A SPECTRAL CLIMATOLOGY FOR ATMOSPHERIC COMPENSATION OF HYPERSPECTRAL IMAGERY

by

John H	I. Powell
A Dis	sertation
Submit	ted to the
Gradua	te Faculty
	of
George Mas	son University
in Partial F	Fulfillment of
The Requireme	nts for the Degree
	of
Doctor of	Philosophy
Computational Scie	ences and Informatics
<u>F</u>	
Committee:	
	Dr. Kirk Borne, Dissertation Director
Dr. Ronald Resmini, Committee Membe	
Dr. Mike Summers, Committee Member	
Dr. Ruixin Yang, Committee Member	
Dr. Maria Dworzecka, Interim Director, School of Physics, Astronomy, and Computational Sciences	
	Dr. Donna M. Fox, Associate Dean, Office of Student Affairs & Special Programs, College of Science
	Dr. Peggy Agouris, Dean, College of Science
Date:	Spring Semester 2015 George Mason University Fairfax, VA

A Spectral Climatology for Atmospheric Compensation of Hyperspectral Imagery

A Dissertation submitted in partial fulfillment of the requirements for the degree of Doctor of Philosophy at George Mason University

by

John H. Powell Master of Science Naval Postgraduate School, 1996 Bachelor of Science University of Colorado, 1986

Director: Kirk Borne, Professor School of Physics, Astronomy, and Computational Sciences

> Spring Semester 2015 George Mason University Fairfax, VA

Copyright 2015 John H. Powell All Rights Reserved

DEDICATION

This work is dedicated to my parents, Nathalie Powell and John W. Powell, for instilling in me a lifelong love of learning.

ACKNOWLEDGEMENTS

First and foremost, I would like to thank my committee members for their inspiration and support, both in the classroom and through the course of the research. I am grateful to Drs. Ruixin Yang and Mike Summers, for their sage advice and willingness to help. Special thanks to Dr. Kirk Borne, for his steadfast support throughout the process, and for helping me understand the power of scientific databases and big data before it was "big". Finally, I would like to express my endless gratitude to Dr. Ron Resmini, for introducing a brilliant idea and letting me run with it, for keeping me on the right track, and for helping a weatherman become a "seasoned spectral practitioner."

I also gratefully acknowledge the support and assistance of my supervisors, friends and colleagues who have helped me along the way. Thanks to Dr. Ernie Reith, Mr. Wayne Hallada, Mr. John "Boom" Minto, Ms. Connie Gray, Mr. John Cole, and Ms. Elizabeth Sweet. Special thanks to Ryan DaRe, David Adams, Adam Edson, Elizabeth White, Tom Smith and Shannon Jordan for their technical advice and assistance.

I would like to thank the administration of the College of Science and the George Mason Office of the Provost for their support of my work.

Lastly, thank you to Sara, Joseph, Laura, Ryan, Eric and Quinn for your patience, love and encouragement during this long journey.

TABLE OF CONTENTS

Page
List of Tables
List of Figures
List of Abbreviations xi
Abstractxii
1. Introduction and Literature Review
1.1 Hyperspectral remote sensing
1.2 Atmospheric Compensation
1.2.1 Empirical Methods11
1.2.2 Physics-Based Methods
1.2.3 Alternative Methods
1.3 Comparisons
1.4 Environmental Databases
1.5 Research Objectives
2. Data Description
2.1 HYDICE Sensor
2.2 Ground Truth Measurements
2.3 Imaged Sites
3. Methodology
3.1 Remote Sensing Model
3.2 ELM Process
3.3 QUAC Process
3.4 ELM Postprocessing
3.5 Coefficient Modeling
3.6 Coefficient Standardization
4. Scientific Database
4.1 Research Need

	4.2	Application Overview	72
	4.3	Database Design	79
	4.4	User Interface	80
	4.5	Validation and Application Examples	83
	4.5.1	Data Validation	84
	4.5.2	Application Examples	86
5.	Anal	ysis	92
	5.1	Modeled Coefficient Optimization	92
	5.2	Illumination Adjustment Validation	103
	5.3	Geometric Adjustment Validation	106
	5.4	Seasonal Adjustment Validation	108
	5.5	Global Coefficient Standardization	110
	5.6	Coefficient Classes	117
	5.6.1	K-means clustering	118
	5.6.2	Spectral similarity clustering	123
	5.6.3	Cross validation	136
	5.7	Example Application	139
6.	Conc	clusion	148
A	ppendix	۲	151
	A1. EN	IVI (Environment for Visualizing Images)	151
	A2. MA	ATLAB	152
	A3. OP	PeNDAP	156
	A4. M0	ODTRAN	157
	A5. My	VSQL	158
R	eference	es	160

LIST OF TABLES

Table		Page
Table 1	Summary of site characteristics and imaging. Climate and land cover	
categories a	re described in the following tables.	33
Table 2	USGS Land Cover Institute (LCI) land cover class definitions [74]	33
Table 3	Köppen–Trewartha climate classification categories.	34
Table 4	RMSE residuals for linear, nonlinear, and adjusted linear regressions,	
averaged ac	ross all wavelengths and all image coefficients in the study	55
Table 5	Spectral similarity metrics for the adjustment procedure reflectance	
validation.		63
Table 6	List of major radiative transfer driver parameters (common to all	
MODTRAN	V runs)	95
Table 7	MODTRAN Atmospheric and scattering input parameters for each site	. 102
Table 8	List of standard reference coordinates for coefficient comparison	. 111
Table 9	Summary of cluster statistics.	. 136
Table 10	Statistics for site membership identification	. 138
Table 11	Summary of reflectance retrieval similarity metrics across all four panels	5.147
Table 12	Listing of significant custom MATLAB scripts and functions used	. 154

LIST OF FIGURES

Figure Page
Figure 1 Alunite spectrum as measured by multi- and hyperspectral instruments. Source:
Clark
Figure 2 HSI schematic diagram
Figure 3 Solar illumination at TOA and Earth surface for typical atmospheric conditions.
Major absorption bands are labeled. Source: Rhode [13]
Figure 4 Schematic diagram of reflected and scattered radiance paths detected by an
airborne sensor. The ellipse represents the sensor's IFOV 10
Figure 5 Geographic locations of sites. The marker in the California Sierra Nevada
Mountains represents three separate sites
Figure 6 Current images of a sample of the sites showing representative land cover
categories. Images captured from Google Earth
Figure 7 Color composite images of calibration panels (left) and with ELM ROIs overlaid
(right) for approximate flight levels of 5000 ft (upper figures) and 20,000 ft (lower
figures) AGL
Figure 8 Typical ground truth reflectance spectra of six gray shade calibration panels 40
Figure 9 Example of typical ELM coefficients (Site 4, 1998-08-25, 17:14z). The left y-
axis scales the gain (<i>m</i> , in black) and the right y-axis scales the offset and RMSE residual
(<i>b</i> and <i>err</i> , in red)
Figure 10 Example of ELM coefficients (Site 9, 1998-03-23, 14:13z) exhibiting
background vegetation features at 0.7 to 0.9 µm
Figure 11 Example of ELM coefficients (Site 5, 1995-10-21, 20:20z) exhibiting strongly
negative offset values from 0.55 to 1.8 µm
Figure 12 Generalized nonlinear correction function (solid red) and base linear function
(dashed blue)
Figure 13 ELM coefficients (left, Site 9, 1998-03-23, 14:13z) and corresponding
estimated offset spectra (right) for the four algorithms described in the text. The blue
markers on the left indicate the wavelengths corresponding to Figure 14 through Figure
16
Figure 14 Regression at 0.544 µm for linear, nonlinear, and adjusted linear methods
(plots overlay exactly at this wavelength)
Figure 15 As Figure 14 but for 0.875 µm wavelength. Right hand panel is a detail
enlargement of the lower reflectance region
Figure 16 As Figure 14 but for 1.548 µm wavelength
Figure 17 Adjusted ELM coefficients (upper left) corresponding to the unprocessed
coefficients in Figure 13. Upper right and lower panels show the linear, nonlinear, and

adjusted linear coefficients and residuals. The adjusted linear offset is identical to the
nonlinear offset
Figure 18 Modeled simulation of an "ELM" relationship. The upper plot shows the
modeled radiance points for decreasing reflectance values against a linear fit. The lower
figure shows the residuals from the linear regression
Figure 19 Mean image reflectance values derived linear ELM, modified ELM, Adjusted
ELM, and QUAC (see text) for sample images (Site 1, 1997-09-23, 17:03z (upper), Site
1, 1997-09-24, 18:57 (lower)) containing negative offsets
Figure 20 Mean image reflectance values derived linear ELM, modified ELM, Adjusted
ELM, and QUAC (see text) for sample images (Site 4, 1997-12-10, 20:42z (upper), Site
8, 1995-08-24, 13:26 (lower)) containing negative offsets
Figure 21 As Figure 20, but showing maximum reflectance values (upper) and ROI
reflectances for a ground truth target material (lower), (Site 8, 1995-08-24, 13:26 (both
panels))
Figure 22 High-level system view of RDBMS75
Figure 23 Example XML file used to populate metadata into the database
Figure 24 Entity-Relationship Diagram for the scientific RDBMS
Figure 25 Database query interface
Figure 26 Validation for the site elevation automated retrieval. The image shows the
manually downloaded digital elevation model with the site location labeled. The inset
shows the automated retrieval of the elevation
Figure 27 Validation for the AIRS integrated column water vapor (left) and water vapor
profile profiles (right) automated retrieval. The image shows the manually downloaded
data with the site location labeled. The inset shows the automated retrieval of the column
water vapor. The plots show the comparison between the automated (OPeNDAP) and
manual (FTP) retrievals
Figure 28 Variability of AIRS water vapor profiles. The left plot shows the mixing ratio
for all observations. The center plot shows the mean fractional difference between
ascending/descending node pairs. The right plot shows the mean absolute difference and
standard deviation of the pairs
Figure 29 Mean ozone burden from AIRS data for example locations plotted versus land
cover type
Figure 30 Example of MODTRAN modeled total radiance output at 1 cm-1 and 5 cm-1
band model resolution
Figure 31 Example of MODTRAN modeled total radiance output at 1 cm-1 and 5 cm-1
band model resolution, resampled to HYDICE the spectral response function. Inset shows
the detail from 1.1 to 1.5 μ m
Figure 32 Modeled gain (upper) and offset (lower) coefficients compared to
corresponding ELM coefficients
Figure 33 Mean differences between modeled and ELM gain (left) and offset (right)
coefficients, averaged over all images (/) of the site; generated using baseline scattering
input parameters
Figure 34 As in Figure 33 but with input parameter VIS=25 km
Figure 55 As in Figure 55 but with input parameter VIS=40 km

Figure 36 As in Figure 33 but with input parameter IHAZE=0;)1
Figure 37 View of the data used to select the final ELM coefficients to be used to derive	;
the standardized coefficients. Upper panels overly the empirical and modeled coefficient	ts
for comparison (gain, offset, and RMSE statistics, left to right). Lower panel shows	
ground truth and ELM notes (left) and the final ELM coefficients (right))3
Figure 38 Gain (left) and offset (right) ELM coefficients for beginning and ending times	
in image sequence) <u>4</u>
Figure 39 As in Figure 38 but after 1714z coefficients are adjusted to 1842z 10) -)5
Figure 40 Fractional errors for 90-minute time standardized coefficients trial)5)5
Figure 40 Fractional errors for 90-minute time standardized coefficients trial	כו דו
Figure 41 Oan (left) and offset (light) ELM coefficients for 5 and 10 kit mages 10	,,
Figure 42 As in Figure 41 but after 5 kit coefficients are aujusted to 10 kit altitude AOL.	י. דינ
$\Gamma_{\text{result}} = \frac{1}{2} \sum_{i=1}^{n} \frac{1}{2} $)/ \0
Figure 43 Fractional errors for 5 kft altitude standardized coefficients trial	18
Figure 44 Standardized gain (left) and offset (right) coefficients for the seasonal	
adjustment validation. The heavy black lines are the Dec reference ELM coefficients; the	e
colored lines are the standardized August coefficients)9
Figure 45 Fractional errors for Aug – Dec seasonal standardized coefficient trial 10)9
Figure 46 Site 4 gain (left) and offset (right) coefficients standardized to global reference	e
coordinates (JJA) 11	11
Figure 47 Site 4 mean standard gain (upper) and offset (lower) coefficients (JJA).	
Standard deviation is plotted in black11	12
Figure 48 Site 3 standardized gain coefficients (JJA) 11	13
Figure 49 Site 4 mean standard gain (upper) and offset (lower) coefficients (JJA).	
Standard deviation is plotted in black11	4
Figure 50 Site 9 mean standard gain coefficients ensemble distribution (left) and mean	
(right) (MAM). Standard deviation is plotted in black	15
Figure 51 Site 9 clustered mean standard gain (upper) and offset (lower) coefficients	
(JJA). Dashed lines show the standard deviation of the clusters	16
Figure 52 Means of the standard gain (upper) and offset (lower) coefficients for all sites	
and imaging events	17
Figure 53 Clusters of all standard gain coefficients, with cluster means in black	18
Figure 54 Clusters of all standard offset coefficients, with cluster means in black 11	19
Figure 55 Clustered gain means from Figure 53 (black) with site means in color	19
Figure 56 Clustered offset means from Figure 54 (black) with site means in color 12	20
Figure 57 Aggregate site membership of the k-means coefficient clusters. Green indicate	25
the primary site orange indicates a site with similar characteristics and other colors	20
indicate dissimilar sites. Legend refers to the site number)1
Figure 58 As in Figure 57 but with cluster members manned to Können _Trawartha	~ 1
climate categories	$, \gamma$
Figure 50 Similarity metrics SAM and ED for site mean gain and offset coefficients	
notted against each other to show relationshing. Colors indicate the sites' elimete	
produce against each other to show relationships. Colors indicate the sites cliffiate	<i>۱</i> (
classification categories	<u>_</u> 4

Figure 60 Similarity metrics SAM and ED for site mean gain and offset coefficients
plotted against each other to show relationships. Colors indicate the sites' land cover
classification categories
Figure 61 Similarity metrics SAM and ED for all standard gain and offset coefficients
plotted against each other to show relationships. Colors/symbols represent differing sites
and collections
Figure 62 Full set of standard coefficients plotted as ED_m vs. SAM_m . Upper left
colors/symbols represent sites and collection groups of source coefficients. Other plots
show clustering with increasing numbers of clusters; colors/symbols represent cluster
membership
Figure 63 As in Figure 62 except that the metrics platted are SAM _b vs. ED _m
Figure 64 Aggregate site membership of the ED_m vs. SAM_m coefficient clusters. Green
indicates the primary site, orange indicates a site with similar characteristics, and other
colors indicate dissimilar sites. Legend refers to the site number
Figure 65 Aggregate climate category membership of the ED _m vs. SAM _m coefficient
clusters. Green indicates the primary climate category, orange indicates a category with
similar characteristics, and other colors indicate dissimilar climate categories
Figure 66 Aggregate site membership of the SAM _b vs. ED _m coefficient clusters. Green
indicates the primary site, orange indicates a site with similar characteristics, and other
colors indicate dissimilar sites. Legend refers to the site number
Figure 67 Aggregate climate category membership of the SAM _b vs. ED _m coefficient
clusters. Green indicates the primary climate category, orange indicates a category with
similar characteristics, and other colors indicate dissimilar climate categories
Figure 68 As in Figure 66, but for 11-cluster case
Figure 69 Result of the cross validation of the k-means clusters, from n=3 to n=20
clusters, showing the percentage of samples correctly matched to site classes using the
SAM similarity metric
Figure 70 Ground truth reflectance spectra for the calibration panels
Figure 71 ELM reflectances of the three panels, obtained by performing ELM using
library reference spectra
Figure 72 Reflectance spectra of the panels retrieved by using the climatology
standardized coefficients
Figure 73 Reflectances of the three panels from QUAC144
Figure 74 FLAASH derived reflectances for the four panels145
Figure 75 Summary of SAM similarities to the ground truth reflectance spectra for the
4%, 15%, 40%, and 60% panels (left to right) 146
Figure 76 Summary of ED similarities to the ground truth reflectance spectra for the 4%,
15%, 40%, and 60% panels (left to right)

LIST OF ABBREVIATIONS

Above ground level	AGL
Active Server Page	ASP
Atmosphere Removal Algorithm	ATREM
Bi-directional reflectance distribution function	BRDF
Empirical Line Method	ELM
Fast Line-of-sight Atmospheric Analysis of Spectral Hypercubes	FLAASH
Full width at half maximum	FWHM
Grid Analysis and Display System	GrADS
Ground sample distance	GSD
Hyperspectral Digital Imagery Collection Experiment	HYDICE
Hyperspectral imagery	HSI
Long-wave infrared	LWIR
Mid-wave infrared	MWIR
Moderate Resolution Transmittance	MODTRAN
Multispectral imagery	MSI
National Land Cover Database	NLCD
Open Geospatial Consortium	OGC
Open-source Project for a Network Data Access Protocol	OPeNDAP
Physics-based	PB
Quick Atmospheric Correction	QUAC
Region of interest	ROI
Relational database management system	RDBMS
Root mean squared error	RMSE
Second Simulation of the Satellite Signal in the Solar Spectrum	6S
Short-wave infrared	SWIR
Signal-to-noise ratio	SNR
Structured Query Language	SQL
Top of atmosphere	TOA
Uniform Resource Locator	URL
Visible/near infrared	VNIR

ABSTRACT

A SPECTRAL CLIMATOLOGY FOR ATMOSPHERIC COMPENSATION OF HYPERSPECTRAL IMAGERY

John H. Powell, Ph.D. George Mason University, 2015 Dissertation Director: Dr. Kirk Borne

Most Earth observation hyperspectral imagery (HSI) detection and identification algorithms depend critically upon a robust atmospheric compensation capability to correct for the effects of the atmosphere on the radiance signal. Atmospheric compensation methods typically perform optimally when ancillary ground truth data are available, e.g., high fidelity in situ radiometric observations or atmospheric profile measurements. When ground truth is incomplete or not available, additional assumptions must be made to perform the compensation. Meteorological climatologies are available to provide climatological norms for input into the radiative transfer models; however no such climatologies exist for empirical methods. The success of atmospheric compensation methods such as the empirical line method suggests that remotely sensed HSI scenes contain comprehensive sets of atmospheric state information within the spectral data itself. It is argued that large collections of empirically-derived atmospheric coefficients collected over a range of climatic and atmospheric conditions comprise a resource that can be applied to prospective atmospheric compensation problems. This research introduces a new climatological approach to atmospheric compensation in which empirically derived spectral information, rather than sensible atmospheric state variables, is the fundamental datum. An experimental archive of airborne HSI data is mined for representative atmospheric compensation coefficients, which are assembled in a scientific database of spectral observations and modeled data. The empirical techniques for extracting the coefficients and correcting for small nonlinear features, the modeling methods used to standardize the coefficients across varying collection and illumination geometries, and the resulting comparisons of adjusted coefficients are presented. The resulting climatological database is analyzed to show that common spectral similarity metrics can be used to separate the climatological classes to a degree of detail commensurate with the modest size and range of the imaging conditions comprising the study. The study closes with a notional application example and a discussion of the potential benefits, shortfalls and future work to fully develop the new technique.

1. INTRODUCTION AND LITERATURE REVIEW

Most Earth remote sensing applications are designed to characterize the Earth's surface by analyzing reflected or emitted electromagnetic energy detected from a satellite or aircraft borne instrument. The features of interest are separated from the sensor by the Earth's atmosphere, through which all energy must pass before being measured. The energy reaching the sensor, referred to as "at-aperture radiance," is the result of a complex series of interactions between the incident solar radiation, atmospheric constituents, and material surfaces contained in the scene. Interactions include absorption, scattering, and emission, and affect not only the direct path illumination, but also include contributions from indirect path, multiply scattered, and emissive radiance as well. The effects are dependent upon the solar and observation geometry, optical path length, scene content, and atmospheric conditions. All of these can vary greatly, but the atmospheric conditions are often the most difficult to quantify. While most quantitative remote sensing applications must account for the atmospheric effects to some degree, the requirement for rigorous compensation varies by discipline. The treatment of these effects is referred to generally as atmospheric compensation, although the application can involve non-atmospheric interactions as well.

The fundamental problem addressed in this research is atmospheric compensation in an Earth remote sensing context. For an imaging sensor at some altitude above the

Earth's surface, atmospheric compensation is the process of deriving the surface reflectance values from the at-aperture radiance images recorded by the sensor. The magnitude of the atmospheric effects on measured electromagnetic energy can be strongly wavelength dependent, varying across the absorption regions of water vapor, and major and trace gas constituents in the atmosphere. Scattering by molecules and suspended aerosol particles is also wavelength dependent. Successful analysis of remotely sensed hyperspectral imagery (HSI) is particularly dependent upon a robust atmospheric compensation capability. HSI applications rely on precise relationships between spectral bands and virtually any quantitative HSI analysis must therefore begin with an inversion problem to derive the surface reflectance or emittance from the measured at-aperture radiance.

1.1 Hyperspectral remote sensing

Earth observation sensors can be classified as sounders or imagers. Sounders measure incoming radiance from individual, non-contiguous points on the ground, whereas imagers measure a two-dimensional array of adjacent ground sample points, or pixels. Imagers can further be categorized by the sensor's spatial, spectral, temporal, and radiometric resolution. Spectral remote sensing refers to the class of sensors that measure radiance in multiple wavelength bands at each ground point. Multispectral sensors have several to tens of wavelength bands across the visible and infrared regions of the spectrum. Examples of current multispectral instruments are the Landsat Thematic Mapper (seven spectral bands) [1] and the Moderate Resolution Imaging Spectroradiometer (MODIS) (36 spectral bands) [2].

Hyperspectral imagery is the sub discipline of spectral remote sensing that employs very high spectral resolution imaging spectroradiometers to produce essentially continuous spectra for each pixel. (The use of the term "pixel" to refer to a sample point containing many spectral observations is technically imprecise, but is widely used in the literature and can be used in an HSI context without confusion.) HSI sensors have much greater spectral resolution and more spectral bands than multispectral sensors. There is no accepted precise definition to separate HSI, but HSI sensors typically have spectral bandwidths of no more than 10 nm, and have a sufficient number of spectral bands (usually 100 or more) to produce a contiguous spectrum over the range of the sensor [3]. The salient feature of HSI is that it is possible to derive a complete reflectance or emissivity spectrum from each pixel [4]. Figure 1 illustrates the difference between multispectral imagery (MSI) and HSI spectral measurements. Figure 2 shows a schematic diagram of an HSI collector. The HSI dataset is a three-dimensional data cube, in which the third dimension is the spectral component. The x- and y-axes form an image in each spectral band, and each pixel contains a complete spectrum over the spectral range of the instrument. Thus each pixel provides spectral information related to the materials in the field of view, including the reflecting/emitting surface as well as the atmospheric constituents through which the light passes.

Hyperspectral sensors are in use over the range of the electromagnetic spectrum from 0.4 to 14 mm wavelengths, but a single focal plane/instrument is usually limited to a subset of wavelengths – visible/near infrared/short-wave infrared (VNIR/SWIR, 0.4-2.4 mm), mid-wave infrared (MWIR, 3-5 μ m), or long-wave infrared (LWIR, 8-14 μ m).



Figure 1 Alunite spectrum as measured by multi- and hyperspectral instruments. Source: Clark [5].

The principal reason for such high spectral resolution is the ability to perform spectroscopic analysis of each pixel or group of pixels. Through comparison with known spectral signatures of materials, usually measured in a laboratory, materials can be identified with remotely sensed imagery. Identification is based not upon the size, shape, or appearance of spatially resolved features, nor is it limited to broad classes of materials such as vegetation, road, urban, etc. Rather it is derived from the spectroscopic identification of material properties and includes identification of specific mineral and chemical composition to separate different types of vegetation, roofing materials, and paints, for example. The applications of HSI are many and include geology and mineralogy, plant ecology, environmental hazards, crop assessment, and military applications. Other HSI applications are concerned not with Earth surface properties, but with the properties of the atmosphere/water column such as climate studies, air quality monitoring, ocean color, and littoral water composition and bathymetry [3].



Figure 2 HSI schematic diagram.

