*.