Many airborne HSI instruments have been developed and operated by government, academic, and commercial institutions over the last two decades. The most widely used datasets in the scientific literature are from the Airborne Visible-Infrared Imaging Spectrometer (AVIRIS) and the Hyperspectral Digital Imagery Collection Experiment (HYDICE) sensors. AVIRIS acquires 224 spectral bands with a bandwidth of approximate 0.01µm in the VNIR/SWIR range [6]. HYDICE measures 210 spectral bands with a bandwidth of approximately 0.01µm in the VNIR/SWIR range [7]. Other instruments such as the Spatially Enhanced Broadband Array Spectrograph System (SEBASS) operate in the MWIR/LWIR region [8].

HSI instruments are classified by the type of scanner employed. Whiskbroom sensors such as AVIRIS measure a single pixel at a time and optically scan along the cross track direction. The along track direction is surveyed by the forward motion of the observing platform. Pushbroom sensors such as HYDICE use a two-dimensional detector array (one spatial and one spectral dimension) to collect one line at a time – in all spectral bands – in the cross track dimension; here, too, the along track direction is measured by the forward motion of the platform. A third class of sensors called staring sensors collects the full two-dimensional spatial scene, in all spectral bands, simultaneously [9].

The spatial resolution of HSI sensors varies depending on the optics and altitude flown. The HYDICE sensor was typically flown at relatively low altitudes of 1.5-8 km, producing a ground sample distance (GSD) of 1-4 m. Conversely, AVIRIS often flies on NASA ER-2 aircraft at altitudes of 20 km, producing a GSD of 20 m. The high data volume of HSI images and signal-to-noise ratio (SNR) concerns make placing HSI sensors on satellite platforms technically daunting. Two Earth surface imaging HSI instruments are currently in operation. The Hyperion sensor on NASA's Earth Observer-1 satellite collects 220 spectral bands in the VNIR/SWIR at a GSD of 30 m [10]. The

Hyperspectral Imager for the Coastal Ocean (HICO) was deployed on the International Space Station in 2009. HICO collects 102 spectral bands between 0.38 and 0.96 µm at a GSD of 92 m at nadir [11].

The sensitivity of hyperspectral remote sensing to atmospheric constituents is used to advantage by atmospheric scientists. NASA's Atmospheric Infrared Sounder (AIRS) is a satellite-borne hyperspectral instrument designed to measure accurate temperature and water vapor vertical profiles. It uses 2378 spectral bands in the infrared region of the spectrum. The very high spectral resolution in the water absorption bands coupled with temperature and pressure effects on emission allow a radiative transfer inversion to infer the temperature and water vapor vertical structure. The instrument accuracy is ± 1 degree K in 1 km layers for temperature, and $\pm 20\%$ in 2 km layers for water vapor. That corresponds to 24 standard pressure level values for temperature, and 12 for water vapor. The spatial resolution of each sounding is 13 km at nadir to 40 km at the swath edge. The sensor also measures trace gas concentrations, and NASA currently provides ozone, carbon monoxide, and methane concentrations. Other trace gas retrievals are being researched but are not yet operational [12].

1.2 Atmospheric Compensation

The fundamental problem addressed in this research is atmospheric compensation in an Earth remote sensing context. For an imaging sensor at some altitude above the Earth's surface, atmospheric compensation is the process of deriving the surface reflectance or emittance values from the at-aperture radiance images recorded by the sensor. The magnitude of the atmospheric effects on measured electromagnetic energy

can be strongly wavelength dependent, varying across the absorption regions of water vapor and the major and trace gas constituents in the atmosphere. Scattering by molecules and suspended aerosol particles is also wavelength dependent. Figure 3 shows the top of atmosphere (TOA) solar illumination spectrum against the illumination at the Earth's surface for typical atmospheric conditions [13].



Figure 3 Solar illumination at TOA and Earth surface for typical atmospheric conditions. Major absorption bands are labeled. Source: Rhode [13].

However, some remote sensing applications require no compensation for atmospheric effects whatsoever. Broadband "panchromatic" images, which depict a single broad spectral band, are used to visually interpret features and the variation of the signal across the observed spectrum has little effect on their utility. Similarly, MSI spectral bands are relatively broad and often placed in atmospheric "window" regions where atmospheric effects are small – many applications, even band ratio algorithms such as Normalized Differential Vegetation Index (NVDI) are relatively insensitive to atmospheric effects. HSI applications, however, do rely on precise magnitudes of spectral bands. In particular, the ability to perform spectral matching with laboratory spectra for material identification depends critically on the ability to compensate for atmospheric path radiance and transmissivity. The application of atmospheric compensation considered here is for HSI, so high fidelity compensation for detailed spectral analysis is required.

As shown in Figure 3, the TOA solar spectrum is, to first order, a blackbody emission curve at the temperature of the sun's surface. The figure shows the attenuation effect of atmospheric gases across the VNIR/SWIR spectral region; note that water vapor is a strong, prevalent absorber. This single-path absorption is only one effect that must be considered, however. Multiple paths of electromagnetic propagation are detected by the sensor. The significant reflective radiation interactions of interest are shown schematically in Figure 4. The largest component of the measured signal comes from single path radiance (path A in the figure), which propagates through the atmosphere, reflects off the surface in the sensor's instantaneous field of view (IFOV) or pixel, and propagates directly back to the sensor. Radiation also scatters from interaction with molecules (Raleigh scattering) and aerosol particulates (Mie scattering). Some of this radiance is scattered back directly to the sensor (path B), or after reflecting off the surface target (path C). Lastly, radiance reflecting off nearby surface areas outside the sensor's FOV can be scattered back to the sensor (path D) or, in complex terrain, through multiple surface reflections (path E). These are referred to as adjacency effects. Even more complex paths occur through multiple scatters, but the magnitude of the signal is reduced with each scattering event and these paths can often be neglected.



Figure 4 Schematic diagram of reflected and scattered radiance paths detected by an airborne sensor. The ellipse represents the sensor's IFOV.

Atmospheric transmission along each of the paths is subject to wavelength dependent absorption losses. The nature of the atmospheric scattering is determined by the wavelength of the radiation relative to the particle size, and is thus highly wavelength dependent. The scattering cross section is also directional, especially in the Raleigh scattering regime, causing the atmospherically scattered paths (B, C, and D) to have different spectra than the direct path illumination.

The atmosphere and terrestrial objects also emit radiation in the mid- and longwave infrared that reaches the sensor (emissive paths are not shown in Figure 4). For nominal terrestrial temperatures (~300K), the wavelength of peak emission given by Wien's law is ~10 μ m, compared to ~0.5 μ m for solar temperatures. Thus, the problem of atmospheric compensation can be effectively separated into emissive and reflective effects over limited wavelength ranges of interest. For the VNIR/SWIR spectral ranges considered in this research, emissive radiance is negligibly small and is neglected in this work (as it is by VNIR/SWIR HSI practitioners). The problem is then reduced to deriving surface reflectance values from the measured at-aperture radiance values.

The body of processing techniques used to separate and account for the multiple atmospheric paths comprises the discipline of atmospheric compensation. Having compensated for atmospheric effects, the nature of the surface interactions with light can then be studied. Many methods of atmospheric compensation are used in the HSI analysis community, and many more have been developed in the literature. The following sections provide a survey of these techniques. Section 3.1 derives in detail the approach used in this research.

1.2.1 Empirical Methods

HSI datasets contain complete sets of spectral measurements of light passing through the atmosphere at each pixel; therefore, the information about the atmospheric transmission is present in the measured radiance signal. Empirical methods use this information along with some additional information about the scene to statistically derive the relationship between radiance and reflectance. It is almost always assumed that a linear relationship exists for each observed wavelength. The set of multiplicative (gain or slope) and additive (offset or intercept) coefficients defines the empirical relationship. The additional scene information can range from an assumption about the statistics of the scene to detailed ground-truth in situ spectral measurements.

The simplest approaches are purely statistical and are sometimes referred to as image-based methods. The Internal Average Relative Reflectance (IARR) method normalizes each pixel by the scene average spectrum. The assumption about the scene is that it is large enough and spectrally diverse enough so that the mean spectrum is representative of only the atmosphere and relatively featureless. Gao, Davis and Goetz [14] describe IARR and a similar method called the Flat Field Correction (FFC), which normalizes the scene based on the average spectrum of a specific area in the image assumed to have a flat, featureless spectrum [15, 16]; it is essentially a one-point empirical line method. Alternatively, the method of log residuals described by Green and Craig [17] normalizes the scene by dividing by the geometric mean taken over all wavelengths. Like IARR, the method is dependent on the spectral diversity of the scene. These methods do not require any ground-truth information, but are valid only for scenes having the assumed spectral characteristics. They also return only relative reflectances, and do not correct for the average offset, dominated by atmospheric scattering. The dark object subtraction (DOS) method evaluated by Campbell [18] accounts for this term by

subtracting global (or defined regional) radiance minima across the spectrum. Crippen's regression intersection method (RIM) [19] is a scene based, empirical technique that estimates the absolute path radiance term. In RIM, spectrally contrasting pixels of homogeneous areas are selected and bispectral regression lines are projected to intersection points, which are assumed to represent zero reflectance values. Several pairs of points are used to determine the path radiance.

To obtain the best absolute reflectance values requires more information from the scene, specifically the spectra of one or more objects in the scene. With this information, exact coefficients can be obtained by regression over the known targets, and applied over all the pixels in the scene. This technique is called the Empirical Line Method (ELM), first developed in detail by Conel [20] and validated by numerous studies of spectral remote sensing [21, 22]. With only one reference spectrum, the regression line is assumed to pass through the origin, representing a theoretical zero reflectance pixel, but this does not account for the atmospherically scattered radiance (path B in Figure 4). At least two points are normally used, with widely varying reflectance values (one bright, one dark) to obtain the most accurate regression line. Smith and Milton [22] showed that three or more reference spectra further reduces the error in ELM retrievals. This could in theory account for nonlinearity in the retrieval coefficients, but the improvement is more likely due to a more accurate determination of the linear relationship. Several studies have demonstrated the linearity of the ELM method in the VNIR/SWIR range [22]. Baugh and Groeneveld [23] studied ELM utilizing Landsat Thematic Mapper and over 2600 ground truth spectra to show linear regression coefficients of determination of 0.90-0.99.

ELM accounts for the transmission losses and path radiance effects (paths A, B and C in Figure 1). The method will incorporate local adjacency effects in the reference spectra, but cannot account for non-local adjacency effects. For this reason, large reference targets are needed so that edge pixels that may be mixed or contaminated by adjacent pixels are minimized. Further theoretical assumptions inherent in the method are that the illumination and atmospheric effects are constant across the image. Thus, broken cloud and topographic shadowing is ignored [7]. While trace gases and aerosols are likely nearly homogeneous across a scene, water vapor can vary substantially on scales of several kilometers [15]. Finally, variations in optical path due to varying topographic elevations across a scene are not treated.

One advantage of ELM is that it is resilient to calibration errors as the method inherently corrects for any radiance errors and sensor artifacts so long as they are uniform across the scene. The biggest limitation of the method is the requirement for ground truth spectra. Image studies typically deploy large calibration panels or use well calibrated ground measurements for the best results. In applications where no in situ measurements are available, ELM can still be used, perhaps with reduced accuracy. If some materials can be identified visually in a scene (asphalt, concrete, dry lakebeds, etc.), they can be regressed against associated laboratory reflectance spectra to determine the gain and offset terms. On the plus side, no external information about the atmosphere is required, although some studies have proposed improving on ELM results by adding data from additional ground sensors [24]. From a practical standpoint, ELM is simple and

computationally trivial, but it has not yet been fully automated. Currently a human analyst is required to select the reference target pixels in the image.

1.2.2 Physics-Based Methods

The physics of radiative transfer is well understood and can be accurately modeled using radiative transfer algorithms. Physics-based (PB) methods use radiative transfer codes to estimate the atmospheric effects on transmission and determine the surface reflectivity from the model. No in-scene spectral ground truth is required, but detailed knowledge of the illumination and atmospheric conditions is required to accurately model the scene. Some of this information, such as atmospheric water vapor, can be obtained from the hyperspectral image data itself. Other necessary information, such as the illumination geometry and aerosol characteristics must be either supplied externally or calculated from other data. With a detailed radiative transfer model and sufficient scene information (e.g., topography, building layout, ground-truth spectra, etc.), nonlinear effects such as the adjacency effects (paths D and E in Figure 4) can, theoretically, be modeled.

All PB approaches have at their core the algorithms that describe the transmission and scattering effects occurring along the optical path. The great majority of current PB methods use either the Moderate Resolution Transmittance (MODTRAN) developed by Berk et al [25, 26] or the Second Simulation of the Satellite Signal in the Solar Spectrum (6S) [27] code to perform the radiative transfer calculations. The models differ in their molecular band model parameterizations, aerosol and scattering models, but each produce comparable solutions over a range of atmospheric conditions. Both models are the

subjects of continuing research and incremental improvements are made with periodic releases.

Many atmospheric compensation programs have been created based on either MODTRAN or 6S radiative transfer algorithms. These include Atmosphere Removal Algorithm (ATREM) [14], Fast Line-of-sight Atmospheric Analysis of Spectral Hypercubes (FLAASH) [28, 29], MODTRAN Full (MODFULL) [30], Atmospheric Correction Code (ACC) [31] and Atmospheric Correction Now (ACORN) [14]. Of these, FLAASH and ATREM are the most prominent in the literature owing to their availability and widespread use (ATREM is no longer supported). Each of the programs has unique strengths and shows superior performance in certain applications. FLAASH and ACC offer the ability to estimate visibility and aerosol loading from the HSI data, correction for the adjacency effect, and spectral polishing, a post processing technique that aims to improve the returned reflectance spectra. MODFULL employs a statistical method to select the best pixels in the scene to use for water vapor and aerosol retrieval, and can be executed entirely from the scene without additional user input. The PB methods generally use a variation of the three-band ratio technique of Gao et al. [32, 33, 34] to determine water vapor. Aerosol characteristics and methods to determine aerosol loading vary more appreciably between the models. The High-accuracy Atmospheric Correction for Hyperspectral Data (HATCH) [35] includes several unique features, including a custommade radiative transfer algorithm and a method to derive water vapor from differences in the spectral smoothness of the atmospheric component. Another method, Atmospheric Correction via Simulated Annealing (ACSA), forms a constrained optimization problem

to solve a spectral smoothness criterion designed to reduce residual effects on reflectance values retrieved from a PB model in low SNR environments [36].

PB methods are computationally intensive. While increases in computer processing speed have made it possible to run the programs on desktop computers, lookup tables, reduced resolution computations, and modified calculations for multiple scattering are still employed in order to speed up the execution as illustrated by Cairns et al [37]. Even with trade-offs, PB models often calculate the most important factors, like water vapor and transmission losses, on a pixel-by-pixel basis. Given a detailed enough topographic elevation model, topographic effects on optical path length and illumination can be accounted for in the models. In theory, full bi-directional reflectance effects could be modeled, but these have yet to be demonstrated. Some approaches perform a large set of PB model runs to cover the expected range of conditions in advance. The results are stored in look-up tables that are accessed during the scene compensation [38, 39]. The look-up tables are often very large, and can take days of computation time to generate, but allow rapid execution of later scene compensations.

The ability to compensate for atmospheric and topographic variability across a scene is a powerful feature of the PB approach. However, the results suffer when imperfect atmospheric conditions are modeled. The water vapor retrieval schemes have been shown to produce accurate results over many conditions (but not all, for instance over low reflectance backgrounds like water), but aerosol retrieval is still an evolving science and the specification of aerosols is a weakness of the models [40, 41, 42]. Climatological atmospheric values are often used when there are no in situ measurements

of water vapor or aerosol content and in-scene retrievals are not suitable. These climatologies are typically not very sophisticated; the FLAASH climatological default has only six climate regimes, based on latitude and season [43].

PB approaches also demand highly accurate radiometric calibration, since any error in the radiance values will propagate to the calculated reflectance values. Matteoli, Ientilucci, and Kerekes [44] compared results in target detection in real and modeled HSI scenes using Digital Imaging Remote Sensing Image Generation (DIRSIG [45]), FLAASH and ATREM. The results indicated the need for a better understanding of the major sources of uncertainty in PB approaches. From a diagnostic point of view, PB models can help to diagnose atmospheric transmission by decomposing their underlying model components, whereas empirical methods may not fully separate the contributions of the various constituents. However, with proper assumptions, it is possible to relate empirical coefficients to physical radiative transfer parameters. Conel [20] developed the physics formulation of ELM that is used and expanded upon in this research.

1.2.3 Alternative Methods

Other methods of atmospheric compensation do not strictly fit in the categories discussed so far. The first is the set of hybrid methods that employ PB modeling along with empirical techniques. For example, ELM can be used at a single calibration site, then the results extended over varying topography using a PB model [14]. Spectral polishing techniques such as the Empirical Flat Field Optimal Reflectance Transformation (EFFORT) applied to the results of PB models are also examples of the hybrid approach [15].

A number of emerging methods use statistical methods to derive surface reflectance solely from the data in the scene. These are sometimes called invariant methods because they are designed to work without explicitly defining the illumination and atmospheric conditions. These include the coupled subspace model of Chandra and Healey [46] and the Quick Atmospheric Correction (QUAC) model of Bernstein et al. [47, 48, 49]. QUAC assumes a linear radiative transport equation like ELM, but uses a ratio of scene-derived statistics to reference scene statistics to calculate the gain coefficients. The reference scene is a spectrally diverse collection of laboratory reflectance measurements. These methods show results comparable to PB models under many conditions but without the requirement for exact atmospheric and illumination information. They also involve less calculation than PB methods and are not dependent on accurately calibrated radiance data. Other work suggests the use of covariance matrix statistics, which are frequently used in HSI noise analysis and detection algorithms, to separate atmospheric effects from the reflectance signal. The Covariance Matrix Method (CMM) uses statistics to estimate the path radiance contribution [50]. Monte Carlo methods have been developed to analyze the effects of atmospheric features in detail and to validate PB codes under a range of conditions as in the study by Nardino et al[51].

An important class of analytical methods is those that use information from prior analyses in addition to the current dataset. Atmospheric compensation methods, with few exceptions, use only the current information in the scene. The only prior information routinely used is historical climatologies of sensible atmospheric quantities (temperature, water vapor, aerosol concentration, etc.). An exception is the Surface Prior-Information

Reflectance Estimation (SPIRE) algorithms of Viggh and Staelin [52, 53]. SPIRE uses prior reflectance information from the region imaged along with current radiance and basic statistics of the inversion gain and offset terms to estimate the current reflectance. Variants use spatial or spectral filtering, or a combination of both, to isolate the degrees of freedom of the terms and estimate the reflectance. Results compare favorably to ELM and ATREM in many cases (high clouds, haze, etc.), but the method is limited to situations in which there are multi-temporal collections and the change in reflectance from the prior image is small.

It is worth noting that one option is to forego atmospheric compensation entirely and conduct the analysis in radiance units. For some algorithms such as anomaly detection this approach may be suitable, but it is much less effective for detailed spectral analyses such as material identification algorithms. Yuen and Bishop [54] demonstrated that high fidelity atmospheric compensation improved the spectral contrast in the HSI scenes by 20%, resulting in classifications that were markedly (up to a factor of ten times) more accurate than those that were based upon the radiance image. Another option is to determine an estimated atmosphere for the scene, convert the library spectra of interest into radiance, and apply the spectral matching algorithms in radiance space. These methods are advocated from a computational standpoint when using PB models, as far fewer computations are required compared to running the model at each pixel. However, nearly the same computational benefit could be achieved with equal accuracy if one or several pixels were modeled numerically, and the resulting transformation applied across the rest of the scene. Either choice opens the possibility of computing an ensemble of atmospheres and selecting either the best match to observed conditions or some statistical derivative of the ensemble. This approach assumes constant illumination and atmospheric conditions across the scene.

Experienced spectral analysts rarely rely on a single method for atmospheric compensation. In detailed analyses, several methods can be performed hierarchically to arrive at the best result. For example, an automated PB or statistical routine can provide a "first guess" compensation to allow a material identification. That result can then be used to perform an ELM compensation with more accurate results than would be obtained using either method alone.

1.3 Comparisons

Each of the approaches and individual algorithms discussed above has strengths over the other methods in certain applications – specific classes of scene or atmospheric conditions. When comparisons to other methods are made in the literature, they typically compare direct quantitative results from one or two other methods on one to several scenes. Broader comparisons are confined, qualitatively, to the general strengths and weaknesses that have been summarized in the previous sections.

One of the more extensive comparisons in the literature was conducted by Stewart, Bauer, and Kaiser [55]. They compared five different algorithms, including both empirical and PB methods, over 15 HYDICE datasets. The data included scenes representative of five differing environments: desert, jungle, alpine, littoral, and forest. The study found that overall, ELM usually performed the best (calibration panels were present in the scenes). ELM was less effective in scenes that had large spatial or temporal
atmospheric variations. Of the PB approaches tested, ACC performed the best overall. It was noted that spectral polishing post processing did not always improve the results. A few studies include detailed error analysis. Kerekes [56] performed an error analysis on ELM and ATREM reflectance retrievals from two HYDICE datasets. He found random errors of 1-2% and systematic bias errors of 1-4% in the retrieved reflectance data compared to in situ radiometric observations. ELM errors were smaller than those of ATREM, but showed large errors near water vapor absorption regions when the water vapor concentration varied significantly between reference (calibration panels were in the scene) and target areas. The PB modeling approach was shown to be considerably more sensitive to sensor noise and calibration errors. These results are consistent with the errors quoted in other sources; under the best conditions, retrievals accurate to 2-4% are commonly reported.

One difficulty in extensive comparisons of the atmospheric compensation methods is that the algorithms are constantly improving, so results begin to lose value after only a few years. It is also difficult to obtain large numbers of HSI scenes suitable for such a study. No recent, large-scale comparisons involving the latest PB and alternative methods have been published. HSI atmospheric compensation remains somewhat ad hoc, with researchers and analysts using the methods they are most comfortable with based on the available tools and the characteristics of the scene. The best compensations are obtained through trials of variety of methods to find the best method for the individual dataset. Still, the general characteristics of the established methods are well understood, especially the empirical methods, which are relatively

stable. When high quality ground truth is available and the appropriate assumptions are valid, ELM remains the gold standard for atmospheric compensation, returning the best reflectance values over the range of applicability.

1.4 Environmental Databases

Scientific databases have demonstrated the power to organize and facilitate knowledge mining from vast scientific datasets. Kruger et al. [57] show how a relational database and distributed web data server can organize a multi-terabyte collection of Next-Generation Radar (NEXRAD) meteorological radar data. The system frees researchers from tedious data management tasks, allowing them to focus on the research and work more efficiently. When enabled with data collections from other scientific disciplines, new synergies can emerge. The astronomy community has established a set of data and metadata standards that are extensive enough to enable a "Virtual Observatory", in which numerous projects across the globe have linked their astronomical archives and databases together with analysis tools and computational services to effectively function as a single archive. The success of these efforts is documented in the research literature [58,59]. Data intensive, technology enabled scientific collaborations such as these are referred to as e-Science [60].

Most scientific communities, however, produce and distribute scientific data in numerous and often nonstandard formats. While advances in storage, communications and Internet technologies have enabled access to large data stores online, the standardization of data formats has lagged. The need to have specific software or to write

custom computer routines to access and analyze data from each source has been a significant obstacle to data sharing across disparate user communities. For singlediscipline scientific applications, one may simply adopt the format of the dominant data provider. For applications requiring a variety of data types and sources, however, the problem of multiple data formats and Application Program Interfaces (APIs) can become intractable.

Standard formats have been accepted and used by major providers within the Earth sciences communities for some time. Examples include the Common Data Format (CDF) [61], Hierarchical Data Format (HDF) [62], and Gridded Binary (GRIB) format [63]. Within a given standard, however, exist variants such as HDF4, HDF5, and HDF-EOS formats, and GRIB1 and GRIB2 formats, with limited compatibility between them. These formats are optimized for specific classes of scientific data, and are therefore well suited to disseminating complex datasets. However, they are designed to serve data as packaged by the provider; that is, there is no server-side provisioning of the data. The data package that is offered might not be optimal for a particular science application. For instance, it is not efficient to download a full granule of satellite observations to populate a small number of database fields.

Implementing e-Science requires web-based data services that allow users to retrieve the desired data across a variety of sources using a single API. Great progress has been made in specific disciplines. The geospatial community has adopted a series of standards under the purview of the Open Geospatial Consortium (OGC), including the Web Map Service (WMS), Web Coverage Service (WCS), and Web Feature Service

(WFS) standards [64]. These standards provide the backbone for community wide sharing of maps and geospatial information, and have enabled the proliferation of geospatial web service applications on the Internet. Bai et al. [65] present a taxonomy of geospatial data services and show how the classification scheme applies to an emerging network of Earth science data providers, the Global Earth Observation System of Systems (GEOSS) [66].

Providers of scientific data in other communities have been slower to implement these OGC standards, in part because of the challenges in implementing the existing WCS standard for complex grids and data structures [67]. In some cases, downloading fixed files remains the norm for scientific data exchange. In others, such as the Virtual Observatory example, community-specific standards are accepted and applied within the community to great benefit.

Powell et al. [68] showed how a relational database management system (RDBMS) can provide researchers with a single, easily accessible repository of meteorological data to aid in researching atmospheric compensation. This work is summarized in Chapter 4, which describes how the database system was adapted to support this research.

1.5 Research Objectives

Understanding the properties of the atmosphere, particularly aerosol, water vapor, and trace gas content, is key to quantitative HSI analysis. The success of empirical atmospheric compensation methods suggests that remotely sensed HSI scenes contain

comprehensive sets of atmospheric state information within the spectral data. This information is most effectively used in its native spectral form encapsulated in the observed radiance data. ELM directly uses the native spectral information in the form of gain and offset coefficients. Conversely, PB methods use techniques such as band ratios to extract the information and convert it to conventional meteorological parameters (water vapor mixing ratios, aerosol concentrations, etc.). The PB models then use radiative transfer algorithms to translate the meteorological information back into spectral effects during the reflectance inversion. Undesirable artifacts are inevitably introduced into the data with each translation between domains. When in situ reflectance measurements of ground truth targets are not present, the accuracy of ELM results also varies depending on the availability of representative natural reference targets in the scene. In operational settings, such information is rarely available, leaving the analyst to apply empirical methods using in-scene sources or PB models. Similarly, in situ meteorological information is rarely available to provide the best parameters for the PB models.

For images without ground truth or detailed meteorological information, the best information available about the atmospheric state is likely climatological data. The data, however, are compiled in terms of sensible meteorological parameters. The native spectral information about the atmosphere captured in empirical atmospheric compensations is not compiled for use outside of the scenes from which they were derived. This research develops a new paradigm for HSI atmospheric climatology, using a statistical and PB approach in which empirically derived spectral information is the

fundamental datum rather than sensible atmospheric state variables. Furthermore, large experimental HSI archives are shown to comprise an untapped resource that can be applied to a wide range of atmospheric compensation problems.

The primary contribution to the state of the science is the demonstration and use of the atmospheric information contained in the empirical coefficients, which has to date not been explored in the literature. The methodology used to standardize the coefficients is a new application of PB radiative transfer modeling. In the process of developing these methods, the nature and content of the empirical coefficients is revealed in greater depth than in previous studies, which almost without exception focus only on the direct use of the coefficients to perform a reflectance retrieval within the scene itself. Finally, a sizable collection of empirical coefficients generated over a range of climate regimes is analyzed and presented to the community.

The remainder of this document describes the methodology, results, and conclusions of the research and is organized as follows. Chapter 2 describes the HSI image and ground truth data sets that were used in the study. The imagery model and ELM procedure is described in Chapter 3. This section covers how the ELM-derived coefficients were processed and standardized for inter-comparison, and how the scientific database contributed to the research. Chapter 5 describes the development and analysis of climatological classes, and Chapter 6 reviews the notable results and conclusions of the work.

2. DATA DESCRIPTION

The hyperspectal data used in this research was collected by the HYDICE sensor between the years 1995 and 2000 over a range of climatic regions, backgrounds and seasons. Each collection was accompanied by ground truth information to characterize the scene. This section describes the sensor characteristics, imaging geometries, ground truth, and ancillary data used. It also discusses the characteristics of the imaged sites included in the research.

2.1 HYDICE Sensor

HYDICE was a pushbroom hyperspectral sensor with a spectral range of 0.4 to 2.5 micrometers (VNIR/SWIR). It used a Schmidt prism spectrometer with a single indium antimonide (InSb) focal plane. HYDICE collected 210 spectral bands with a nominal bandwidth of 10 nm and 320 spatial samples [7]. The sensor was mounted on a nadir-viewing commercial stabilization platform on a Convair CV-580 aircraft. A 0.5 mrad instantaneous field of view (IFOV) produced ground sample distances (GSD) ranging from approximately 1 m to 4 m at typical operating altitudes (5000 to 20,000 ft above ground level (AGL), respectively). The HYDICE sensor employed an onboard tungsten-halogen calibration source for in-situ calibration measurements. Error sources have been well studied and absolute radiometric uncertainty is approximately 5% [69,70].

Average band spacing is 10.0 nm and average bandwidth (full width at half maximum (FWHM)) is 13.4 nm.

The HYDICE experimental collections of the mid- to late 1990's included images of large ground calibration panels to assist in sensor characterization and atmospheric compensation, typically on each flight line. Images containing the calibration panels were selected for use this in this research. Images were collected at three standard flight levels – approximately 5,000, 10,000, and 20,000 ft AGL. A total of 181 HSI images were used in the research, 127 of which produced atmospheric compensation coefficients. Most of the remaining 54 images lacked acceptable ground truth information or metadata to permit full analysis. A few of these were corrupt or adversely affected by varying illumination conditions.

2.2 Ground Truth Measurements

The HSI images used in this study are part of a series of research collections that were meticulously ground-truthed. Specifically, the calibration panels used in the analysis were measured spectrally across the range of the HYDICE sensor to ensure an accurate in-scene reflectance target was known. Spectral measurements were taken with a highresolution field spectrometer such as the GER Mark IV to 3700 models (precursors to current SVC HR models [71]) or equivalent Analytical Spectral Devices (ASD) models [72]. The field spectrometers were calibrated daily using a field spectral radiance standard source. Three to twelve spot measurements from various locations on the panel were averaged together in each observation to account for non-uniformity. Typically the radiance of a known standard "Spectralon" calibration material was measured alternately between measurements of the target panels to derive the ratios used in the reflectance measurements. Spectralon is a commercial name for polytetrafluoroethylene (PTFE), an environment-resistant material having stable and flat spectral reflectance properties [73]. The in situ Spectralon-derived ratios provided the most accurate reflectance measurements, but in a small minority of measurements, laboratory instrument calibration factors were applied directly to the target panel radiance measurements to obtain the panel reflectance. Detailed logs were kept of each measurement, including the time and number of measurements, sky conditions and any anomalies noted. Similarly, flight logs were maintained to record the time, altitude and location of each HYDICE image. With few exceptions, the calibration panels were spectrally measured each flight day, and often at times chosen to correspond closely with the HYDICE overflights.

Other related data were collected from the imaging sites, depending on the specific experiment goals. These sometimes included weather observations, downwelling radiance, photographs and spectral measurements of other target and background materials.

Despite the detailed records collected in the original experiments, nearly 20 years has elapsed since some of the data were collected and considerable effort was required to piece together the relevant ground truth data for this research. The loss of records correlating filenames to content, obsolete digital formats and undocumented processing provenance all contributed to difficulties in assembling complete ground truth records of the collects. Ground truth problems were responsible for the bulk of images that could

not be used in the study. Incomplete ground truth data, such as measurements from a previous day or measurements of only a partial set of panels, were used only when acceptable ELM results were obtainable from the data.

2.3 Imaged Sites

Ground sites for the HYDICE experiments were selected for diversity of climate, ground cover and season. The images included in the research encompass 14 collections over nine geographic sites. The geographic locations are depicted in Figure 5. These sites include continental plain and mountainous terrain as well as littoral regions. Environments range from tropical to mid-latitude temperate to arid conditions. Ground elevations range from sea level to nearly 10,000 ft. Background land cover includes bare earth, open shrub, agricultural vegetation, forest and urban environments. Each collection event occurred over a period of two to five days, and three to five images from each day are included in this research. Imaging times are from 9 am to 3 pm local time, with the majority of the images nearer to local noon. Imaging conditions were mostly clear skies but occasionally included broken clouds or thin high clouds.



Figure 5 Geographic locations of sites. The marker in the California Sierra Nevada Mountains represents three separate sites.

Table 1 summarizes the characteristics of the sites and images used in the research; Table 2 and Table 3 define the associated climate and land cover classes. Land cover classes are based on USGS National Land Cover Database (NLCD) 2001 categorization [74]. The classification scheme was developed for remote sensing-based land classification, specifically the LANDSAT Thematic Mapper. In this research, categories are assigned by manual review of the HSI imagery in true color and infrared false color composites compared to the descriptions in the classification system. Two categories are assigned for each site. The primary category refers to the immediate

environment surrounding the calibration panels (within approximately 10 pixels). The

secondary category describes the dominant land cover of the region (within

approximately 500 m).

tonowing tac	Showing tables.						
Site	Climate Category	Elevation (ft)	Land Cover Category (pri, sec)	Seasons Imaged (MMM)	Number of Images		
1	Н	6810	31, 51	SON	19		
2	Н	8498	31, 42	SON	12		
3	Н	9754	31, 51	SON	12		
4	BW	5267	31, 51	JJA, DJF	8		
5	BW	786	31	JJA, SON, DJF	24		
6	Cs	62	23, 51	JJA	4		
7	Cf	1025	71, 82	JJA	11		
8	Dc	18	71, 41	JJA	22		
9	Ar	193	71, 41	MAM	15		

Table 1 Summary of site characteristics and imaging. Climate and land cover categories are described in the following tables.

 Table 2
 USGS Land Cover Institute (LCI) land cover class definitions [74].

Land Cover Category	Level I Class	Level II Class
23	Developed	Commercial/Industrial/Transportation
31	Barren	Bare Rock/Sand/Clay
41	Forest Upland	Deciduous Forest
42	Forest Upland	Evergreen Forest
51	Shrubland	Shrubland
71	Herbaceous Upland	Grasslands/Herbaceous
82	Herbaceous Planted/Cultivated	Row Crops

Climate Category	Climate Class	Туре	
Ar	Tropical	Rainy	
BW	Dry	Arid or desert	
Cf	Subtropical	Humid	
Cs	Subtropical	Dry-summer maritime subalpine	
Dc	Temperate	Continental	
Н	Highland	N/A	

 Table 3
 Köppen–Trewartha climate classification categories.

Figure 6 shows representative land cover images from the sites. These are recent images captured from Google Earth and are therefore not, in general, representative of the imaging conditions or seasonal growth present in the HSI images used in the study.

The modified Köppen climatic classification system as described by Trewartha [75, 76] is used to characterize the climatic zones of the sites. The Trewartha-Köppen system considers prevailing atmospheric factors (e.g., temperature, humidity, rainfall) as well as geographic and geological aspects in defining the classes. This produces climatic regions that are better related to dominant flora and soil types, and thus better related to the typical atmospheric characteristics than are strictly geographic/morphological categorization systems. The system further subdivides the climate types listed in the table, but these discriminators are considered too fine for the purposes of this study.



a. Category 31 Barren (Site 6)



c. Category 41 Forest Upland (Site 8)



b. Category 31 Shrubland (Site 5)



d. Category 71 Herbaceous Upland (Site 7)



e. Category 82 Herbaceous Cultivated (Site 7)



f. Category 23 Developed (Site 6)

Figure 6 Current images of a sample of the sites showing representative land cover categories. Images captured from Google Earth.

3. METHODOLOGY

The central argument of this study is that empirically derived reflectance inversion coefficients can be used to characterize the atmosphere. Accordingly, the derivation of the coefficients and their subsequent processing is of the utmost importance to the research. This section provides the mathematical model used to represent the observed radiance in the form of the empirical gain and offset coefficients. The methodology to derive the coefficients is presented, as well as the analysis employed to evaluate and correct the coefficients. An alternative, automated empirical method (QUAC) is described for comparison.

ELM coefficients are specifically tied to the illumination and imaging geometry present at the time of the image from which they are derived, as all the factors affecting the radiative transfer are encapsulated within the coefficients. In order to isolate the atmospheric effects and to compare coefficients amongst different images, a method to account for varying illumination and geometry is required. This methodology is described in the latter part of this section.

3.1 Remote Sensing Model

Referencing the imaging model represented in Figure 4, the radiance reaching the sensor, $L_{s}(\lambda)$, can be written as:

$$L_{S}(\lambda) = L_{dir}(\lambda) + L_{sky}(\lambda) + L_{path}(\lambda)$$
(1)

where $L_{dir}(\lambda)$ is the direct path reflected radiance (Path A in Figure 4), $L_{sky}(\lambda)$ is the indirect sky-illumination reflected-radiance (Path C) and $L_{path}(\lambda)$ is the path radiance (Path B). Adjacency and multiple surface scatter effects are neglected, as are thermal emissive radiance contributions. The direct path term is given by:

$$L_{dir}(\lambda) = E_0 \tau_d(\lambda) \tau_u(\lambda) \frac{\rho(\lambda)}{\pi} \cos \sigma$$
(2)

where E_0 is the solar irradiance at the top of the atmosphere, $\tau_d(\lambda)$ is the downward path transmittance, $\tau_u(\lambda)$ is the upward path transmittance, $\rho(\lambda)$ is the surface reflectance factor, and σ is the incident angle to the surface. Implicit in the reflectance factor term is the assumption of a Lambertian surface. In the more general case, $\rho(\lambda)$ would be replaced by the bi-directional reflectance distribution function (BRDF). The indirect reflected term is written as:

$$L_{sky}(\lambda) = E_s \tau_u(\lambda) \frac{\rho(\lambda)}{\pi}$$
(3)

where E_s is the skylight irradiance at the surface. Here it is assumed that the entire hemisphere of the sky is visible to the surface and again, the surface is Lambertian.

Equation 1 can then be written as a linear relationship (dropping the wavelength dependence notation for clarity):

$$L_{\rm S} = m\rho + b \tag{4}$$

where *m* and *b* are the gain and offset vectors given by:

$$m = \left(E_0 \tau_d \cos \sigma + E_s\right) \frac{\tau_u}{\pi} \tag{5}$$

$$b = L_{path}.$$
 (6)

In the ELM model, the direct path and sky illumination reflected radiance represented by the gain coefficient *m* and the path radiance is represented by the offset *b* in equation 4. The coefficients *m* and *b* are assumed constant across the image, and are therefore one-dimensional vectors in wavelength space. They are determined empirically by selecting two or more groups of pixels for which the reflectance values are known (or assumed known) and performing a linear regression of the measured radiance against the ground truth reflectance. These vectors are then applied against each pixel to estimate its reflectance. When adjacency and multiple scattering effects are non-negligible, they are absorbed into the coefficients. To the extent that these effects are relatively consistent across the image, the reflectance values obtained by ELM will not suffer; however the relationship to the radiometric partitioning described by this simplified model will be degraded. Similarly, any sensor calibration problems or other degradations to the absolute radiance values measured will also be absorbed into the derived coefficients.

3.2 ELM Process

The ELM coefficients were generated manually using the image processing software ENVI (Environment for Visualizing Images) [77]. Details on this and other software packages used in the research are included in the Appendix. The pixels used to perform the ELM regression were selected to obtain an average radiance value for each calibration panel while minimizing the adjacency effects. The pixels were selected interactively one by one, to form a contiguous "region of interest" (ROI) consisting of pure pixels from each panel. The interactive display allows the user to see the spectral profile at each pixel, such that adjacency effects are easy to identify in the pixel spectrum. Pixels showing spectral contributions from the neighboring background were not included, resulting in a buffer of one or more pixels on the outside of the panel that were not included in the ROI. The calibration panels were sized to ensure at least one full pixel on the panel given the maximum flight level and instrument IFOV. At typical flight levels, between 4 and 40 pixels were averaged for each calibration panel. Figure 7 shows the imaged calibration panels and ELM ROIs for minimum and maximum flight levels. Each image contained 4-6 gray-shade panels of varying brightness (typically from 2% to 64% average reflectance) whose materials were chosen to have relatively featureless spectral profiles. Figure 8 shows typical ground truth spectra for the calibration panels.



Figure 7 Color composite images of calibration panels (left) and with ELM ROIs overlaid (right) for approximate flight levels of 5000 ft (upper figures) and 20,000 ft (lower figures) AGL.



Figure 8 Typical ground truth reflectance spectra of six gray shade calibration panels.

The ENVI Empirical Line Method algorithm was used to calculate the average radiance spectrum of each panel ROI and perform a linear ELM regression against the panels' ground truth reflectance spectra. The output consists of three vectors containing the gain, offset, and root mean squared error (RMSE) residuals for each wavelength band. An example of the coefficients is shown in Figure 9. HYDICE radiance cubes were provided in scaled spectral radiance units of $(1/75) * W / m^2 sr \mu m$. All radiance quantities presented in this study are converted to "microflick" units ($\mu W / cm^2 sr \mu m$), and all gain coefficients are scaled such that they produce standard radiance units when applied to normalized reflectance values (ranging from 0 to 1 for 0-100% reflectance). The gain and offset coefficients are plotted on different axes in order to overlay them on the same plot. These coefficients comprise the basic observational quantity used in the study.



Figure 9 Example of typical ELM coefficients (Site 4, 1998-08-25, 17:14z). The left y-axis scales the gain (*m*, in black) and the right y-axis scales the offset and RMSE residual (*b* and *err*, in red).

3.3 QUAC Process

The Quick Atmospheric Correction (QUAC) algorithm has become an extremely prevalent method for VNIR/SWIR HSI atmospheric compensation among practicing HSI analysts. It is fast, needs no external ground truth or atmospheric information, is tolerant of radiometric uncertainty, and highly robust. Even when other methods are ultimately used, QUAC often serves as a baseline for comparison. QUAC's main deficiency is in absolute accuracy of the reflectance returns, shown to be accurate within approximately 15% compared to the best FLAASH results (best meaning highly accurate radiometric input and characterization of the atmosphere) [48]. HSI analysts will often accept that degradation in absolute accuracy because many statistical algorithms employed in detection and identification problems are dependent on relative versus absolute accuracy, and when the data conditions for FLAASH are not ideal, the difference in accuracies is reduced or even reversed.

Like ELM, QUAC relies on in-scene information to derive the reflectance from the source radiance data. Instead of using ground truth reflectance measurements, however, QUAC performs the linear regression against an internal library of laboratory signatures. The assumption is, that for spectrally diverse HSI scenes, there will be an adequate number of material classes represented whose mean reflectance spectra match the material class library reflectance spectra. These library reflectances can then be used in place of the ELM ground truth spectra to perform the regression. ELM, in fact, is often employed in this manner for scenes without ground truth measurements; the difference is that QUAC chooses the scene endmembers and library signatures automatically. The approach is effective over a wide range of scenes, but does degrade in scenes that do not contain adequate spectral diversity [49].

QUAC is included in this research because it is a widely used, well-documented approach that yields consistent relative accuracy. QUAC gain and offset coefficients are analogous to ELM coefficients and are used for comparison. ELM depends only on the target spectra, so they were derived from small image "chips" such as those in Figure 7. To maximize the spectral diversity, QUAC was run on the full HSI data sets, typically 320 by 1280 spatial pixels. The QUAC coefficients were inverted to compare directly to

ELM coefficients because QUAC defines the gain to be $\rho/(L-b)$ rather than equation 4. QUAC does not provide statistics on the regression residuals so no RMSE is included with the QUAC coefficients.

3.4 ELM Postprocessing

Figure 9 depicts a classic, physically reasonable shape for the ELM coefficients – rising to a peak in the visible region and exponentially decaying toward the longer wavelengths – with relatively low residuals. The majority of the images produced ELM coefficients with similar spectra as expected. A sizable number of images, however, produced coefficients with unexpected offset spectra. Normally when ELM is used to obtain a reflectance image, the shape and values of the coefficients are of little consequence as long as they produce an acceptable reflectance inversion. For that reason the coefficients themselves have been rarely studied in the literature. In this application, however, we attribute the gain and offset to a physical partitioning of radiative transfer effects, so non-physical anomalies are important to address.

The most common offset coefficient anomaly observed is a clear "red edge" effect – a pronounced peak in amplitude in the near infrared region (Figure 10). This is commonly seen in spectra of vegetation but is not expected in the flat spectral response of the calibration panels. Despite efforts to avoid mixed target/background pixels, these effects are present in many of the vegetative scenes and can be attributed to scattering of light off the background vegetation into the imaging path. This scattering, known as the

adjacency effect, is part of the imaging environment and is considered as part of the consolidation of the environmental effects incorporated in the coefficients.



Figure 10 Example of ELM coefficients (Site 9, 1998-03-23, 14:13z) exhibiting background vegetation features at 0.7 to 0.9 μ m.



Figure 11 Example of ELM coefficients (Site 5, 1995-10-21, 20:20z) exhibiting strongly negative offset values from 0.55 to 1.8 μ m.

More concerning are artifacts such as those in Figure 11. The areas of negative offset values are completely nonphysical from the standpoint of our imaging model, i.e., the offset representing the path radiance. No literature was found that addresses offset coefficients such as these. Potential causes of the artifacts are: a) poor radiometric quality of the HSI data; b) poor ground truth reflectance spectra; or c) mis-partitioning of the radiance with respect to the radiative transfer model. The radiometric fidelity of the HYDICE sensor is well characterized and vetted in the literature so this is an unlikely source of error. Errors and marginal quality ground truth measurements are present in some of the collects, but these are usually noted in the ground truth documentation and can be discovered through careful review of the data. Furthermore, the artifacts appear episodically; for instance, several images might produce normal coefficients, and subsequent images on the same flight using the same ground truth produce coefficients with the negative gain artifacts. For these reasons, the artifacts are attributed to the ELM linear regression not fully partitioning the radiative transfer in the way described by the image model.

It should perhaps not be surprising that this sort of variability is observed in the offset coefficients, particularly when the RMSE residuals are large. ELM is fitting a line to a set of 4-6 radiance-reflectance pairs (one set per wavelength), and scatter in the data can easily shift the intercept when the values are so small relative to the gain (commonly well over an order of magnitude smaller). In a preliminary paper on this work, the issue was minimized by leaving one or more reflectance panel observations out of the ELM analysis, trying various combinations to reduce RMSE and improve the offset coefficients [78]. This was largely effective at producing more physical offset coefficients, but there was no objective measure to indicate whether the regression was improved or harmed by the reduction of observation points because the goodness of fit is often improved simply by reducing the number of data points from four to three, for example.

To understand the causes of the negative offset coefficients, it is necessary to scrutinize the wavelength-by-wavelength linear regressions in the ELM process. This was done by performing the linear regression at each wavelength in sequence, animating the line fit against the data points, and summing up the residuals at each point. Every scene

has a very dark panel, approximately 4% reflectance, that was useful to show how well the linear fit is "anchored" at low reflectance value near the y-axis. It was observed that while there were regions of scatter-related negative offsets (i.e., the linear fit passing below the dark panel observation point due to scatter in the data), that was not a consistent pattern in the large regions of the spectrum where the negative values occur. This was born out by the arithmetic sums of the dark panel residuals, which did not in general show a negative bias in the problem regions. The pattern that emerges is a slight apparent nonlinearity in the dark tail of the distribution, such that the slope of the line is reduced near the y-axis.

Second and third order fits were compared, but over the full range of data either worsened the overall fit or over-fit the data in images containing only four panels, in agreement with the body of work indicating a linear relationship. Instead, a small nonlinear term was added to equation 4 of the form:

$$\frac{h}{g\rho+1}.$$
 (7)

This term has value *h* at $\rho = 0$, decaying toward zero with an effective width defined by the parameter *g*. Adding to the right hand side of equation 4 gives:

$$L_{\rm S} = m\rho + b + \frac{h}{g\rho + 1} \,. \tag{8}$$



Figure 12 Generalized nonlinear correction function (solid red) and base linear function (dashed blue).

The resulting curve has the shape shown in Figure 12. To constrain the parameters *g* and *h*, we require that the slope be positive at $\rho = 0$. Differentiating equation 8 gives a slope at $\rho = 0$ of *m* - *hg* (where *m* is the *linear* slope of the line), which must be greater than or equal to zero, requiring:

$$g \le \frac{m}{h}.$$
 (9)

To express the constraints in terms of departure from linearity as shown in figure, we introduce a new parameter α and set:

$$g = \frac{m}{\alpha h}.$$
 (10)

Then the slope at $\rho = 0$ is given by:

$$slope_0 = m \left(1 - \frac{1}{\alpha} \right) \tag{11}$$

so α can be used to constrain the departure of the slope from *m* and *h* gives the departure of the intercept from *b*.

The correction can be incorporated in the ELM regression as follows. Equation 8 is linearized as:

$$L_{\rm S} = m\rho + b + h\delta \tag{12}$$

where δ is an estimated nonlinear term given by:

$$\delta = \left(\frac{m\rho}{\alpha h_{est}} + 1\right)^{-1}.$$
 (13)

Here h_{est} is an initial estimate of h. The linearized form (equation 12) can then be fitted to the data using conventional linear regression. The problem remains of how to determine h_{est} from the data. An algorithm was developed to calculate the estimated intercept in four ways: 1) linear regression (conventional ELM); 2) second order polynomial fit; 3) third order polynomial fit; and 4) using the linear regression value of the slope but forcing the line to exactly intersect the lowest reflectance/radiance pair. These four vectors (across all wavelengths) were evaluated by a set of objective criteria to reject sets that deviated too greatly from the ELM result (in peak magnitude), and those that had significant negative values. The estimated intercepts from the remaining methods were averaged to determine h_{est} . In practice, the algorithm sometimes failed by rejecting all candidates or by selecting a mean that was too far from the ELM result (despite meeting the criteria for selection). Each set of estimated offsets was plotted and reviewed manually to confirm or override the algorithmic result. The α parameter controls how quickly the nonlinear term falls to zero and was set empirically to a constant value that matched the bulk of the data well.

Equation 12 was then fitted to the data using linear regression and a set of nonlinear ELM coefficients was generated. Because the end purpose of the coefficients is to Empirical *Line* Method coefficients, it is desirable to keep the improved offset values but still form a linear fit to the data. The linear regression was therefore repeated by forcing the offsets to the nonlinear values and optimizing the gain to minimize the RMSE residuals. These adjusted coefficients were used in the end analysis (with a few exceptions to be described at the end of this section).

The constraints on the numerical adjustment were intentionally set conservatively in favor of the original coefficients, and even in cases showing large negative offset regions, the adjustment to the overall fit (and to the gain values) was small. The following series of figures illustrates the process. Figure 13 shows a set of ELM coefficients with areas of negative offset values, and the estimated intercept vectors for

the four estimation algorithms described above. In this case the linear and second order adjustments were rejected on the basis of large negative values and the third order and linear fit delta estimates were averaged to determine h_{est} . Figure 14 through Figure 16 show the linear, nonlinear, and adjusted linear fits to the data at several key wavelengths.



Figure 13 ELM coefficients (left, Site 9, 1998-03-23, 14:13z) and corresponding estimated offset spectra (right) for the four algorithms described in the text. The blue markers on the left indicate the wavelengths corresponding to Figure 14 through Figure 16.



Figure 14 Regression at 0.544 µm for linear, nonlinear, and adjusted linear methods (plots overlay exactly at this wavelength).



Figure 15 As Figure 14 but for 0.875 µm wavelength. Right hand panel is a detail enlargement of the lower reflectance region.



Figure 16 As Figure 14 but for 1.548 µm wavelength.

The figures illustrate the mild correction to the linear fit caused by the adjustment method. Figure 15 and Figure 16 are observations at the wavelengths of maximum correction; at the majority of the wavelengths the correction is negligible, as in Figure 14. Note that the y-axes in Figures 13 through 15 are labeled in native scaled HYDICE radiance units. The overall impact to the regression is demonstrated in Figure 17, which overlays the linear, nonlinear, and adjusted linear plots of the gain, offset, and RMSE residuals. The adjusted linear regression will, by definition, increase the residuals compared to the ELM least squares fit, but impact to the overall goodness of fit is small. Despite the very small increase in the residuals, the correction to the offset coefficients is significant.



Figure 17 Adjusted ELM coefficients (upper left) corresponding to the unprocessed coefficients in Figure 13. Upper right and lower panels show the linear, nonlinear, and adjusted linear coefficients and residuals. The adjusted linear offset is identical to the nonlinear offset.

The coefficient adjustment procedure was applied to all coefficients for consistency; in cases having high linearity and low scatter, the adjustment effect was

negligible. A few cases required larger adjustments than those shown in the example. The difference in the goodness of fit due to the adjustment procedure averaged across the entire collection is shown in Table 4. Over all images used in the study, the nonlinear function provided a slightly better fit than the linear regression. The adjusted linear regression produced larger residuals as expected, but only 5.6% greater than the linear case on average.

Table 4 RMSE residuals for linear, nonlinear, and adjusted linear regressions, averaged across all wavelengths and allimage coefficients in the study.

	Linear Regression	Nonlinear Regression	Adjusted Linear Regression
RMSE (μ W / cm ² sr μ m)	65.50	62.18	71.13

The improved fit of the nonlinear regression suggests that the functional form shown in Figure 12 more accurately matches the data than the linear form. Care was taken to minimize the departure from linearity, and to confine the nonlinearity to the low reflectance region of the relationship, so it is unlikely to be simply a result of over-fitting with a higher order polynomial. The specific form of the nonlinear term in fact has a physical basis. Conel [20] develops the form as an effect of the background spherical albedo, or the portion of incident radiance that is scattered isotropically. Using the MODTRAN radiative transfer partitioning described in section 3.5, it is possible to model the nonlinear effect. Figure 18 shows the results of a series of MODTRAN simulations. The blue data markers show the modeled radiance values for a series of decreasing surface reflectances. The red line is the best linear fit, and the lower plot shows the residuals from the regression. The variation from linear is very slight, but the dark tail nonlinearity is indicated by the increasing residuals as the reflectance is reduced from 3% to zero. A moderate (non-zero) background albedo was specified in this simulation, to be described more fully in the next section.



Figure 18 Modeled simulation of an "ELM" relationship. The upper plot shows the modeled radiance points for decreasing reflectance values against a linear fit. The lower figure shows the residuals from the linear regression.

One additional study was completed to verify that the adjustment procedure did not adversely affect the coefficients. A sample of coefficients that required significant adjustments to the offset spectra were applied to their source HSI radiance data to determine reflectances. The reflectance inversions were completed for the full linear ELM (denoted "ELM"), the linear ELM as modified in the preliminary study by eliminating one or more calibration panels ("ELM-f"), the adjusted linear ELM as described in this section ("ELM-Adj"), and QUAC. The resulting reflectance data were summed over all pixels (ranging between 204,800 and 409,600 points) and compared. The goal is to ensure that the improvements made to the offset coefficients do not negatively affect the reflectance values returned. Figure 19 and Figure 20 show the mean scene reflectance for four sample images, each of which exhibited significant regions of negative offset values in the original linear ELM regression.


Figure 19 Mean image reflectance values derived linear ELM, modified ELM, Adjusted ELM, and QUAC (see text) for sample images (Site 1, 1997-09-23, 17:03z (upper), Site 1, 1997-09-24, 18:57 (lower)) containing negative offsets.



Figure 20 Mean image reflectance values derived linear ELM, modified ELM, Adjusted ELM, and QUAC (see text) for sample images (Site 4, 1997-12-10, 20:42z (upper), Site 8, 1995-08-24, 13:26 (lower)) containing negative offsets.

In Figure 19 (upper), both the modified and adjusted coefficients produced similar results that are very comparable to the reflectances derived from the linear ELM; i.e., the adjustments did not significantly affect the reflectance inversion. In the lower panel, the two ELM modifications again produced similar results, but reflectance values are 0.02-0.05 reflectance units less than the full linear ELM. The new adjustment procedure is slightly less detrimental than the panel omission procedure. Figure 20 (upper) shows very close agreement up to approximately 1.9 μ m, above which the new adjustment method is clearly better. The lower panel seems to indicate that the panel omission procedure performed poorly in the 0.5 to 1.3 μ m region; both ELM-Adj and QUAC reflectances are very close to the original linear ELM. However, further study revealed that ELM-f actually produced better reflectance spectra.

Figure 21 (upper) shows the maximum reflectance values returned across all pixels rather than the mean. The original linear ELM, adjusted ELM, and QUAC all return reflectances well above 1.0 over sizable regions of the spectrum. The lower panel shows the mean reflectance spectra for an ROI over a specific target material (white canvas) found in the scene. The material had a ground truth reflectance spectrum measured, shown as the purple trace in the figure. This confirms that the panel omission procedure produces a superior reflectance return. In this case, one of the calibration panel ground truth spectral measurements was suspect, and it was appropriate to omit the data point from the ELM regression.



Figure 21 As Figure 20, but showing maximum reflectance values (upper) and ROI reflectances for a ground truth target material (lower), (Site 8, 1995-08-24, 13:26 (both panels)).

The observations about the plots in the preceding paragraph are confirmed objectively in Table 5. Two spectral similarity metrics were calculated for the cases described – Spectral Angle Mapping (SAM) [3] and Euclidean Distance (ED) [79]. SAM represents the spectra as n-dimensional vectors in wavelength space and measures the angle between them; ranging from 0 (perfect match) to $\pi/2$. ED is the l₂ norm of the difference in reflectance units across all wavelengths. The measures are similar, but as an additive distance metric, ED is sensitive to magnitude differences whereas SAM is invariant to absolute magnitude differences [80]. Image crf06m027 is similar in both measures. In image crf31m300, the ELM-f produced a slightly better spectral shape (lower SAM) but worse magnitudes (larger ED) than ELM-Adj. In crf43m046, the ELM-Adj result was vastly better in both measures. As discussed, crf07m82 ELM showed closer agreement to the ELM-Adj reflectances, but the ELM results were shown to be poor in this case, indicated by the ED scores of the reflectance retrievals referenced to the ground truth spectrum of the canvas material. Across the 127 images in the study, only 5% used less than the full set of available calibration panels to generate the ELM coefficients. The remainder used all available panels with the adjustment procedure outlined in this section.

Image	Figure	Spectral Angle Mapping (SAM)			Euclidean Distance (ED)		
		ELM-f	ELM-Adj	QUAC	ELM-f	ELM-Adj	QUAC
crf06m027	17 (upper)	0.0347	0.0396	0.0610	13.6779	12.8798	64.4738
crf31m300	17 (lower)	0.0437	0.0500	0.1302	45.9464	34.8019	39.3301
crf43m046	18 (upper)	0.2025	0.0279	0.1954	48.4526	12.1025	54.8883
crf07m82	18 (lower)	0.0731	0.0092	0.1069	51.6044	2.7050	31.7417
crf07m82 (GT)	19 (lower)	0.0566	0.0684	0.1832	44.5012	221.5883	260.1221

 Table 5
 Spectral similarity metrics for the adjustment procedure reflectance validation.

3.5 Coefficient Modeling

ELM derived coefficients correct for illumination and geometric factors as well as atmospheric effects. In order to compare coefficients derived from images collected under differing imaging conditions, it is necessary to modify the coefficients for variations in illumination and altitude of the observation. To diagnose these effects, the scenes were modeled using the MODTRAN radiative transfer software. MODTRAN models the atmosphere by treating it as a series of homogeneous layers characterized by their temperature, pressure and molecular composition. MODTRAN models the absorption, scattering and emission for each of the molecular constituents along a specified optical path, from the ultraviolet to far infrared range of the spectrum at up to 0.1 cm⁻¹ wavenumber resolution. MODTRAN can also provide solar illumination based on geographic position, date, and time. Thus the terms in equation 1 can be modeled explicitly.

In preliminary work, a "two stream method" was used to model the radiative transfer terms in equation 1, in which two MODTRAN simulations were run. Input parameters were identical except for the ground surface albedo, which was set to zero in one run and one in the other [78]. This isolated the path radiance term and allowed the other terms to be calculated using ratios of the MODTRAN outputs. The method was not optimal for two reasons. The method tries to force the PB model to provide radiative partitioning to the terms in equation 1, but the MODTRAN output fields are not consistent with that model and the results are therefore inexact. Second, setting the ground albedo to one across the entire field of view is an artificial extreme that amplifies any errors or artificialities present in the scattering model.

Fortunately, MODTRAN 5.3 provides access to radiometric fields that directly align with the model described in section 3.1. The derivation is fully developed in Appendix G of the MODTRAN 5.3 User's Manual [81]; the salient relationship is equation 11 from [81]:

$$L_{S} \approx L_{0} + \frac{A\rho + B\overline{\rho}}{1 - \overline{\sigma}\overline{\rho}} \,. \tag{14}$$

where ρ is target (pixel) reflectance, $\overline{\rho}$ is the average reflectance of the surrounding area, and $\overline{\sigma}$ is the spherical albedo. L_S is the sensor radiance and L_0 is the sensor radiance for the zero reflectance case. *A* and *B* are numerically derived coefficients described below. (The variable names have been changed from the reference to be consistent with the model developed in section 3.1, and the wavelength notation has been dropped. All variables except $\overline{\sigma}$ are wavelength dependent, and radiances L_S and L_0 are integrated over the instrument spectral channel.) Each of the variables on the right hand side of equation 14 is either an input to the simulation or can be calculated from MODTRAN output data. Coefficient *A* is defined as the product of the total transmitted solar irradiance and the sensor-to-ground direct transmittance, and *B* is the product of the total transmitted solar irradiance and the sensor-to-ground diffuse transmittance, each convolved with the channel spectral response function (SRF) [81]. Spherical albedo is the fraction of the incident irradiance that is reflected by the surface in all directions, summed over all wavelengths. This radiative transfer construct can be related to the ELM model noting that equation 14 can be written as equation 4:

$$L_{\rm S} = m^* \rho + b^* \tag{15}$$

where:

$$m^* = \frac{A}{1 - \overline{\sigma\rho}} \tag{16}$$

$$b^* = L_0 + \frac{B\overline{\rho}}{1 - \overline{\sigma}\overline{\rho}} \,. \tag{17}$$

The star superscript is introduced to specify modeled vice empirical coefficients. Referencing the remote sensing model developed in this study, equation 15 can be thought of as modifying the gain by accounting for the contribution of the spherical albedo term, and equation 16 as separating the offset into direct and diffuse parts. Using equations 15 and 16, PB modeled gain and offset coefficients can be computed from the output of a single MODTRAN run with ρ =0. This approach was shown to be superior to two- or three-stream methods [81]. The computation of the coefficients is simple as the variables are readily available in the MODTRAN output (details of specific MODTRAN output cards and fields used are given in the appendix). MODTRAN allows a spectral filter function (.FLT file) to be applied to the raw output, consisting of the spectral response function for the HSI sensor. MODTRAN then performs the spectral channel convolutions on the PB band model output automatically. In addition to the input parameters required to specify the radiative transfer operating modes, atmospheric profile and scene geometry, computation of m* and b* requires the input of $\overline{\rho}$, the nearby background mean reflectance.

The ability to model the gain and offset coefficients is the key enabler for this research. The PB model allows us to simulate coefficients under a variety of observational conditions, which can be used to standardize ELM coefficients so they can be compared and studied together as a whole, despite being derived under varying observational conditions.

3.6 Coefficient Standardization

One of the strengths of ELM is that the method implicitly accounts for illumination conditions, imaging geometry in addition to correcting for the atmospheric propagation and secondary surface interactions. However, the goal of this research is to use the empirical coefficients to compile a broad climatology relevant to HSI remote sensing. Therefore the illumination and geometric factors implicit in the coefficients must be factored out before they can be compared. This section presents the methodology for standardizing the coefficients for comparison across the range of conditions under which the HSI data were collected.

Any radiometric variability present in the data or ground truth is also folded in to the ELM coefficients. The absolute radiometric uncertainty for the HYDICE sensor is estimated to be less than 5% (see section 2.1), which is considerably lower than the variability caused by differences in the environment (illumination, geometric, and atmospheric). Included in the radiometric uncertainty are some sensor biases that would vary little between HYDICE collects, so the relative radiometric uncertainty between images in the study is likely lower. The radiometric uncertainty will be a more significant factor when extending the study to other sensors, but certainly within the HYDICE data, the radiometric consistency is adequate for the purposes of the study. Ground truth spectral measurement variability is a larger concern. The objective accuracy of the ground truth spectral measurements was 95%, but observational factors occasionally reduced the fidelity of the ground truth measurements, chiefly short time scale variability in illumination during the measurements, which would make the reflectance standard normalization unreliable. Fortunately, detailed ground truth collection logs were kept and any degraded conditions were noted in the reports, allowing suspect ground truth to be flagged and investigated as described in the case illustrated in Figure 21.

Physic based models provide an excellent tool for understanding and compensating for the effects of varying illumination and sensor geometry. Given the

geographic location, date and time, MODTRAN calculates the incident solar (and lunar, if desired) illumination present during the image collection. Imaging geometry (sensor altitude, look angles, and ground elevation) input is used by MODTRAN to calculate the path lengths and angular relationships. As described in the preceding section, modeled gain and offset coefficients can be computed for the imaging parameters relevant to each set of empirical coefficients derived from the imagery.

The modeled coefficients can then be used to estimate corrections to the ELM coefficients for differing illumination and sensor altitude above ground. All images are nadir-looking, so geometry is completely described by the altitude above ground level. For a set of ELM coefficients m_1 and b_1 derived from one image, modeled coefficients m_1^* and b_1^* are computed for the imaging conditions using equations 16 and 17 as described above. To estimate the ELM coefficients under differing conditions, at a later time, for instance, modeled coefficients m_2^* and b_2^* are computed for the new conditions. Scale corrections are then computed from the modeled coefficients:

$$M_{21} = m_2^* / m_1^* \tag{18}$$

$$B_{21} = b_2^* / b_1^* \,. \tag{19}$$

Then the ELM coefficients at the new time are estimated as:

$$m_2' = M_{21} * m_1 \tag{20}$$

$$b_2' = B_{21} * b_1 \tag{21}$$

where the primes delineate estimated (standardized) ELM coefficients.

Using this procedure it is possible to model the effects of differing illumination and geometry and to scale the ELM coefficients accordingly. This provides an estimate of what the ELM coefficients would be under differing imaging conditions and therefore allows comparison of coefficients across varied times and geometries. The atmospheric inputs to the PB model are identical in the two runs, so the scale correction will only account for changes in the illumination and imaging geometry. The variability that remains is assumed to be due to differences in the atmosphere. This method leverages strengths of the PB approach, e.g., calculating precise illumination and high fidelity atmospheric propagation, but by applying the *ratios* of modeled coefficients, it has the advantage of offsetting any systematic errors in the modeled results. Any artifacts caused by artificialities in the scattering models are present in both model runs and, to first order, cancel each other out. The accuracy of the estimated coefficients will decrease as the magnitude of the change in imaging conditions increases, but within some bounds of variability, the method produces accurate estimates.

4. SCIENTIFIC DATABASE

Rigorous remote sensing research studies may include careful measurement of the atmospheric conditions at the time of the collect. The overwhelming majority of them do not. In these cases, radiative transfer based atmospheric correction routines are run using model default or coarse climatological data for atmospheric inputs. Historically, this may have been the only option, as meteorological data have been difficult to obtain and require specialized knowledge and tools to analyze. With the expanding array of online climatological and observational atmospheric data available today, easy access to worldwide atmospheric data relevant to atmospheric compensation is possible.

Froude [82] describes a web-based system that uses the web-based Open-source Project for a Network Data Access Protocol (OPeNDAP) [83] to obtain meteorological datasets to be used in storm tracking software. The system automatically queries and parses large datasets on the remote server and feeds the input to the analysis software, enabling analysis of datasets that would be intractably large for local processing. In remote sensing research applications, atmospheric information has been compiled in databases for several purposes. Historical data is compiled for statistical reference, such as the Thermodynamic Initial Guess Retrieval (TIGR) database [84], which contains thousands of statistically representative atmospheric profile samples from worldwide rawindsonde observations. Databases are also used to store radiative transfer model

output for a range of atmospheric conditions. These are used as look up tables for atmospheric correction inversions [38, 39]. The use of atmospheric data from distributed data stores for atmospheric compensation research, however, is largely not addressed in the literature. Application of the open, web based protocols used successfully in other disciplines promises greater efficiency and higher fidelity data for atmospheric compensation research.

4.1 Research Need

In previous work, Powell et al. [68] described a relational database management system (RDBMS) developed to be used in atmospheric compensation research. That work focused on automated access and compilation of external atmospheric information relevant to the imagery. This research has much less emphasis on external sources, but a significant need for information organization and automated access to research data. The study used nearly 200 images, each with image metadata, site information, ground truth data and multiple versions of empirical coefficients. Each empirical coefficient has numerous MODTRAN runs associated with it for experimental trials and optimization, which are in turn associated with empirical coefficients from other images or standard spatial/temporal coordinates. Accurate accounting of the MODTRAN input parameters and their association with the related numerical output data and derived estimated coefficients is crucial to developing relationships amongst the data. Even with the relatively modest number of images used in the final results (127), approximately 1200 MODTRAN runs were executed, producing 8 GB of numerical output data. Manually cataloging that volume of data alone would be arduous, and with the requirement to maintain accurate relationships between data it becomes untenable. Furthermore, the need to manipulate subsets of data based on metadata conditions to perform the analysis calls for automated access. These information management needs were met by adapting the database design developed in [68] for this research application. The following sections describe the design and implementation of the RDBMS.

4.2 Application Overview

The purpose of the RDBMS is to provide researchers working on the atmospheric compensation of remotely sensed spectral imagery a single source for environmental data related to the imagery, and to manage the empirical coefficients, numerical output, and experimental trials conducted in the course of the research. The system must automatically query against a number of sources of environmental information based upon the image metadata, without requiring the researcher to search for, download, extract and convert the data from their native formats. This enables the researcher to assess the availability and variability of relevant environmental information. Although the user is insulated from the root data source, all source metadata must be stored to maintain the data pedigree. The system must also support the generation of input parameters for radiative transfer algorithms, and the cataloging of experimental runs and results.

Environmental data of interest to the research problem consist of Earth surface and atmospheric column data, and may include observed, modeled, or climatological

average values. Data sources were limited to those with worldwide coverage and at least five years of historical data (in the final design, those sources that did not cover the full time frame of the HYDICE data were not used). The data sources used are summarized below:

(1) FLAASH climatology: The FLAASH algorithm includes a coarse default climatology of integrated column water vapor values and surface air temperatures [85]. FLAASH default values are widely used, at least in casual atmospheric corrections, so they are considered in this project as the baseline to improve upon in the characterization of the atmosphere.

(2) COLA climate data: The Center for Ocean-Land-Atmosphere Studies (COLA) provides numerous climate-related datasets on the Internet. COLA is "a unique institution which allows earth scientists from several disciplines to work closely together on interdisciplinary research related to variability and predictability of Earth's climate on seasonal to decadal time scales" [86]. As part of its mission to share research and tools with the community, COLA provides numerous climate-related datasets on the internet. Geared toward climate modeling, the data are also well suited to radiative transfer modeling. The data obtained from COLA are surface elevation and meteorological parameters from the National Centers for Environmental Prediction's (NCEP) Reanalysis project [87].

(3) AIRS data: NASA's Atmospheric Infrared Sounder (AIRS) is a satelliteborne hyperspectral instrument designed to measure accurate temperature and water vapor vertical profiles. The AIRS derived atmospheric properties retrieved in this project

are surface elevation, surface air and skin temperature, integrated column water vapor, cloud fraction, total ozone burden, and profiles of water vapor mixing ratio, temperature, and geopotential height [12]. All AIRS data were provided by the NASA Goddard Earth Sciences Data and Information Services Center (GES DISC) [88].

The high-level design is shown in Figure 22. The system is designed around modular components that are loosely coupled to pass data and information. This design was chosen so the components could be developed and upgraded independently, with a minimum of integration effort required. The host file system stores the actual image files, MODTRAN data, and ancillary information and is not managed by the RDBMS. The file system is directly accessible by the user, and the database contains file system pointers to maintain the relationship to the data files. This introduces the possibility of inadvertently moving or deleting files without updating the database, potentially resulting in a loss of database integrity. However, it allows the system to catalog data that is stored in convenient, human readable storage volumes, and facilitates processing tools to access data files independently from the RBDMS. The external environmental data stores are not duplicated in the system; rather tailored relevant data is retrieved and stored in the system.



Figure 22 High-level system view of RDBMS.

These requirements called for a loosely coupled system using file based data transfer. For numerous reasons, the Extensible Markup Language (XML) [89] format was selected for the data transfer. XML is a markup language in the same family as HyperText Markup Language (HTML), but was designed for data transfer rather than data formatting and display. Like HTML, XML uses tags to identify information elements, but in XML the tags are not predefined and are tailored to the individual application. XML is simple and self-describing, so changes and potential errors introduced are easy to diagnose by inspection. It is text based and therefore platform/language independent. XML easily mimics relational database table schema and is very suitable for passing data to be loaded into the RDBMS. There are many open source libraries available to handle XML generation and parsing – a key attribute in a rapid development project. A drawback to XML – inefficient data storage – is not an issue in this system because the volume of data transferred is small. Figure 23 shows an example of an XML file used to populate a database table.

```
<?xml version="1.0" encoding="utf-8" ?>
<AIRS PROVENANCE>
  <IMAGEID>20080020</IMAGEID>
   <TITLE>AIRS Level 3 Daily Gridded Standard Product</TITLE>
   <SHORTNAME>"AIRX3STD"</SHORTNAME>
   <URL>http://acdisc.sci.gsfc.nasa.gov/opendap/Aqua AIRS Level
    3/AIRX3STD.005/2008/AIRS.2008.06.02.L3.RetStd001.v5.2.2.0.G0
    8156125406.hdf</URL>
   <DATETIME>2008-11-21 22:09:56</DATETIME>
   <GRIDLAT>38.5</GRIDLAT>
   <GRIDLON>-98.5</GRIDLON>
   <GRIDRES>1</GRIDRES>
   <GRIDUNITS>Decimal Degrees</GRIDUNITS>
   <DATATYPE>"HDFEOS V2.12"</DATATYPE>
   <HISTORY>"2008-06-04T16:54:06.000Z"</HISTORY>
   <PROCHIST>"L3 Daily products are a gridded output from L2
    retrievals."\n "Processing Level3 Daily PGE with ( 521 ) L2
    files."\n" Starting date is 6-2-2008 E with ( 521 ) L2
    files."\n" (Data starts the day before the one of
    interest.) "</PROCHIST>
   <PROCVER>"5.2.2.0"</PROCVER>
   <GANULEID>"AIRS.2008.06.02.L3.RetStd001.v5.2.2.0.
    G08156125406.hdf"</GANULEID>
   <QAPARAM>"Surface Skin Temperature"\n "Surface Air
    Temperature"\n "Atmospheric Temperature"\n "Water Vapor Mass
    Mixing Ratio"\n "Total Precipitable Water Vapor"\n "Ozone
    Volume Mixing Ratio"\n" ...</QAPARAM>
   <QAFLAG>"Passed"\n ...</QAFLAG>
   <QAEXP>"Based on percentage of product that is good. Suspect
    used where true quality is not known."\n ...</QAEXP>
   <<u>SOURCEDOC>http://disc.sci.gsfc.nasa.gov/AIRS/documentation</u>
    /v5 docs/AIRS V5 Release User Docs/V5 Data Release UG.pdf
    </SOURCEDOC>
</AIRS PROVENANCE>
```

Figure 23 Example XML file used to populate metadata into the database.

The project requirement to insulate the user from the various data sources and servers created a difficult design requirement for the data retrieval portion of the system – completely automated retrieval of environmental scientific data. Furthermore, server-side provisioning was strongly desired to limit the amount of data to be downloaded. One server protocol that supports the above requirements is OPeNDAP [83]. The goal of OPeNDAP is to allow remote access to datasets through the Internet in the same manner that they would be accessed if stored locally. The degree to which this is achieved depends upon the application using the protocol, but in any case, it does allow programmatic data retrieval via simple command line queries. OPeNDAP also allows server-side data provisioning, so if the user only needs one data point out of a large grid, the protocol allows the server to send only the grid point requested. Tremendous savings in communication and client processing costs are possible.

Data from OPeNDAP servers can be obtained by adding OPeNDAP libraries to existing analysis software such as MATLAB or the Grid Analysis and Display System (GrADS) (see the appendix for details). Some clients include full Graphical User Interfaces (GUIs), but data can also be queried through a command line. The OPeNDAP query consists of a function name and a Uniform Resource Locator (URL) that specifies the OPeNDAP server, directory and filename desired. The URL may have appended constraint expressions to identify spatial, temporal and content subsets of the data file requested. The protocol then performs the necessary server-side queries and reformatting to provide the data in the format requested by the client.

The types of data available from OPeNDAP servers are somewhat limited. OPeNDAP was initially developed for the oceanography community, and oceanographic data is well represented. Another community that has embraced the standard to some degree is the climate community, and a significant amount of atmospheric data is accessible via OPeNDAP servers. All of the external data sources were accessed using OPeNDAP.

One complication was discovered in using OPeNDAP servers – the user must know the directory and file structure on the server. OPeNDAP supports discovery of the data structure within the data files, but not the server directory structure. This is normally easy to discover by browsing the server, but in one case, the solution was not so simple. The NASA GES DISC OPeNDAP server providing the AIRS data includes the time that the data file was created in the filename. As there is no way for the remote system to know this filename, a screen scrape routine was employed to find the name of the appropriate file before it could be queried against. There is also the possibility that the server's directory structure might change over time, causing automated scripts to fail. Fortunately, most of the data servers are organized in simple hierarchical structures.

Systematic data movement and queries are handled by the MATLAB [90] processing engine via an Open Database Connectivity (ODBC) connector. The processing engine also drives the MODTRAN runs using an object oriented MATLAB class wrapper for MODTRAN 5 [91]. The database is also accessible through direct Structured Query Language (SQL) query. Details of the database and MATLAB software implementation are contained in the appendix.

4.3 Database Design

Figure 24 is an entity-relationship diagram (E-RD) showing the RDBMS schema. The major entities include the following. The image table contains image metadata and is related to the associated site, ground truth (GT) and sensor tables. The image table has a one-to-one relationship with the ELM and QUAC tables, which contain the file locations and statistics of the respective coefficients. The image table is one-to-many related to MODTRAN run and MODTRAN coefficient, which contain the MODTRAN and modeled coefficient metadata. Experimental trials are defined by the trial and standard coordinate tables, which are associated with a parallel pair of MODTRAN tables containing the metadata for coefficients generated for standard coordinates (as opposed to image coordinates). These relate to the image-based MODTRAN tables and the ELM table to describe the standard coefficient table, which contains the coefficient information after standardized to a given date/time and geometry. Finally, the standard coefficient table relates many-to-one to the atmosphere table, representing the clustering of standard coefficients that define the climatological classes. External environmental information is stored in the NCEP (National Centers for Environmental Prediction) and FLAASH climatology tables. The database is normalized in the first normal form [92] to ensure unambiguous data references. No further normalization was required. Data are loaded and updated if needed as a single unit for each data type, so the risk of database integrity problems from the non-normalization is minimal.

In addition to the technical metadata, the database includes fields to record experiment names, purposes, and observations or comments. The final versions of

processing and coefficients to be used in the research are flagged using Boolean "record" fields. (In one case an associative record trial table is used to relate each site to its record experimental trial. This table was added as a convenience to avoid complicated SQL queries otherwise required to make this link.) Design and implementation of the RDBMS entailed considerable overhead, but that investment was rewarded through improved data assurance and access to query-able records to manage the large volume of data. The ability to easily subset data based on flexible criteria for analysis proved invaluable in the analysis phase of the research. Additionally, the RDBMS establishes the groundwork to make the system available to the research community to expose and expand its holdings.

4.4 User Interface

The User Interface was originally designed for to perform two main functions – inputting image metadata to pass to the Processing Engine, and performing query/display on the RDBMS. These functions were developed as web-based applications that can be accessed remotely from the servers. The image metadata input interface includes all information needed to query the data stores, record the image file location on the file system, and identify the data set through user-oriented information tags such as image title, sensor name, and spectral range. The interface also includes free text fields to enter unique configuration information and spectral information, such as calibration data, about the sensor.



Figure 24 Entity-Relationship Diagram for the scientific RDBMS.

The metadata input interface was developed as an Active Server Page (ASP) file. This format was chosen because of its relative simplicity and an easily adaptable utility was available on the Internet to write the format out in XML. It consists of two scripts to generate the html form and to write out to an XML file. The page was hosted on a PC running Microsoft Internet Information Services 7.0 [93].

The key requirement for the User Interface was to allow simple query and display of the database contents. The database query interface is shown in Figure 25. The query interface was developed collaboratively as a class project [68]. Structured Query Language (SQL) queries are entered in the text box and the resulting data are displayed in the report area at the bottom of the screen. The interface includes functions to enter research notes into the database and links that describe the data sources and protocols used. The query interface was also developed as an active web script, but this portion is implemented as a Java Server Page (JSP). The application was developed using the NetBeans development environment [94], which is freely available. Database drivers are available to allow the scripts to interact directly with the database. Downloaded scripts were adapted from several sources [95, 96, 97]. The page was hosted on a PC running Java Server [98]. Both ASP and JSP technologies offer similar capabilities and the choice of scripting language was driven by individual programmer preference and the availability of easily adaptable utilities to interface with the database and other system modules.

http://127.0.0.1/www.page/FormSubmit.asp - Windows Internet Explorer									
C + ttp://127.0.0.1/www.page/FormSubmit.asp	 								
🚖 Favorites 🏽 🍘 http://127.0.0.1/www.page/FormSubmit.asp									
Database for Atmospheric Compensation Research									
Data Links	Please enter a query:								
COLA Climate Datasets <u>UMD Landcover Classification</u> <u>NASA AIRS Documentation</u> <u>OPeNDAP Protocol</u> <u>FLAASH Module Users Guide</u>		Submit	=						
Experiment run:									
Experiment method:	Submit								
Experiment comment:									
Query: <u>SELECT D.H2OVAP_NCEP_NAME, D.H2OVAP_NCEP, D.H2OVAP_AIRS_NAME, D.H2OVAP_AIRS_A,</u> <u>D.H2OVAP_SDEV_FROM_RSD.DERIVEDFLA_D_WHERE_D.IMAGEID =20080001</u>									
D.H2OVAP_NCEP_NAME D.H2OVAP_NCEP	D.H2OVAP_AIRS_NAME	D.H2OVAP_AIRS_A	D.H2OVAP_SDEV						
NCEP Reanalysis Column Mean Precipitable Water 2.997	AIRS Integrated Column H2O Vapor	2.175	0.879						

Figure 25 Database query interface.

4.5 Validation and Application Examples

This section addresses the usefulness of the RBDMS with respect to the project goals and design requirements. Several application examples are provided to illustrate potential uses of the system. These examples draw from a small sample dataset entered in the database, prior to the main experiment.

4.5.1 Data Validation

Data retrieved by automated scripts must be carefully verified to ensure the correct data are retrieved. Particularly when only a point value or small area is retrieved, errors in the positioning, scale or parameters retrieved might not be obvious by inspection, as they would be in the context of a large geographical area. Each of the data parameters retrieved were accordingly validated by comparing against the same dataset retrieved by another, non-automated method. Figure 26 shows the validation for the COLA elevation dataset. A larger area of topographic data was manually downloaded and plotted as a relief map. The image coordinate was selected on the relief map and compared to the values stored by the automated process in the XML file, shown in the inset of the figure.

The same general validation method was applied to random samples of each data type. Figure 27 shows the validation for AIRS integrated column water vapor (left) and water vapor profile profiles (right). In this case, the entire global HDF file was downloaded from a different NASA server by FTP. Using MATLAB's native HDF extraction GUI, the correct data fields and geographic locations were extracted, plotted and compared to the automatically generated XML data fields. The method does not independently validate the dataset, but does verify that the automated retrieval algorithms are functioning properly. This methodology did reveal an error in the retrieval code for this dataset, caused by accessing the incorrect geographical subset directory on the OPeNDAP server.



Figure 26 Validation for the site elevation automated retrieval. The image shows the manually downloaded digital elevation model with the site location labeled. The inset shows the automated retrieval of the elevation.



Figure 27 Validation for the AIRS integrated column water vapor (left) and water vapor profile profiles (right) automated retrieval. The image shows the manually downloaded data with the site location labeled. The inset shows the automated retrieval of the column water vapor. The plots show the comparison between the automated (OPeNDAP) and manual (FTP) retrievals.

4.5.2 Application Examples

The AIRS atmospheric profiles provide worldwide daily measurements that have previously only been available from sparse rawindsonde meteorological observations. These observations have been shown to improve numerical weather prediction, but their utility to atmospheric compensation remains to be demonstrated. As an example of an application of the scientific database, consider a key science question: Is the relatively coarse spatial and temporal resolution of the AIRS data sufficient to characterize the environment given the natural variability of the atmosphere?

To begin to address this question, the RDBMS was queried to compare the variation of the data across one observation day to the overall variability of the dataset. Figure 28 shows the profiles of all water vapor mixing ratios contained in the sample database (left) and the mean relative (fractional) difference between each pair of ascending/descending observations (center). The relative differences show a consistent 20% difference throughout the column, except at high levels where the water vapor content is very low. This is consistent with the 20% design tolerance of the AIRS water vapor retrieval, but the consistent negative bias is unexpected. The 20% average error is a small variation compared to the overall variability, and shown in the right hand plot. Here, the absolute difference is compared to the standard deviation of all observations at each level. This is merely suggestive, and to perform a more relevant comparison would require a larger data set over a single geographic area (or class of climate type).



Figure 28 Variability of AIRS water vapor profiles. The left plot shows the mixing ratio for all observations. The center plot shows the mean fractional difference between ascending/descending node pairs. The right plot shows the mean absolute difference and standard deviation of the pairs.

The effect of trace gases on atmospheric compensation problems is not often studied, in part because the required concentration data are not available. Using the AIRS soundings, such studies are possible. To design a study of trace gas effects requires knowledge of the variability of the gas concentration over a range of target images. Figure 29 shows the variability of ozone concentration over the sample database, categorized by the land cover of the underlying surface. The type of land use could be a contributing factor to total ozone concentration. The figure shows the mean value for each land cover type, and standard deviations where multiple observations are present. In this small data sample, no significant relationship is shown, as all but one category lay within or very nearly within one standard deviation of the other categories. This sampling is too small to draw conclusions, but it shows how the combination of the different data types can be leveraged by the RDBMS to address scientific questions.



Figure 29 Mean ozone burden from AIRS data for example locations plotted versus land cover type.

This chapter has described the goals and design of the RDBMS and provided several examples of how the system can support atmospheric compensation research. The system demonstrates the power of web-based, open data standards with server side provisioning to alleviate the burden of dealing with multiple data formats and large file downloads. These capabilities improve the research involving multiple data types by opening the door to more diverse datasets and allowing the researcher to devote more time to analysis rather than data gathering and grooming.

Several improvements to the system were noted that would provide greater utility. Expanding the User Interface to include web form controls and dropdown menus to construct simple queries would be highly desirable. This would relieve the user of entering SQL commands for often-used queries, and would open the system to users who are unfamiliar with SQL. Better integration of the modules, especially porting the GUI pages to run on the same server scripting language (ASP or JSP) would improve performance and maintainability. The process for loading the data into the database could be further automated. Lastly, the actual data manipulations performed by the MATLAB scripts are rudimentary, and could easily be ported to Java by an experienced programmer. This would eliminate the need to have MATLAB running on the server to retrieve the data.

There are many more datasets that could profitably be added to the system in time, but several additional data needs immediately present themselves. Firstly, aerosol extinction or at least visibility data is badly needed for HSI research in the visible spectral range. Secondly, the Level-2 AIRS products containing the individual soundings would provide higher temporal and spatial resolution to the data. Finally, while HSI collects typically cover geographically small areas, in some cases it might be beneficial to retrieve a small area of environmental data rather than just single point values. This would also provide the ability to measure the spatial variability and estimate the potential temporal variability of the atmosphere at the image site due to advection.

E-Science is typically driven by large research consortia, but this project demonstrates how small research efforts can benefit from e-Science technologies. With minor modification, the RDBMS presented here could be used for meteorological event analysis. For example, locations reporting local meteorological conditions such as poor urban air quality or anomalous wind events could be entered into the system to provide comprehensive atmospheric information associated with the events. Applications need not be constrained to atmospheric topics; given the variety of oceanographic and satellite information available using OPeNDAP, the techniques presented in this study are well suited to climate research as well. Applications beyond these could use a similar database approach but would likely require the use of a different access protocol that supports the specific data types.

The system evolved as the research progressed, primarily because the HSI data sets ultimately used predated the main external source of atmospheric data (AIRS) and because the increased requirement to manage numerically modeled radiative transfer output and experimental trials. Although the demands for external data access and provisioning were greatly reduced, the complexity of the database schema (Figure 24) grew significantly. As described in the following chapters, the research is dependent upon maintaining configuration control over model inputs and multiple versions of modeled and estimated atmospheric compensation coefficients, generated via several distinct workflows. The scientific database was a crucial enabler in keeping the integrity of the many associations among the data. In particular, the programmatic access by the analysis

software environment allowed simple selection of various categories and combinations of data to facilitate exploratory data analysis.

To fully implement the methodologies presented in the work will require opening the database to a larger community to provide a much greater range of empirical observations. The work done in designing and implementing the RDBMS will provide a head start in this work. The core of the future system is already developed, needing only a web front end and user interface to open access to the community for collaboration.

5. ANALYSIS

The standardized ELM coefficients constitute the basis for the climatological analysis. When properly standardized to remove the effects of illumination and scene geometry variability, they are directly comparable are analyzed as proxy atmospheric variables (transformed through atmospheric transmission and the ELM process into spectral space). This chapter first describes the optimization of the PB modeling that generates the modeled coefficients and the final selection of ELM coefficients to be used in the analysis. Then a series of case studies are described to validate the coefficient standardization for variations in time, altitude, and season. Lastly, the standardized coefficients are clustered and analyzed to develop separable climatological classes.

5.1 Modeled Coefficient Optimization

In preliminary work [78], the coefficient standardization method was shown to be resilient to magnitude errors in the modeled coefficients, for the reasons cited in section 3.6. However, MODTRAN provides tremendous flexibility to tailor the environmental inputs and it is desirable to use the most realistic modeling results possible in the analysis. To that end, an optimization study was completed for each site to determine the best set of input parameters to use. The standard for evaluation of modeled coefficients was the equivalent ELM coefficients, but "best" is not completely objective here because

in most cases there are trade-offs in the modeled results, e.g., better mean magnitudes versus overall spectral shape, or accuracy in the SWIR versus visible wavelengths.

Most of the defining radiometric model parameters were unchanged throughout the trials because the imaging model was constant. The exception is the spectral model band resolution, which defines the number and bandwidths of spectral bands over which the MODTRAN simulates the radiative transfer. The highest resolution supported by MODTRAN 5 is 1 cm⁻¹ (expressed in wavenumber) provides the best resolution of fine gas absorption features, but has a cost in execution time. The next coarser resolution supported is 5 cm⁻¹. Given that minimum band spacing of the HYDICE sensor is approximately 13 cm⁻¹ in the \sim 2.5 µm bands, the 5 cm⁻¹ resolution model is adequate to capture the broad features in the modeled coefficients. Figure 30 shows the difference in the output MODTRAN radiances for the 1 and 5 cm⁻¹ resolution models. Figure 31 shows the same MODTRAN output convolved with the HYDICE SRF. The plots overlay exactly over most of the wavelength range, but noticeable differences are evident near the water and oxygen absorption bands. The lower resolution model is acceptable for broad comparison and was used for the optimization runs, but the full resolution 1 cm⁻¹ model was used for all record runs used in the final analysis.

The list of major radiative transport driver parameters used in the all trials is given in Table 6. Complete descriptions of all input parameters and their use are contained in the MODTRAN Users Manual [81]. Multiple scattering and the first-principle, plane parallel atmosphere discrete ordinate multiple scattering algorithm (DISORT) are required to produce the atmospheric compensation outputs used to compute the modeled
coefficients, as are the related DSAZM and DSALB parameters. Surface reflectance is specified as Lambertian, and the surface skin temperature is set to 1 K to eliminate thermal emission in the computation of the reflected/scattered fluxes. LBMNAM controls the band model resolution, "f" means the default 1 cm⁻¹ resolution is run, otherwise the lower resolution parameter data file name is specified.





Figure 30 Example of MODTRAN modeled total radiance output at 1 cm-1 and 5 cm-1 band model resolution.



Figure 31 Example of MODTRAN modeled total radiance output at 1 cm-1 and 5 cm-1 band model resolution, resampled to HYDICE the spectral response function. Inset shows the detail from 1.1 to 1.5 μ m.

I dole 0	Dist of major radiative transfer artver parameters (common to an wrod river radiative).				
Input card	Field name	Value	Meaning		
1	MODTRN	Μ	Use MODTRAN band model		
1	ITYPE2	2	Slant path between two altitudes		
1	IEMSCT	2	Spectral thermal plus solar/lunar radiance		
1	IMULT	1	Multiple scattering mode		
1	SURREF	LAMBER	Lambertian surface approximation		
1	TPTEMP	1	Boundary temperature of image pixel (K)		
1A	DISORT	t	DISORT scattering algorithm is used		
1A	DISAZM	t	Azimutal depedence enabled for DISORT		
1A	DISALB	t	Sperical albedo calculated		
1A	LLFLTNM	t	Apply instrument SRF filter		
1A	H2OAER	t	Aerosol properties modified per water vapor specification		
1A	LBMNAM	See text	Band model data file name (defines model resolution)		
1A	CO2MX	390	CO2 mixing ratio		

 Table 6
 List of major radiative transfer driver parameters (common to all MODTRAN runs).

The MODTRAN input parameters defining the sensor geometry and illumination conditions are populated from the image metadata, specifically geographic location, date, time, sensor altitude, and ground elevation. Those related to atmospheric profiles (aerosol and water vapor models), seasonal models, scattering models and surface reflectance are also specified based on the characteristics of the site and environment, but defining the values to be used is not straightforward. The multiple scattering mode includes radiance that scatters off of the background surface (not in pixel IFOV) and into the sensor. MODTRAN input CSALB allows specification of the background reflectance as a function of wavelength. MODTRAN includes numerous default generic background reflectance profiles from spectral libraries. Experimentation with the generic backgrounds showed the modeled offset coefficients to be sensitive to these values, and best values did not always correspond to the type of land cover in the images. Representative HSI scenes were analyzed to extract mean reflectance signatures from each site; these were written to the MODTRAN library data and used as the background reflectance spectra. The spatial extent of the background that will contribute to the image will vary depending on altitude and atmospheric conditions, but an average effective radius of ~25 pixels was used. This is consistent with the background smoothing kernel size used in FLAASH [99].

Figure 32 shows an example of the modeled coefficients using the generic background reflectance signature (green) versus the scene-derived reflectance signature (red), plotted against the ELM-derived signature for Site 5 (barren desert). In this case the generic signature is fairly representative of the scene-derived signature, and neither modeled coefficients are clearly better. At other sites, especially those with vegetative backgrounds, the difference is larger and in general, the coefficients modeled using scene-derived background signatures are superior and these were used for all sites.

96



Figure 32 Modeled gain (upper) and offset (lower) coefficients compared to corresponding ELM coefficients.

Much more optimization was required to set the scattering related input parameters. An initial run was executed for image at each site, using the most appropriate settings for the known environments. With that run as a baseline, various parameters were modified to find the best set of input parameters. At each iteration, the differences between the modeled and corresponding actual ELM coefficients, averaged over all images at the site, served as the metric for suitability for that site. (Sites that were imaged in multiple seasons were each treated separately.) This optimization was performed empirically, but guiding the parameter adjustments based on knowledge of the site characteristics and presumed MODTRAN effects from the changes, rather than using an automated, systematic iteration of input parameter values. The parameters are too numerous to do the latter, especially since the goal is merely to get close to the ELM result rather that fully optimizing the parameters.

The modeled coefficients were found to be very sensitive to scattering related input specifications. The following figures illustrate the variation in modeled coefficients for the low elevation desert site (Site 5) as the input parameters are varied. The baseline parameters are:

- MODEL=2: Seasonal atmospheric profile, mid-latitude summer;
- IHAZE=10: Desert aerosol extinction model;
- VIS=0: Use visibility from default profile/extinction model determined by WSS for desert model;
- WSS=0: Wind speed = 0 m/s;
- ISEASN=1: Spring/summer.

The resulting modeled coefficient differences are shown in Figure 33. The differences are indicative of too much scattering – too much radiance is scattered out of the direct path (gain is too low), and too much diffuse radiance contributing to the path radiance (offset is too high). The spectral shapes are consistent, however, suggesting that the scattering model is correct.



Figure 33 Mean differences between modeled and ELM gain (left) and offset (right) coefficients, averaged over all images (7) of the site; generated using baseline scattering input parameters.



Figure 34 As in Figure 33 but with input parameter VIS=25 km.

Figure 34 shows the modeled coefficients with the default desert visibility overridden and set to VIS=25 km. Since the default value sets VIS based on an internal algorithm, that value is unknown. Based on the increase in scattering indicated compared to the default case, we can infer that setting VIS=25 km reduced the visibility, and worsened the scattering. For the next trial (Figure 35), VIS was increased to 40 km, which had the desired effect of reducing the scattering. The scattering remained too high, however, and further increasing VIS had little effect on the coefficients. The scattering model was then changed by setting the aerosol model to zero (IHAZE=0). Figure 36 shows the results - a much closer agreement with the ELM coefficients, indicating that the scattering model and input parameters are close to optimum. These input parameters were used for the final (record) run for use in the study.

This iterative process was repeated for each site to arrive at the best input parameters to us in the MODTRAN runs. Table 7 details the final input parameters used for each site to perform the MODTRAN runs for the remainder of the experiment. With the final modeled coefficients available, a final comparison was made of all the coefficients (ELM, adjusted ELM, QUAC, and modeled coefficients) along with available error statistics from the regressions and any notes from the ground truth reports or written during the ELM procedure. This display helped identify any gross inconsistencies in the coefficients and identify any unacceptable ground truth. Once the coefficients for each image were confirmed as final, they were annotated as record coefficients in the database, to be used in the development of the standard atmospheres.

100





Site	Input Card	Field	Value	Meaning
1	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	5	Urban extinction
	2	VIS	40	km visibility
	2	WSS	0	Default
2	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	10	Desert extinction
	2	VIS	18.5	km visibility
	2	WSS	0	Default
3	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	10	Desert extinction
	2	VIS	40	km visibility
	2	WSS	0	Default
4	1	MODEL	2/3	Mid-Latitude Summer/Winter
	2	IHAZE	10	Desert extinction
	2	VIS	0	Default
	2	WSS	0	Default
5	1	MODEL	2	Mid-Latitude Summer
(summer)	2	IHAZE	0	No aerosol extinction
	2	VIS	0	Default
	2	WSS	0	Default
5	1	MODEL	3-Jan	Mid-Latitude Winter
(winter)	2	IHAZE	10	Desert extinction
	2	VIS	0	Default
	2	WSS	0	Default
6	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	5	Urban extinction
	2	VIS	30	km visibility
	2	WSS	0	Default
7	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	5	Urban extinction
	2	VIS	40	km visibility
	2	WSS	0	Default
8	1	MODEL	2	Mid-Latitude Summer
	2	IHAZE	1	Rural extinction
	2	VIS	40	km visibility
	2	WSS	0	Default
9	1	MODEL	1	Tropical
	2	IHAZE	4	Maritime extinction
	2	VIS	17	km visibility
	2	WSS	0	Default

 Table 7
 MODTRAN Atmospheric and scattering input parameters for each site.



Figure 37 View of the data used to select the final ELM coefficients to be used to derive the standardized coefficients. Upper panels overly the empirical and modeled coefficients for comparison (gain, offset, and RMSE statistics, left to right). Lower panel shows ground truth and ELM notes (left) and the final ELM coefficients (right).

5.2 Illumination Adjustment Validation

To validate that the coefficient standardization can adjust for illumination variations, the method is tested against a time series of images at a single site. The series of seven images were collected from similar altitudes (~10,000 ft AGL) over a 90-minute period. The ELM coefficients from the first image (10:14 Local time) ELM coefficients were adjusted to the time of the last image (11:42L) using the method described in

section 3.6. Figure 38 shows the ELM coefficients before the illumination adjustment, and Figure 39 shows the coefficients standardized to 11:42L. The figures show that the gain variability from the change in illumination is almost completely corrected. The offset adjustment is not as complete, but the difference between the coefficients is reduced by the standardization by approximately 50%. Figure 40 shows the fractional error in the standardized coefficients. The gain is corrected to within 3% RMS and the offset to within 27% RMS (low SNR wavelengths in the broad water absorption bands are not included in the plots or statistics). The offset coefficient error is almost entirely negatively biased, suggesting that the modeled scattering is too low to fully compensate for the increase scattering effects at the later time. This could be because the scattering model is not quite correct, or it could be caused by the aerosol concentration changing from 10:14 to 11:42L.



Figure 38 Gain (left) and offset (right) ELM coefficients for beginning and ending times in image sequence.



Figure 39 As in Figure 38 but after 1714z coefficients are adjusted to 1842z.



Figure 40 Fractional errors for 90-minute time standardized coefficients trial.

5.3 Geometric Adjustment Validation

Since all images in the study are nadir looking, the only major geometric variables are the sensor altitude and site elevation, i.e., the sensor altitude referenced to ground level (AGL). The validation method used above was repeated, but where the two observations differ in altitude in addition to time. In this example, the first scene was imaged at 12:48 PM local time from an altitude of 5.0 kft AGL. The second scene was imaged 26 minutes later from an altitude of 10.4 kft AGL. Figure 41 shows the ELM coefficients before the illumination adjustment, and Figure 42 shows the coefficients standardized to 10,407 ft AGL. Figure 43 shows the fractional errors (low SNR wavelengths in the broad water absorption bands are not included in the plots or statistics).

The pre-adjustment gain coefficients changed little between the two observation times. This is because the altitude change (higher altitude means greater transmission loss) and the time change (higher sun angle means greater illumination) have opposing influences on the gain coefficient. The standardization routine decreases the gain for altitude change but increases gain for the time difference. With both geometry and illumination changing, the error is expected to be higher than in the illumination only case, but the size of the difference suggests that the geometric correction has higher uncertainty. This makes physical sense, because in the altitude correction, the PB model is simulating transmission path that is not present in the ELM coefficients that are being standardized. The results are therefore very sensitive to variability in the atmospheric with altitude. The gain is corrected to within 10% RMS and the offset to within 60%

106

RMS. Again, low SNR wavelengths in the broad water absorption bands and where the offset values are near zero are excluded.



Figure 41 Gain (left) and offset (right) ELM coefficients for 5 and 10 kft images.



Figure 42 As in Figure 41 but after 5 kft coefficients are adjusted to 10 kft altitude AGL.



Figure 43 Fractional errors for 5 kft altitude standardized coefficients trial.

5.4 Seasonal Adjustment Validation

The goal of the study is to identify and separate distinct climatic regimes, so the coefficients must be compared across differing seasons and geographic regions. An example of a seasonal comparison is shown in Figure 44 and Figure 45. Coefficients from a summertime collection of seven images were standardized to an image collected at the site in the wintertime. To the extent that the coefficient adjustment procedure adequately corrected for the differing illumination and scene geometry conditions, the differences in the coefficients are due to differing environmental states. In Figure 44 the wintertime ELM coefficients are shown (thick black line) against the collection of

standardized summertime ELM coefficients (colored lines). The fractional differences between the reference wintertime coefficients and the mean of the standardized summertime coefficients is shown in Figure 45.



Figure 44 Standardized gain (left) and offset (right) coefficients for the seasonal adjustment validation. The heavy black lines are the Dec reference ELM coefficients; the colored lines are the standardized August coefficients.



Figure 45 Fractional errors for Aug – Dec seasonal standardized coefficient trial.

109

The ensemble of standardized summertime coefficients is clearly separate from the wintertime ELM coefficient, suggesting that the environmental variability is larger than the variability remaining after the illumination and geometric adjustments. This is confirmed by the fractional errors, which are significantly larger than the variability among the ensemble members, which are represented in the case depicted in Figure 40. RMSE for the gain is 40% in the inter-seasonal case versus 3% for the summer ensemble. RMSE for the offset is 42% versus 27%. The images in the summertime ensemble were all collected from the same altitude as the wintertime image, so although the time of day varied by up to 3.5 hours, the consistent altitude gives this case less variability than one would expect on average. Still, the study suggests that the seasonal atmospheric signal can be large enough to detect through any residual error from the standardization procedure.

5.5 Global Coefficient Standardization

Up to this point, all examples shown have standardized coefficients from some number of images to correspond to the imaging conditions of another image to facilitate direct comparison. In order to form classes from the global collection of coefficients, it is necessary to standardize them all to a single reference point. The exact geospatial and temporal coordinates of the reference point can be somewhat arbitrary, but to minimize error, the reference point was chosen to lie near the median of the various imaging observational conditions comprising the study. Thus all coefficients are standardized relative to the coordinates listed in Table 8.

Coordinate	Reference Value
Geographic location	35° N, 95° W
Altitude	10,000 ft AGL
Date	8/15/1997
Time	17:00z

 Table 8
 List of standard reference coordinates for coefficient comparison

Figure 46 shows the coefficients for Site 4, standardized to the global standard reference coordinates. Site 4 is a high desert environment and all seven images are from the same altitude, so the standard coefficients make a tight grouping. The mean coefficients and standard deviations are plotted in Figure 47.



Figure 46 Site 4 gain (left) and offset (right) coefficients standardized to global reference coordinates (JJA).



Figure 47 Site 4 mean standard gain (upper) and offset (lower) coefficients (JJA). Standard deviation is plotted in black.

Not all sites are as uniform. Figure 48 shows the distribution of the standardized gain coefficients for Site 3, a high mountain location. The blue plot below the main cluster in the visible wavelength region had been flagged for poor data quality in earlier analysis. The remaining two low-gain coefficients had not been identified. Investigation in this case revealed that the ground truth was suspect; two of the calibration panel reflectances were apparently mislabeled and there was no definitive way to know if they were the correct spectra. These two were therefore rejected in the analysis. The resulting mean coefficients (averaged over the remaining nine coefficients) are shown in Figure 49.



Figure 48 Site 3 standardized gain coefficients (JJA).



Figure 49 Site 4 mean standard gain (upper) and offset (lower) coefficients (JJA). Standard deviation is plotted in black.

Among the sites with the highest variability is Site 9, the tropical site. The standardized gain coefficients have a wide range of values across the full spectrum (Figure 50). Both gain and offset (not shown) ensemble distributions, however, show a bimodal pattern, especially in the visible wavelength range. The coefficients were clustered using k-means algorithm [100 with two clusters specified, resulting in the cluster means shown in Figure 51. These appear to be distinct clusters, with separation in the visible wavelength region well beyond two standard deviations. There are no apparent indicators in the ground truth or image metadata to suspect an observational reason for the bimodal pattern, so it is likely caused by an environmental change during the five-day imaging timeframe.



Figure 50 Site 9 mean standard gain coefficients ensemble distribution (left) and mean (right) (MAM). Standard deviation is plotted in black.

Figure 52 shows the distribution of the mean standard gain and offset coefficients from all sites and imaging events. Separate means are included for the same site imaged in different seasons or years, and for the two separate means analyzed in the tropical site images (Site 9).



Figure 51 Site 9 clustered mean standard gain (upper) and offset (lower) coefficients (JJA). Dashed lines show the standard deviation of the clusters.



Figure 52 Means of the standard gain (upper) and offset (lower) coefficients for all sites and imaging events.

5.6 Coefficient Classes

The mean standard coefficient spectra in Figure 52 visually appear to have some groupings, and it is desirable to determine classes of coefficients from the data to

compare to climatological characteristics of the corresponding sites and images. The groupings of coefficients are explored by clustering the coefficients based on band-toband spectral magnitude, and then by clustering based on the spectral similarity metrics introduced in section 3.4. These analyses are described in the following sections.

5.6.1 K-means clustering

A k-means clustering was performed on the set of all coefficients to compare to the site means. Figure 53 shows the clusters of the gain coefficients along with the cluster means in black; Figure 54 shows the offset clusters. Figure 55 and Figure 56 show the means of these clusters superposed against the site means from Figure 52.



Figure 53 Clusters of all standard gain coefficients, with cluster means in black.



Figure 54 Clusters of all standard offset coefficients, with cluster means in black.



Figure 55 Clustered gain means from Figure 53 (black) with site means in color.



Figure 56 Clustered offset means from Figure 54 (black) with site means in color.

The clusters were formed by concatenating each standard gain coefficient with its corresponding offset, the computing the k-means cluster analysis for k=7. The number of clusters (7) was chosen as an approximate number of climate classes represented in the data. Figure 55 shows some correspondence between the cluster means found from the full set of gain coefficients and the means of the site-specific gains. Figure 57 maps the cluster members to their sites; e.g., 66% of the coefficients in cluster 1 are from Site 2, 7% from Site 5, etc. The dominant site is colored green; sites that have similar characteristics, either climate or land cover, are colored orange. Overall, 53% of the coefficients are clustered together with coefficients from the same site to form the majority, and 74% are clustered with the same site or a site with similar characteristics.



Figure 57 Aggregate site membership of the k-means coefficient clusters. Green indicates the primary site, orange indicates a site with similar characteristics, and other colors indicate dissimilar sites. Legend refers to the site number.

Sites 1, 2, 3, 5, 8, and 9 each make up the majority of at least one cluster. Sites 4 and 5 are both desert environments and are nearly equally distributed in cluster 3. Similarly, the tropical and temperate forest sites (8 and 9, respectively) have 30-40% overlap in clusters 3 and 4. Site 7, the high plains agricultural site, appears in 4 of 7 clusters with no majority in any.

Figure 58 shows the same clusters, but with coefficient membership mapped to their respective Köppen –Trewartha climate categories. The cluster membership is slightly more cohesive in this mapping, with 67% of the coefficients clustered together

with coefficients from the same category to form the majority, and 77% clustered with the same or a similar category. However, some discrimination ability is lost because some of the sites that are separable in Figure 57 are combined into a single category; notably the three mountainous sites (Sites 1, 2, and 3) which are all classified as highland (H). The defining characteristic of class H is that the climate is substantively affected by the elevation, which in fact covers a wide range of conditions. The sites' elevations range from 6810 to 9754 ft, and they could reasonably be separated into separate climate categories (although not defined by the Köppen –Trewartha system).



Figure 58 As in Figure 57 but with cluster members mapped to Köppen – Trewartha climate categories.

5.6.2 Spectral similarity clustering

Because k-means uses a distance metric to define the clusters, the clustering is heavily dominated by the contribution of the gain coefficients, which are typically an order of magnitude larger than the offset coefficients. The effects of aerosol scattering and background adjacency will have the greatest effect on the empirical offset coefficients, so it is preferable to have a contributing input signal to the clustering algorithm.

To better explore the contributions from and relationships between the gain and offset coefficients and the empirical clustering of the data, the two spectral similarity metrics described in Section 3.4, SAM and ED, are introduced into the cluster analysis. These metrics are complimentary because SAM is more sensitive to spectral features such as depth of absorption features and relative shape of the spectra, while ED is heavily weighted by the magnitudes of the spectra. The mean of all gain and offset coefficient site means was used as the reference spectra against which to apply the similarity metrics. SAM and ED were applied to each site mean gain and offset coefficient, denoted SAM_{sm}, SAM_{sb}, ED_{sm}, and ED_{sb}. Figure 59 shows the metric scatter plots with the site climate category denoted by the color scale. Figure 60 shows the same plots but with colors corresponding to the land cover class of the site.



Figure 59 Similarity metrics SAM and ED for site mean gain and offset coefficients plotted against each other to show relationships. Colors indicate the sites' climate classification categories.



Figure 60 Similarity metrics SAM and ED for site mean gain and offset coefficients plotted against each other to show relationships. Colors indicate the sites' land cover classification categories.

The scatter plots show some structure and clustering of points, particularly in the ED_{sm} vs. SAM_{sm} and ED_{sb} vs. SAM_{sb} plots. The colors comprising the clusters suggest similarity between the components, for example, the proximity of highland (H) and arid (BW) climate classes in lower left of these two plots in Figure 59, and the proximity of shrub and barren land cover classes in the same area of Figure 60. These correlations make physical sense since the classes have similar characteristics. The figures suggest that both types of site characteristics, climate class and land cover, affect the remote sensing radiance to reflectance retrieval. A seasonal grouping was also constructed, but is not included because the seasons collected are not well distributed across the climate classes; the only wintertime collections were in the arid (BW) classes.

Similar scatter plots were constructed using the full set of standard coefficients (rather than the site means) to provide more samples for clustering. As before, the mean of the site means was used as the reference spectra. This was chosen over the full coefficient ensemble mean because the sites do not each have the same number of coefficients and the ensemble mean would be weighted toward the sites with more samples. The same similarity metrics were computed for each individual gain and offset coefficient, denoted SAM_m, SAM_b, ED_m, and ED_b. Figure 61 shows the possible combinations. The colors and symbols represent differing sites and collection seasons. A dark background is used in the following figures to better show the data points.



Figure 61 Similarity metrics SAM and ED for all standard gain and offset coefficients plotted against each other to show relationships. Colors/symbols represent differing sites and collections.

Two of these are selected for cluster analysis, the ED_m vs. SAM_m and ED_m vs. SAM_b plots. Because of the vastly different scales for SAM and ED, the metrics are normalized prior to clustering. An agglomerative hierarchical clustering scheme [100] is used in which a dendrogram is created from closest proximity pairs, which are then consolidated pair-wise in successive iterations until the tree is reduced to only two classes. The clustering criterion can then be applied to any level of the dendrogram desired to produce the appropriate number of clusters. Without a training set or existing data to reference, no objective criteria were established to determine the best number of clusters. The number of clusters is set subjectively, considering the density of data and number of likely environmental classes represented. The distance metric used is Euclidean Distance. Many differing objective functions can be used to define the clusters. The best results were obtained using Ward's minimum variance criterion [101], which minimizes the variance of the points within the cluster at each step.

Figure 62 shows the full set of standard coefficients mapped as ED_m vs. SAM_m . The upper left plot colors/symbols represent the sites and collection groups of the source coefficients (same as the ED_m vs. SAM_m plot in Figure 61). The other plots show representative levels from the hierarchical clustering from 6 to 16 clusters. In these the colors/symbols represent cluster membership. Note that the specific color and symbol combination assigned to a coefficient is arbitrarily assigned and does not map from one plot to another. Rather, the groupings of color/symbol combinations are the points of comparison. As the number of clusters increases, the cluster membership becomes more aligned to the site membership in the upper left panel.

Figure 63 is identical to Figure 62 except that the metrics plotted are SAM_b vs. ED_m . In both figures it can be seen that the similarity metrics are, to a significant degree, able to separate the coefficients in accordance with their source image collections. Single clusters generally represent cohesive groupings coefficients from the same or similar sites (i.e., groups of same colored symbols in the upper left plot), particularly as the number of clusters increases. Those site member coefficients that are not grouped together in the site grouping scatter plots are of course not separable in this analysis. From visual inspection, it is not obvious which metric pair better separates the image sites into cohesive clusters.



Figure 62 Full set of standard coefficients plotted as ED_m vs. SAM_m . Upper left colors/symbols represent sites and collection groups of source coefficients. Other plots show clustering with increasing numbers of clusters; colors/symbols represent cluster membership.



Figure 63 As in Figure 62 except that the metrics platted are SAM_b vs. ED_m.


Figure 64 Aggregate site membership of the ED_m vs. SAM_m coefficient clusters. Green indicates the primary site, orange indicates a site with similar characteristics, and other colors indicate dissimilar sites. Legend refers to the site number.

To compare the spectral similarity clustering to the k-means clustering, the site and climate class membership of the clusters is analyzed as in section 4.6.1. Figure 64 and Figure 65 show the relationship of the cluster members to the site number and climate category, respectively, for the 7-cluster case. Overall, the ED_m vs. SAM_m clustering is not as effective (at 7 clusters) than the k-means clustering at clustering the primary sites together, showing 37% of the cluster members belonging to the primary cluster site, and 60% belonging to the primary or a similar site.



Figure 65 Aggregate climate category membership of the ED_m vs. SAM_m coefficient clusters. Green indicates the primary climate category, orange indicates a category with similar characteristics, and other colors indicate dissimilar climate categories.

Unlike the k-means clustering, all highland sites (Sites 1, 2, and 3) are grouped in a single cluster. The tropical site (Site 9) is more distinct from the temperate forest site (Site 8) in this clustering. As in the k-means clustering, Sites 6 and 7 are not the primary component of any cluster, and in most cases are evenly distributed across the clusters. Figure 65 shows the mapping to member climate categories. As in the k-means case, the aggregation of sites into the same climate categories improves the percentage of members clustered together, indicating that the distribution of differing sites in the clusters is to some degree driven by the climate characteristics of the sites.



Figure 66 Aggregate site membership of the SAM_b vs. ED_m coefficient clusters. Green indicates the primary site, orange indicates a site with similar characteristics, and other colors indicate dissimilar sites. Legend refers to the site number.

Figure 66 and Figure 67 show the same cluster membership mapping, but for the $SAM_b vs. ED_m$ clustering. The overall results are very similar to those of the $ED_m vs.$ SAM_m clustering, with 40% in the primary site group and 62% in the primary or similar site group. There are some notable differences, however. In this case, the highland sites are split between two clusters (3 and 4), with all three sites represented in each, and a substantial number of dissimilar sites in each. On the other hand, the vegetation dominated sites (Sites 7, 8, and 9) are considerably better separated here than in the ED_m vs. SAM_m or the k-means clustering. This suggests that the offset coefficients (and therefore path radiance) are more important for differentiating between the sites that are more influenced by the vegetative land cover. This supports the visual interpretation of the dominant features in the offset clusters of the forested sites in the near infrared region of the spectrum (Figure 54).



Figure 67 Aggregate climate category membership of the SAM_b vs. ED_m coefficient clusters. Green indicates the primary climate category, orange indicates a category with similar characteristics, and other colors indicate dissimilar climate categories.



Figure 68 As in Figure 66, but for 11-cluster case.

Figure 62 and Figure 63 appear to show increasing fidelity with greater numbers of clusters. The cluster membership mapping to the sites is shown in Figure 68 for the 11-cluster case. Eleven clusters produced a 10% improvement over the 7-cluster case in number of primary or similar sites clustered together (68% versus 62%). It produces similar, but slightly less fidelity than the 7-cluster k-means case. The differences between

the SAM_b vs. ED_m and ED_m vs. SAM_m clustering noted in the 7-cluster case hold true in the 11-cluster case as well – greater distinction between the vegetation dominated sites using the SAM_b vs. ED_m metrics (the ED_m vs. SAM_m case for n=11 is not shown), but slightly less distinction between the highland sites and between the desert sites. In addition to improving the overall site separation, the increased number of clusters showed improved ability to separate the seasonal variation. The clusters containing wintertime desert coefficients (cluster 10) are nearly completely separate from those containing summertime desert coefficients (clusters 2 and 6). The 11-cluster ED_m vs. SAM_m case additionally had separate clusters for the high and low desert sites, which were not well separated in any of the 7-cluster trials.

Table 9 summarizes the results of the cluster member identification. Aside from the improvements in differentiating some of the fine details noted above, the overall performance with respect to numbers of coefficients clustered with same-site coefficients was not markedly improved by using the spectral similarity metrics; the percentage was ~40-50% in most cases. In fact the greatest correlation to site occurred in an 11-cluster k-means trial (60% of the primary members of the clusters were from the same sites). The results suggest that the gain coefficients are the primary drivers for classifying the remote sensing environment. The offset coefficients add discrimination power in cases where the background scattering is a large component, like those sites containing dense vegetation. Despite the pains taken to ensure the offset coefficients were as consistent as possible, the relatively high noise in the offset coefficient signal, likely caused by transient, synoptic

scale variations in aerosol content, makes it difficult to characterize in all but the strongest cases.

Clustering Scheme	# Clusters	Identity Mapping	Members Mapped to					
			Primary Category	Primary or Similar Category	Dissimilar Category	Clusters w/Pri Majority	Clusters w/Pri/Sim Majority	Notes
k-means	7	site	53%	74%	26%	57%	100%	Highland sites separated; vegetative sites are not
		climate cat.	67%	77%	23%	86%	100%	Vegetative sites not separated
ED_m vs. SAM_m	7	site	37%	60%	40%	43%	71%	Highland sites in single cluster
		climate cat.	57%	69%	31%	57%	100%	Desert sites mixed among 3 clusters
ED_m vs. SAM_b	7	site	40%	62%	38%	43%	86%	Vegetative sites separated
		climate cat.	56%	67%	33%	43%	100%	Vegetative sites separated
ED_m vs. SAM_b	11	site	47%	68%	32%	45%	82%	High and low deserts separated; summer and winter deserts separated

Table 9Summary of cluster statistics.

5.6.3 Cross validation

In the preceding sections, the dataset used to form the clusters was also used to "score" the results in terms of homogeneity of clusters and ability to identify component climate or site membership. Without an independent test dataset, the results are likely positive biased, because the "training" data are not independent of the validation dataset. To reduce this bias, an "n-1" cross validation [100] was performed on the data. In this technique, one coefficient was held out as the validation sample, and the other n-1

coefficients used to form the clusters or means, where n is the total number of samples. Then the next coefficient was held out and the other n-1 coefficients (including the first validation sample) were used to form the clusters or means. The process was repeated n times such that each sample was used once as the validation dataset and each validation trial was independent of the data used to form the means, and the results averaged across all n trials.

First the coefficient membership by site was validated by performing the cross validation against the site means. In this case, the site means were calculated without the validation samples, and each validation sample was matched to a site based on spectral similarity to the site mean coefficients. In other words, this tests site membership in the case of perfect site clusters (because the "clusters" are formed by perfect knowledge of the site membership). Table 10 shows the results of the validation, in the cases of using either SAM or ED as the similarity metric. Since the "clustering" is perfect with respect to site membership, these numbers represent the upper bound on accuracy when using actual statistical clustering.

The cross validation was then performed on the clusters resulting from the kmeans clustering and repeated for k=3 to k=20 clusters. Site membership was used as the test criterion, but in order to stress the limits of separability, the sites were sub-divided in two ways. First, the site with images from multiple seasons (Site 5) was separated into three subclasses – summer, autumn and winter. Site 4 also encompassed summer and winter, but had only a single image in winter, so the cross validation would not be able to resolve a single sample class and Site 4 was therefore not separated. Secondly, the

bimodal distribution of coefficients noted in Site 9 (see Figure 51) were separated and treated as subclasses. This resulted in a total of twelve site classes.

		Members Mapped to			
Identity Mapping	Similarity Metric	Primary Category	Primary or Similar Category	Dissimilar Category	
site	ED	53%	74%	26%	
	SAM	6 7 %	77%	23%	

Table 10 Statistics for site membership identification.



Figure 69 Result of the cross validation of the k-means clusters, from n=3 to n=20 clusters, showing the percentage of samples correctly matched to site classes using the SAM similarity metric.

Figure 69 shows the result of the cross validation. The figure confirms that seven is too few clusters, at least in the expanded twelve-class case. The maximum score of with respect to site identification occurs between 13 and 18 clusters, suggesting that the number of separable classes is in that range. A range is given because there is random element in the k-means clustering (in the initial cluster seeding). Each validation trial was repeated five times and averaged, but there still appears to be some random scatter in the plot. The maximum values of 48-51% primary identification and 70% primary or similar identification compare well with the ideal results in Table 10. The scores are expected to be lower than the corresponding values in Table 9 because the expanded site classes contain 33% more classes to separate.

It is worthwhile to note that while the site membership (or climate category) is the best success criterion available at this time, it relies on the assumption that the sites as defined represent unique, separate atmospheric states. In general this is not the case; Site 9 is an extreme example but others showed lesser, but still significant, variability within the site, and could certainly overlap with other sites to some degree. Ultimately, validation should be based on accuracy of resulting reflectance retrievals, which is the true measure of effectiveness for the spectral climatology. The following section describes an example of how the spectral climatology could be applied in practical atmospheric compensation, and suggests the reflectance accuracy metrics that might be used to perform a more encompassing validation of the climatology.

5.7 Example Application

A notional case study was developed to illustrate how the spectral climatology could be used to aid in atmospheric compensation. The study represents the use case of determining surface reflectance spectra from an HSI radiance dataset that has no accompanying ground truth. A representative HSI dataset was selected from a site that was not used in the development of the climatology. The site is located near Site 8 and was imaged on September 22, 1999. The site was not included in the climatology because the ground truth data was incomplete. In this study, the spectral climatology is used as a lookup table that provides reflectance retrieval coefficients to be used on the new image. The scene contains calibration panels, but they are not used in the atmospheric compensation; they are used only to validate the reflectance values. The resulting reflectances are compared to ground truth spectra and a number of other standard methods of atmospheric compensation to gauge the effectiveness of the method.

There are a number of artificialities in this case study. While the image is independent of the data that comprises the climatology, the site is very similar to Site 8 so the error involved in selecting the site or coefficient cluster that best applies to the image is greatly reduced compared to the general case. (This offsets the fact that the range of environments in the climatology is relatively limited.) The image is typical but does contain greater than average spectral diversity, which is favorable for the other empirical methods. In more difficult, uniform environments both QUAC and ELM would be degraded. Lastly, no optimization has been done in selecting the coefficients from the climatology; the site mean of the closest environment was standardized for the new image conditions and applied directly to the radiance data.

The following figures show the results of the reflectance retrievals. In all figures, the low SNR atmospheric absorption regions of the spectra are omitted from the plots and the statistics. Figure 70 shows the reflectance spectra of the calibration panels in the HSI image. These spectra were not measured at the time of the collect, but were measured under similar conditions and should be sufficiently representative for this notional study. A 60% reflectance panel was present in the scene, but is not shown here because the mean brightness of the panel was at least 30% greater than any other pixel in the image, and none of the methods represented the reflectance spectrum well. Figure 71 shows the ELM derived reflectances, obtained by using identifiable materials in the scene and reference library reflectance spectra to perform the ELM. The scene contained asphalt roads and an airport runway with bright white painted areas. These provided a good range of light and dark reference spectra to perform the ELM retrieval. The plot illustrates the remarkable ability of the ELM procedure to remove noise and produce clean, smooth spectral shapes. In this case the dark panel is poorly represented, probably because the asphalt was not as aged or dusty as that in the library samples. The high reflectance regions are not high enough, for similar reasons.



Figure 70 Ground truth reflectance spectra for the calibration panels.



Figure 71 ELM reflectances of the three panels, obtained by performing ELM using library reference spectra.

Figure 72 shows the reflectances derived by using the spectral climatology "lookup table" coefficients against the image. High frequency noise is evident in the visible region of the spectra, especially in the higher reflectance spectra. This is likely caused by small errors in the resampling of the spectra; when even small features are spectrally misaligned, the ELM algorithm's inherent correction ability can work against itself, producing high frequency noise. The magnitudes are generally better than in the other methods, though, and notably better for the low reflectance panels. Figure 73 shows the reflectances derived from QUAC. The shapes of the spectra match those of the ground truth very well, however the magnitudes are considerably too high in all regions.



Figure 72 Reflectance spectra of the panels retrieved by using the climatology standardized coefficients.



Figure 73 Reflectances of the three panels from QUAC.

Figure 74 shows the reflectance spectra computed by FLAASH. The input parameters were not optimized any more than was necessary to get the program to execute, and atmospheric parameters were taken from the appropriate FLAASH default atmospheric models (which are the same as those in MODTRAN). DISORT multiple scattering was used with the 5 cm⁻¹ band model. No water vapor retrieval was used and spectral polishing was applied with a 9-band kernel [28, 29]. FLAASH also greatly overestimates the reflectances in most regions of the spectra for the 40% panel. The 4% and 15% panel magnitudes are closer to ground truth; similar to the QUAC retrievals.



Figure 74 FLAASH derived reflectances for the four panels.

Figure 75 and Figure 76 and Table 11 summarize the reflectance retrieval results in terms of the SAM and ED similarity metrics. Despite the high frequency noise in the Lookup ELM retrieval, the climatology method is similar to the Library ELM and slightly better than the FLAASH results in terms of average SAM value. The Lookup ELM retrieval performed the best in terms of total ED. QUAC outperformed all in terms of SAM, illustrating why it is so popular in whitened signal processing applications, where the magnitudes are less important. FLAASH had significant trouble calculating the higher reflectance values. The Library ELM trial also performed better in the medium reflectance values; the results are particularly degraded in the very dark regions.



Figure 75 Summary of SAM similarities to the ground truth reflectance spectra for the 4%, 15%, 40%, and 60% panels (left to right).



Figure 76 Summary of ED similarities to the ground truth reflectance spectra for the 4%, 15%, 40%, and 60% panels (left to right).

	Library ELM	Lookup ELM	QUAC	FLAASH
Mean SAM (radians)	0.251	0.259	0.119	0.270
Total ED (norm. refl.)	8.54	5.99	7.69	6.63

 Table 11 Summary of reflectance retrieval similarity metrics across all four panels.

In practice, the Library ELM results would be improved by iteratively applying other methods, perhaps DOS or RIM, to better ground the low reflectance retrievals. Similarly, the FLAASH results could likely be improved by iteratively optimizing the input values to the algorithm. In both cases, however, considerable time and expertise is required to obtain the best results. The most appealing feature of the spectral climatology lookup approach and, to a lesser degree, QUAC is that they are quick and (potentially) push-button easy. The spectral climatology lookup methodology needs to be refined to reduce the noise in the reflectance spectra. This could be an even greater problem when sensors other than HYDICE are applied to the database.

Notwithstanding the caveats discussed above regarding this single, notional example, the results are encouraging. The method's excellent performance in the low signal part of the example suggests that the path radiance term is well represented in the climatology. The example suggests a hybrid approach that leverages QUAC's ability to retrieve spectral shape while using the spectral climatology to ground the magnitudes and make up for lack of diversity in difficult scenes.

6. CONCLUSION

The empirical line method of atmospheric compensation has been one of the fundamental tools in the analysis of HSI data for many years, but the extension of the derived atmospheric coefficients to use outside of the source image has not been studied. In this research it has been shown that environmental information embedded within the ELM coefficients can be standardized, cataloged, and used to form a climatological compendium in spectral space. Common spectral similarity metrics were used to show that the climatological classes are separable to a degree of detail commensurate with the relatively modest size and range of the imaging conditions comprising the study. A notional application example was presented that showed competitive performance with other atmospheric compensation methods, and improvements in some areas.

Largely unstudied aspects of ELM process and coefficients are revealed in the work. The small magnitude but highly leveraged trades between gain and offset relative magnitudes that are not significant to ordinary ELM retrievals become problematic in this context because the imaging model equates the offset coefficient with path radiance, which should physically never be negative. In addition to random scatter, a small nonlinear effect of the spherical albedo was discovered in the data and a model developed to correct the coefficients to account for the effect.

A new method was developed to extract the atmospheric effects from the larger effects of varying solar illumination and observation altitude, which are also captured in the ELM coefficients, so that the coefficients can be compared across the range of imaging and solar geometries. The MODTRAN PB radiative transfer program was used to standardize the coefficients for analysis. This involved considerable study of the effects and sensitivities of the MODTRAN model to input parameters to optimize the radiative transfer simulations. The method was validated to correct for changes in solar illumination and sensor altitude using study data.

A scientific database was designed and implemented to manage the variants of coefficients, modeled radiative transfer data, experimental trials, and the interrelationships among them. While adding some overhead to the work, the database proved to be of tremendous value in providing automated query and access to the data for analysis application programs and maintaining information assurance. The RDBMS developed also provides the groundwork for exposing the data to the research community in the future for greater exposure and expansion of the data holdings.

Some of the imaged sites and climatological conditions were very conducive to the model and aligned nearly exactly with the standardization process. Others showed considerable variability that was not explained by the model. This is not surprising; although the body of conditions represented in the study data is relative large for an HSI collection, it is tiny compared to the universe of environmental conditions that can be present on the Earth. Much more data will be required to more fully characterize

climatological means and variability. Even so, the data in the study is over broad enough conditions to demonstrate separable climatological classes.

In addition to expanding the volume and breadth of data in the database, future work will focus on expanding the applicability of the database to other atmospheric compensation methods. The method used to standardize the empirical coefficients for comparison demonstrates that the ELM results can be used to improve other empirical and PB methods. In the same way that MODTRAN was used to generate the broad-brush conditions expected over a range of conditions and the ELM coefficients used to adjust the fine details, so they could be used to fine-tune other methods. In particular, QUAC is known to produce consistent, reasonable reflectance retrievals over a range of conditions, but with some deficiencies in absolute reflectance and spectral details over some parts of the spectrum. The database developed in this study could be correlated to corresponding QUAC coefficients and relationships developed to improve the absolute fidelity of the QUAC retrievals. This would have far-reaching consequences because of the widespread use of QUAC. If applied to QUAC or other fully automated empirical methods, the number of observations available to the database would increase many-fold, with a commensurate increase in its ability to characterize the environment.

APPENDIX

This appendix provides additional information on the software packages, custom code and methods used in the research.

A1. ENVI (Environment for Visualizing Images)

ENVI [77] is a commercial software package that is widely used for hyperspectral imagery analysis. It contains native functions to manage the hyperspectral data cubes, as well as many spectral analysis algorithms and functions, including ELM, image subsetting, interactive spectral plots, histogram/stretching, gain and offset application, plotting capability and statistics generation. The extra atmospheric compensation module adds FLAASH and QUAC routines. ENVI v4.8 with the atmospheric compensation module was used in this study.

ENVI was used for general scene familiarization and for the initial ELM coefficient generation. ROIs for the in-scene calibration panels were defined using the interactive ROI tool in conjunction with a dynamic z-profile of the pixel spectrum. The "point" type of ROI was used so that the spectrum of each pixel could be examined during ROI construction to minimize adjacency effects; mixed pixels near the edges of the panels are evident in the spectra. The ROIs were saved, along with ROI statistics, and imported into the ELM tool to perform the ELM retrieval. The ground truth reflectance

spectra were imported as ASCII text files and matched to the respective ROIs for the panels. ENVI outputs the gain and offset coefficients along with the RMSE residual in a text ".cff" file. These ".cff" files comprise the basic observational data for the research.

The ENVI ROI tool was also used to find the mean background reflectance values used in the MODTRAN runs. Derived ELM coefficients were first applied to representative scenes to convert the scene to reflectance. Then a suitable ROI was selected around the panels (exclusive of the panels themselves and any other targets) and the mean reflectance calculated. Similarly, ENVI was used to perform reflectance retrievals for comparison between the various methods of coefficient generation.

ENVI is packaged with IDL (Interactive Data Language), which allows for easy generation of batch processing scripts using ENVI calls. IDL was used to generate appropriate "header" files for the HYDICE imagery to allow ENVI to import the HSI data, and to generate and convert QUAC coefficients for each of the images.

A2. MATLAB

MATLAB [90] is a commercial software package that uses matrices as a fundamental data type. With a large library of matrix and other mathematical functions, it is well suited to imagery and HSI applications. MATLAB is also a full-functioned high level programming language and it was used for the bulk of the data processing and analysis in this research. MATLAB R2012a was used with optional statistics, database, and image processing toolboxes.

MATLAB allows the use of non-compiled textual scripts or user-defined functions called m-files to create processing applications. Unlike ENVI, there are no native HSI functions in MATLAB, so most processing functions were written from scratch. There is, however, a rich set of generalized mathematical, statistical and visualization functions from which to draw. The statistics toolbox contains built-in functions for k-means clustering, hierarchical dendrogram generation, and resampling/convolution that were used in the study. The database toolbox provides programmatic access to the ODBC connector from the MATLAB program environment and was used extensively to populate and query the database in the analysis.

Table 12 provides a summary of the major custom-built MATLAB modules that were used in the research and their respective functions. Three additional applications were harvested from open source file exchanges for use in the research. A low level plotting function called "ticklabelformat.m", (copyright © 2015, Yair Altman), was downloaded from the MATLAB Central file exchange site (http://www.mathworks.com/ matlabcentral/fileexchange/36254-ticklabelformat-set-a-dynamic-format-of-axes-ticklabels). The utility was used to control plot axis formats for the figures.

Also downloaded from the MATLAB Central file exchange site is a MODTRAN 5 class wrapper called "Mod5.m" (copyright © 2011, Derek Griffith). This utility provides a MATLAB class wrapper for the MODTRAN input card deck and utilities for executing MODTRAN functions as methods. The software was extremely useful in providing a way to organize and format inputs to MODTRAN directly from the MATLAB environment. Native MODTRAN 5 inputs are organized in a series of

interdependent hierarchical input "cards" that are textual, character-position defined, contextually varying formats left over from earlier predecessor programs. They are functional but very difficult to manage and the Mod5 wrapper generated all the required input cards behind the scenes.

Lastly, a MATLAB OPeNDAP connector called "loaddap" (copyright © 2014 OPeNDAP, Inc.) was downloaded from http://www.opendap.org/matlab-loaddap. This tool allows MATLAB to make calls to query and retrieve environmental data from OPeNDAP servers. OPeNDAP is described in the next section of the Appendix. The loaddap tool was useful early in the project but beginning with release 2012a, MATLAB supports Network Common Data Format (NetCDF) natively, so OPeNDAP data servers can be accessed directly without using loaddap.

Module name	Function
ReadELM.m	Read in ENVI coefficient file (.cff), convert to microns, resample to mean wavelengths, and calculate FWHM.
ReadQUAC.m	Read in QUAC coefficient file (.cff), convert to microns, resample to mean wavelengths, and calculate FWHM.
CffStat.m	Read in a series of coefficient files, calculate mean differences and variance.
ReadROIStat.m	Read in a set of ENVI ROI statistics spectra and GT spectra. Reformat for use in processing.
ELM2Mov.m	Given input ROI radiance and GT reflectance, calculate ELM regression and create an animation of the regression cycling through all wavelengths.
ELMAdj.m	Given input ROI radiance and GT reflectance, adjusts ELM coefficients to correct non-physical offset values.
DblPlot.m	Produce optimally scaled, double y-axis charts to display gain and offset spectra together.
ReadSummary.m	Import site and image metadata from Excel spreadsheet and

 Table 12 Listing of significant custom MATLAB scripts and functions used.

write to database.

MODRunAtm.m	Set properties and run MODTRAN 5.3 using MOD5 wrapper.
ElmFStat.m	Read Elm coefficients for all images and write to database
ErrorSum.m	Run through ELMs, sum error stats.
MODCoeff.m	Calculate modeled gain/offset coefficients from MODTRAN 5.3 run
ReadMODErr.m	Read in MODTRAN error coefficients.
WriteDB.m	Write a record to the database.
RunTrial.m	Populate trial table, run MODTRAN and generate coefficients* for all images and write to database.
CompCoeff.m	Run through images, display coefficients for comparison with stats and notes.
RecordCoeff.m	Designate set of MODTRAN coefficient as "record" coefficient - the ones that will be used.
MODRunStd.m	Set properties and run MODTRAN 5.3 using MOD5 wrapper. Populates ModRunStd table rather than ModRun.
MODCoeffStd.m	Calculate gain/offset coefficients from MODTRAN 5.3 run using new methodology for standardized coefficient run.
StdCoeff.m	Calculate Standardized Coefficients from ELM, m1* (idModCoeff) and m2* (idModCoeffStd)
RunGlobalStd.m	Populate Trial table for std coefficient trial, populate StdCoord table, Run MODTRAN and generate coefficient* for selected images, and write to database. Generates coefficient for all images assoc with all sites.
Git.m	Retrieve coefficients based on constraints passed (generate queries and read data).
CompStd.m	Compare coefficients to verify standardized coefficient generation
FracErr.m	Compute fractional error for coefficients.
ClusterWrite.m	Compute stats on standard coefficients and write to data directory.
LoadClust.m	Retrieve stored clusters of coefficients.
KCluster.m	Generate k-means clusters of coefficients.
SAM.m	Compute spectral angle mapping similarity metric.
ClusterMeans.m	Compute means of the clusters and cluster the site means.
SpecSim.m	Compute spectral similarity metrics SAM and ED, generate scatterplots.
ClusterID.m	Associate cluster members with image/site metadata.

A3. OPeNDAP

The Open-source Project for a Network Data Access Protocol (OPeNDAP) [83] is a protocol for requesting and transporting data across the web. The goal of OPeNDAP is to allow remote access to datasets through the Internet in the same manner that they would be accessed if stored locally. The degree to which this is achieved depends upon the application using the protocol, but in any case, it does allow programmatic data retrieval via simple command line queries. OPeNDAP also allows server-side data provisioning, so if the user only needs one data point out of a 10,000 point grid, the protocol allows the server to send only the grid point requested. Tremendous savings in communication and client processing costs are possible.

OPeNDAP was used to retrieve the NCEP reanalysis atmospheric data from Center for Ocean-Land-Atmosphere Studies (COLA) servers. A specific type of OPeNDAP compliant server called a GrADS Data Server additionally allows server-side data processing (GrADS is the Grid Analysis and Display System, developed and maintained by COLA). On a GDS server, the user can not only provision the data but also perform data processing such as averaging or differencing over space and time before transferring the result to the client. This capability is present on the COLA GDS server, but was not used in the project.

A4. MODTRAN

MODTRAN [78] is a community standard radiative transfer software package for modeling the transmission of light through the Earth's atmosphere, from the ultraviolet to the far infrared wavelengths. MODTRAN solves the radiative transfer equations from a fundamental physics approach, using a narrow band model of molecular and particulate absorption, emission, and scattering, as well as surface reflection and emission. The software models the solar and lunar illumination based on geographic location, date, and time. It can simulate a number of different remote sensing geometries, including the airborne Earth surface sensing application treated in this study. The atmosphere is modeled as stratified layers that can be user-defined or defaulted to one of several standard climatological profiles. MODTRAN is the standard U.S. Air Force radiative transfer model and is used by many elements of the Department of Defense. It is also commercially available from the ONTAR Corporation (http://www.modtran.org).

MODTRAN provides great flexibility for the user to specify input parameters defining the model operating modes and features, as well as the environmental conditions. So much flexibility is offered that using MODTRAN in a study can be daunting. The use of MODTRAN in this study was not intended to provide the absolute most realistic simulation for a given set of conditions, but rather to generate a reasonable estimate of the effects of changing imaging geometry and illumination. In this research, the parameters that were exercised consist mainly of geometric, spatial and temporal inputs (for illumination definition), and of those describing the default atmospheres listed

in Table 6 and Table 7. The default climatological atmospheres were modified by changing the water vapor content or aerosol model and concentrations, but user-defined profiles were not used. CO2 mixing ratio was set to a modern value of 390 ppmv. Poor results were obtained from the default choices in the off-target surface albedo file, so observed background reflectance spectra were added to the file and used in the simulations. Other surface parameters were set as described in section 3.5.

A Gaussian sampling filter function (".flt" file) was developed for the HYDICE spectral response function and was applied to the MODTRAN output. This was done within MODTRAN, so that convolved output files (".chn" files) were produced with HYDICE channel radiance values in addition to the band model resolution spectral radiance output files. These output files were used directly (after unit conversion) in equations 16 and 17 to calculate the modeled gain and offset coefficients.

A5. MySQL

The database used in this work is MySQL Community Edition Server version 5.6.23 [102]. MySQL was selected for its simplicity of implementation and widespread use. The Community Edition is free for download for non-commercial applications from http://dev.mysql.com/downloads/mysql/. The database runs as a server that is accessible though other downloadable utilities (SQL command line client 5.6, MySQL Workbench 6.1 with Query Browser and Model Builder, MySQL Notifier 1.1.5, MySQL Connector Net 6.8.3) or an ODBC connector. The MySQL Workbench includes a user friendly, graphical E-RD design tool that allows the user to specify all tables and relationships in

the databases schema, and then "forward engineer" the SQL script needed to build the database instance. The MySQL Workbench was used throughout the research to modify the database schema and perform ad hoc SQL queries. Programmatic access was facilitated through the ODBC connector (downloaded from http://dev.mysql.com) and the MATLAB database toolbox. This allowed MATLAB fetch and write commands to interface with the database and provide data directly from/to MATLAB variables.

REFERENCES

- [1] G. Chander, B. L. Markham and D. L. Helder, "Summary of current radiometric calibration coefficients for Landsat MSS, TM, ETM+, and EO-1 ALI sensors," *Remote Sens. Environ.* **113**, 893–903 (2009) [doi:10.1016/j.rse.2009.01.007].
- [2] X. Xiong, V.V. Salomonson, W.L. Barnes, B. Guenther, X. Xie and J. Sun, "An overview of terra MODIS reflective solar bands on-orbit calibration," *Geosci. Remote Sens. Symposium Proceedings*, 2006, pp 1111-1114 (2006) [doi: 10.1109/IGARSS.2006.287].
- [3] F. D. van der Meer and Steven M. de Jong, Imaging Spectrometry Basic Principles and Prospective Applications, pp XXI, 31-41, Kluwer Academic Publishers, Dordrecht, Netherlands (2003).
- [4] A. F. H. Goetz, "Principles of narrow band spectrometry in the visible and IR: instruments and data analysis," in Imaging Spectroscopy: Fundamentals and prospective applications, Kluwer Academic Publishers, Dordrecht, pp 21-32 (1992).
- [5] R. N. Clark, "Spectroscopy of Rocks and Minerals, and Principles of Spectroscopy," in Manual of Remote Sensing, Volume 3, Remote Sensing for the Earth Sciences, John Wiley and Sons, New York, pp 3-58 (1999).
- [6] G. Vane, R.O. Green, T.G. Chrien, H.T. Enmark, E.G. Hansen and W.M. Porter, "The Airborne Visible/Infrared Imaging Spectrometer (AVRIS)," Remote Sens. Environ. 44, 127–143 (1993) [doi: 10.1016/0034-4257(93)90012-M].
- [7] M.E. Kappus, W.S. Aldrich, R.G. Resmini and P. Mitchell, "The flexible HYDICE sensor's first year of operation," Proc. of the 11th Thematic Conference on Geologic Remote Sensing 1, pp 433-441 (1996).
- [8] S. J. Young, B. R. Johnson and J. A. Hackwell, "An in-scene method for atmospheric compensation of thermal hyperspectral data," J. Geophys. Res. 107(D24), 4774 (2002) [doi: 10.1029/2001JD001266].

- [9] R. Gomez, "Technology assessment of remote sensing applications in transportation: hyperspectral imaging (HSI)," National Consortia on Remote Sensing in Transportation, University of New Mexico, Albuquerque, NM, pp 1-14 (2001).
- [10] M. A. Folkman, J. Pearlman, L. Liao and P. Jarecke, "EO-1/Hyperion hyperspectral imager design, development, characterization, and calibration," Proc. SPIE 4151, 40-51 (2001) [doi: 10.1117/12.417022].
- [11] C.O. Davis, M. Corson, R Lucke, R. Arnone, and R. Gould, "The Hyperspectral Imager for the Coastal Ocean (HICO): Sensor and Data Processing Overview," International Ocean-Colour Coordinating Group, http://www.ioccg.org/sensors/Davis_HICO_IOCCG-15.pdf (2015).
- [12] E. T. Olsen, Ed., "AIRS/AMSU/HSB Version 5 Data Release User Guide," Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA (2007).
- [13] R. A. Rhode, "Solar Radiation Spectrum," Global Warming Art, 9 June 2007, http://www.globalwarmingart.com/wiki/Image:Solar_Spectrum_png, (6 May 2009).
- [14] B. Gao, C. O. Davis, and A. F. H. Goetz, "A review of atmospheric correction techniques for hyperspectral remote sensing of land surfaces and ocean color," IEEE Proceedings of the Int. Geosci. and Remote Sens. Symposium 4, 1979-1981 (2006) [doi: 10.1109/IGARSS.2006.512].
- [15] A. F. H. Goetz, J. W. Boardman, B. Kindel and K. B. Heidebrecht, "Atmospheric corrections: On deriving surface reflectance from hyperspectral imagers," Proc. SPIE 3118, 14-22 (1997) [doi: 10.1117/12.283831].
- [16] M. K. Griffin and H. K. Burke, "Compensation of hyperspectral data for atmospheric effects," Lincoln Laboratory Journal 14(1), 29-54 (2003).
- [17] A. A. Green and M. D. Craig, "Analysis of aircraft spectrometer data with logarithmic residuals," Proceedings of the Airborne Imaging Spectrometer Data Analysis Workshop, Jet Propulsion Laboratory Pub 85-41, 111-119 (1985).
- [18] J.B. Campbell, "Evaluation of the dark-object subtraction technique for adjustment of multispectral remote-sensing data," Proc. SPIE 1819, Digital Image Processing and Visual Communications Technologies in the Earth and Atmospheric Sciences II, (1993) [doi: 10.1117/12.142198].

- [19] R.E. Crippen, "The regression intersection method of adjusting image data for band ratioing," Int. J. of Remote Sensing 8, 137-155 (1987).
- [20] J. E. Conel, "Determination of surface reflectance and estimates of atmospheric optical depth and single scattering albedo from Landsat Thematic Mapper data," Int. J. Remote Sensing 11(5), 783-828 (1990) [doi: 10.1080/01431169008955057].
- [21] R.N. Clark, et al., "Surface reflectance calibration of terrestrial imaging spectroscopy data: a tutorial using AVIRIS," Proc of the 10th Airborne Earth Science Workshop, JPL Publication 02-1 (2002).
- [22] G. E. Smith and E. J. Milton, "The use of the empirical line method to calibrate remotely sensed data to reflectance," Int. J. Remote Sensing 20(13), 2653-2662 (1999) [doi: 10.1080/014311699211994].
- [23] W. M. Baugh and D. P. Groeneveld, "Empirical proof of the empirical line," Int. J. Remote Sensing 29(3), 665-672 (2008) [doi: 10.1080/01431160701352162].
- [24] B. Bartlett and J. R. Schott, "Atmospheric compensation in the presence of clouds: an adaptive empirical line method (AELM) approach," J. Applied Remote Sens. 3, 033507 (2009) [doi: 10.1117/1.3091937].
- [25] Berk, A., G.P. Anderson, P.K. Acharya, L.S. Bernstein, L. Muratov, J. Lee, M. Fox, S.M. Adler-Golden, J.H. Chetwynd, M.L. Hoke, R.B Lockwood, J.A. Gardner, T.W. Cooley, C.C. Borel, P.E. Lewis and E.P. Shettle, "MODTRAN5: 2006 Update," Proc. SPIE, Vol. 6233, 62331F, 2006.
- [26] Alexander Berk, Gail P. Anderson, "Impact of MODTRAN®5.1 on Atmospheric Compensation," IGARSS 2008, 127-129 (2008) [doi: 10.1109/IGARSS .2008.4779299].
- [27] S. Y. Kotchenova, E. F. Vermote, R. Matarrese, and F. J. Klemm, Jr., "Validation of a vector version of the 6S radiative transfer code for atmospheric correction of satellite data. Part I: Path radiance," App. Optics 45(26), 6762-6774 (2006) [doi: 10.1364/AO.45.006762].
- [28] T. Perkins, S. Adler-Golden, M. Matthew, A. Berk, G. Anderson, J. Gardner and G. Felde, "Retrieval of atmospheric properties from hyper- and multi-spectral imagery with the FLAASH atmospheric correction algorithm," Proc. SPIE, Remote Sensing of Clouds and the Atmosphere X 5979, 59790E (2005) [doi: 10.1117/12.626526].

- [29] T. Perkins, S.M. Adler-Golden, P. Cappelaere, and D. Mandl, "High-speed Atmospheric Correction for Spectral Image Processing", Proc. SPIE 8390, Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery XVIII, 83900V (2012) [doi: 10.1117/12.918908].
- [30] M. D. Abel, J.M. Zenner, G.A. Petrick, A.T. Buswell, M.L. Pilati, W.R. Czyzewski, L.P. Alessandro and S.K. Weaver, "New approach to atmospheric correction of hyperspectral data," Proc. SPIE, Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery VIII 4725, 72-82 (2002) [doi: 10.1117/12.478738].
- [31] S. Adler-Golden, M.W. Matthew, L.S. Bernstein, R.Y. Levine, A. Berk, S.C. Richtsmeier, P.K. Acharya, G.P. Anderson, G. Felde, J. Gardner, M. Hoke, L.S. Jeong, B. Pukall, J. Mello, A. Ratkowski and H.H. Burke, "Atmospheric Correction for Short-wave Spectral Imagery Based on MODTRAN4," Proc. SPIE 3753, 61-69 (1999).
- [32] B. Gao, K. B. Heidebrecht, and A. F. H. Goetz, "Derivation of scaled surface reflectances from AVIRIS data," Remote Sens. Environ. 44, 165-178 (1993) [doi: 10.1016/0034-4257(93)90014-0].
- [33] V. Carrere and J. E. Conel, "Recovery of atmospheric water vapor total column abundance from imaging spectrometer data around 940 nm sensitivity analysis and application to Airborne Visible/Infrared Imaging Spectrometer (AVIRIS) data," Remote Sens. Environ. 44, 179-204 (1993) [doi: 10.1016/0034-4257(93)90015-P].
- [34] Q. Cheng, D. Pan, D. Wang, J. Chen and T. Mao, "Retrieval of atmospheric water content based on AISA+ data," Proc. SPIE, Remote Sensing of Clouds and the Atmosphere XII 6745, 67450C (2007) [doi: 10.1117/12.731216].
- [35] Z. Qu, B. C. Kindel, and A. F. H. Goetz, "The high accuracy atmospheric correction for hyperspectral data (HATCH) model," IEEE Trans. Geosci. Remote Sens. 41(6), 1223-1231 (2003) [doi: 10.1109/TGRS.2003.812905].
- [36] R. Marion, R. Michel, and C. Faye, "Atmospheric correction of hyperspectral data over dark surfaces via simulated annealing," IEEE Trans. Geosci. Remote Sens. 44(6), 1566-1574 (2006) [doi: 10.1109/TGRS.2006.870408].
- [37] B. Cairns, B. E. Carlson, R. Yinga, and J. Laveignec, "Accuracy vs Speed: Evaluation of tradeoffs in atmospheric correction methods," Proc. SPIE, Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery VIII 4725, 427-437 (2002) [doi: 10.1117/12.478776].

- [38] R. Richter and D. Schläpfer, "Considerations on water vapor and surface reflectance retrievals for a spaceborne imaging spectrometer," IEEE Trans. Geosci. Remote Sens. 46(7), 1958-1966 (2008) [doi: 10.1109/TGRS.2008.916470].
- [39] A. Murray, B. Eng and K. Thome, "Implementation of a very large atmospheric correction lookup table for ASTER using a relational database management system," Proc. SPIE, 2820, 245-258 (1996) [doi: 10.1117/12.258105].
- [40] L. C. Sanders, J. R. Schott and R. Raqueno, "A VNIR/SWIR atmospheric correction algorithm for hyperspectral imagery with adjacency effect," Remote Sens. Environ. 78, 252-263 (2001) [doi: 10.1016/S0034-4257(01)00219-X].
- [41] B. Gross, O. Ogunwuyi, F. Moshary, S. Ahmed and B. Cairns, "Aerosol retrieval over urban areas using spatial regression between V/NIR and MIR Hyperion channels," IEEE Workshop on Remote Sens. of Atm. Aerosols, 2005, 43-50 (2005) [doi: 10.1109/AERSOL.2005.1494148].
- [42] K. Stamnes, W. Lia, H. Eidea and J. J. Stamnes, "Challenges in atmospheric correction of satellite imagery," Proc. SPIE, Ecosystems' Dynamics, Agricultural Remote Sensing and Modeling, and Site-Specific Agriculture 5153, 1-12 (2003) [doi: 10.1117/12.511229].
- [43] Exelis Visual Information Solutions, FLAASH Module User's Guide, Version 5.2, (2014).
- [44] S. Matteoli, E. J. Ientilucci, and J. P. Kerekes, "Comparison of radiative transfer in physics-based models for an improved understanding of empirical hyperspectral data," IEEE Workshop on Hyperspectral Image and Signal Processing: Evolution in Remote Sensing - WHISPERS (2009) [doi: 10.1109/WHISPERS.2009.5288986].
- [45] J. R. Schott et al., "An Advanced Synthetic Image Generation Model and its Application to Multi/Hyperspectral Algorithm Development", Canadian J. Remote Sens., 25, 99-111 (1999).
- [46] K. Chandra and G. Healey, "Using coupled subspace models for reflectance/illumination separation," IEEE Trans. Geosci. Remote Sens. 46(1), 284-290 (2008) [doi: 10.1109/TGRS.2007.909096].
- [47] L. S. Bernstein, S. M. Adler-Golden, R. L. Sundberg and A. J. Ratkowski, "Improved reflectance retrieval from hyper- and multispectral imagery without prior scene or sensor information," Proc. SPIE, Remote Sensing of Clouds and the Atmosphere XI 6362, 63622P (2006) [doi: 10.1117/12.705038].

- [48] L. S. Bernstein, S. M. Adler-Golden, R. L. Sundberg and A. J. Ratkowski, "Inscene-based atmospheric correction of uncalibrated VISible-SWIR (VIS-SWIR) hyper- and multispectral imagery," Proc. SPIE, Remote Sensing of Clouds and the Atmosphere XIII 7107, 710706 (2008) [doi: 10.1117/12.808193].
- [49] L. S. Bernstein, S. M. Adler-Golden, X. Jin, B. Gregor and R. L. Sundberg, "Quick Atmospheric Correction (QUAC) Code For VNIR-SWIR Spectral Imagery: Algorithm Details," IEEE 4th Workshop on Hyperspectral Image and Signal Processing (WHISPERS), 1-4 (2012).
- [50] D. G. Hadjimitsis and C. Clayton, "The application of the covariance matrix statistical method for removing atmospheric effects from satellite remotely sensed data intended for environmental applications," Proc. SPIE, Remote Sensing for Environmental Monitoring, GIS Applications, and Geology VII 6749, 674936 (2007) [doi: 10.1117/12.751887].
- [51] V. Nardino, F. Martelli, P. Bruscaglioni, G. Zaccanti, S. Del Bianco, D. Guzzi, P. Marcoionni and I. Pippi, "McCART: Monte Carlo Code for Atmospheric Radiative Transfer," IEEE Trans. Geosci. Remote Sens. 46(6), 1740-1752 (2008) [doi: 10.1109/TGRS.2008.916464].
- [52] H. E. M. Viggh and D. H. Staelin, "Spatial surface prior information reflectance estimation (SPIRE) algorithms," IEEE Trans. Geosci. Remote Sens. 41(11), 2424-2435 (2003) [doi: 10.1109/TGRS.2003.818874].
- [53] H. E. M. Viggh and D. H. Staelin, "Surface reflectance estimation using prior spatial and spectral information," IEEE Trans. Geosci. Remote Sens. 45(9), 2928-2939 (2007) [doi: 10.1109/TGRS.2007.898497].
- [54] P. W. T. Yuen and G. Bishop, "Enhancements of target detection using atmospheric correction preprocessing techniques in hyperspectral remote sensing," *Proc. SPIE, Military Remote Sensing* 5613, 111-118 (2004) [doi: 10.1117/12.578688].
- [55] A. Stewart, R. Bauer and R. Kaiser, "Performance Assessment of Atmospheric Correction Algorithms on Material Identification for VIS-SWIR Hyperspectral Data II," Proc. SPIE, Imaging Spectrometry VI 4132, 206-217 (2000) [doi: 10.1117/12.406589].
- [56] J. P. Kerekes, "Error analysis of spectral reflectance derived from imaging spectrometer data," Geosci. Remote Sens. Symposium Proceedings, 1998 5, 2697-2701 (1998) [doi: 10.1109/IGARSS.1998.702323].
- [57] A. Kruger, R. Lawrence and E.C. Dragut, "Building a terabyte NEXRAD radar database for hydrometeorology research," Comp. & Geosci. 32, 247–258 (2006). [doi:10.1016/j.cageo.2005.06.001].
- [58] P. Padovani, M.G. Allen, P. Rosati and N.A. Walton, "Discovery of optically faint obscured quasars with Virtual Observatory tools," Astronomy & Astrophysics 424, 545–559 (2004). [doi: 10.1051/0004-6361:20041153].
- [59] P. Tsalmantza, E. Kontizas, L. Cambrésy, F. Genova, A. Dapergolas and M. Kontizas, "Luminous AGB stars in nearby galaxies: A study using virtual observatory tools," Astronomy & Astrophysics 447, 89–95 (2006). [doi: 10.1051/0004-6361:20053142].
- [60] A.J.G. Hey and A.E. Trefethen, "The Data Deluge: An e-Science Perspective," in F. Berman, G.C. Fox and A.J.G. Hey (Eds.) Grid Computing - Making the Global Infrastructure a Reality, Wiley and Sons, New York, pp 809-824 (2003).
- [61] Common Data Format (CDF), Goddard Space Flight Center Space Physics Data Facility, http://cdf.gsfc.nasa.gov (2015).
- [62] The HDF Group, "Welcome: Solutions to Data Challenges," http://www.hdfgroup.org (2015).
- [63] WMO (World Meteorological Organization), "Guide to the WMO Table Driven Code Form Used for the Representation and Exchange of Regularly Spaced Data in Binary Form: FM 92 GRIB Edition 2", World Meteorological Organization Technical Report, Geneva, 99 pp. (2003).
- [64] OCG, Open Geospatial Consortium, Inc., http://www.opengeospatial.org (12 Aug 2009).
- [65] Y. Bai, L. Di, Y. Wei, "A taxonomy of geospatial services for global service discovery and interoperability," Comp. & Geosci. 35, 783-790 (2009) [doi:10.1016/j.cageo.2007.12.018].
- [66] The Global Earth Observation System of Systems (GEOSS), Group on Earth Observations, Geneva, http://earthobservations.org (2015).
- [67] W. Yang, M. Min, Y. Bai, C. Lynnes, D. Holloway, Y. Enloe and L. Di, "Operational interoperable web coverage service for Earth observing satellite data: issues and lessons learned," Eos Transactions of the American Geophysical Union 89(53), Fall Meeting Supplement, Abstract IN11D-07.

- [68] J. Powell, G. Erukulla, M. Buhisi, and B. Velauthapillai, "A relational database management system for atmospheric compensation research," Comp. & Geosci. 37, 588–597 (2011) [doi:10.1016/j.cageo.2010.04.015].
- [69] M.L. Nischan, J. P. Kerekes, and J.E. Baum, "Analysis of HYDICE noise characteristics and their impact on subpixel object detection," Proc. SPIE 3753, Imaging Spectrometry V (1999) [doi: 10.1117/12.366274]
- [70] P.N. Slater, R.W. Basedow, W.S. Aldrich, and J.E. Coiwell, "In-flight radiometric stability of HYDICE for large and small uniform reflectance targets under various conditions", SPIE Proc., 2821, 300-310 (1996) [doi:10.1117/12.257178]
- [71] SVC, Spectra Vista Corporation, "Field Portable Spectroradiometers," http://www.spectravista.com/ground.html (2015).
- [72] ASD Incorporated, "FieldSpec Portable Spectroradiometer," http://www.asdi.com/products/fieldspec-spectroradiometers (2014).
- [73] Labsphere, "Spectralon Targets," http://www.labsphere.com/products/reflectancestandards-and-targets/reflectance-targets/spectralon-targets.aspx (2015).
- [74] C.C. Homer, L. Huang, B. Yang, B. Wylie, and M. Coan, "Development of a 2001 National Landcover Database for the United States," Photogrammetric Engineering and Remote Sensing, 70 (7), 829-840 (2004).
- [75] G.T.Trewartha, L.H. Horn, An introduction to climate, McGraw-Hill, New York (1980).
- [76] B. Baker, H. Diaz, W. Hargrove, and F. Hoffman, "Use of the Köppen–Trewartha climate classification to evaluate climatic refugia in statistically derived ecoregions for the People's Republic of China," Climatic Change 98, 113–131 (2010) [doi: 10.1007/s10584-009-9622-2].
- [77] EXELIS Visual Information Services, ENVI Environment for Visualizing Images, http://www.exelisvis.com/docs/using_envi_Home.html (2015).
- [78] J.H. Powell and R.G. Resmini, "A spectral climatology for atmospheric compensation," Proc. SPIE 9088, Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery XX (2014) [doi: 10.1117/12.2050596].
- [79] C.-I. Chang, Hyperspectral Imaging, Techniques for Spectral Detection and Classification, pp 20-21, Kluwer Academic Publishers, New York (2003).

- [80] N. Keshava, "Distance metrics and band selection in hyperspectral processing with applications to material identification and spectral libraries," IEEE Transactions on Geoscience and Remote Sensing, 42 (7), 1552-1565 (2004).
- [81] A. Berk, G.P. Anderson, and P.K. Acharya, "MODTRAN 5.3.2 User's Manual," Air Force Research Laboratory (2013).
- [82] L.S.R. Froude, "Storm tracking with remote data and distributed computing," Computers & Geosciences 34, 1621-1630 (2008) [doi:10.1016/j.cageo.2007.11.004].
- [83] P. Cornillon, J. Gallagher and T. Sgouros, "OPENDAP: Accessing data in a distributed, heterogeneous environment," Data Science Journal 2, 164-174 (2003) [doi:10.2481/dsj.2.164].
- [84] F. Chevallier, F. Chéruy, N.A. Scott and A. Chédin, "Neural network approach for a fast and accurate computation of the longwave radiation budget," J. of App. Meteorology 37, 1385–1397 (1998) [doi: 10.1175/1520-0450(1998) 037<1385:ANNAFA>2.0.CO;2].
- [85] FLAASH Module User's Guide, Version 4.5, ITT Visual Information Solutions, Boulder, CO, 42 pp (2008).
- [86] Institute of Global Environment and Society (IGES), "Center for Ocean-Land-Atmosphere Studies", http://iges.org/aboutcola.html (2015).
- [87] E. Kalnay, M. Kanamitsu, R. Kistler, W. Collins, D. Deaven, L. Gandin, M. Iredell, S. Saha, G. White, J. Woollen, Y. Zhu, A. Leetmaa, R. Reynolds, M. Chelliah, W. Ebisuzaki, W. Higgins, J. Janowiak, K. Mo, C. Ropelewski, J. Wang, R. Jenne, D. Joseph, "The NCEP/NCAR 40-year reanalysis project," Bulletin of the American Meteorological Society 77(3), 437-470 (1996) [doi: 10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2].
- [88] NASA Goddard Earth Sciences Data and Information Services Center, "Atmospheric InfraRed Sounder," http://daac.gsfc.nasa.gov/AIRS/AIRS/overview/ (2015).
- [89] T. Bray, J. Paoli, C. Sperberg-McQueen, E. Maler, F. Yergeau (Eds.), "Extensible Markup Language (XML) 1.0 Specification," Fifth Edition, World Wide Web Consortium (W3C), http://www.w3.org/TR/2008/REC-xml-20081126/ (2008).
- [90] MathWorks, MATLAB The Language of Technical Computing, http://www.mathworks.com/products/matlab/ (2015).

- [91] D. Griffith, "Matlab Class Wrapper for MODTRAN 5", MATLAB Central File Exchange, http://www.mathworks.com/matlabcentral (2012).
- [92] A. Silberschatz, H.F. Korth, S. Sudarshan, Database System Concepts, McGraw-Hill, New York, 1142 pp (2006).
- [93] Microsoft, "Internet Information Services (IIS) Home," http://www.iis.net (2015).
- [94] NetBeans, "NetBeans IDE 6.7," http://www.netbeans.org (2009).
- [95] M. Parsian, "JDBC Recipes: A Problem-Solution Approach," Apress, Berkeley, 635 pp. (2005).
- [96] J. Heaton, "JSTL: JSP Standard Tag Library," Sams Publishing, Indianapolis, IN, 412 pp. (2003).
- [97] M. Matthews, J. Cole, J. Gradecki, "MySQL and Java Developer's Guide", Wiley Publishing, Inc., Indianapolis, 408 pp. (2003).
- [98] Sun, "Java Platform, Enterprise Edition," http://java.sun.com/javaee/ (2009).
- [99] S.M. Adler-Golden, M.W. Matthew, A. Berk, M.J. Fox, J. Lee and A.J. Ratkowski, "Improvements in aerosol retrieval for atmospheric correction," Geoscience and Remote Sensing Symposium, IGARSS (2008).
- [100] D.S. Wilks, "Statistical Methods in the Atmospheric Sciences," Academic Press, San Diego (1995).
- [101] J.H. Ward, "Hierarchical Grouping to Optimize an Objective Function", J of the American Stat Assoc, 58, 236–244 (1963).
- [102] MySQL, "MySQL Community Downloads," http://dev.mysql.com/downloads/ (2015).

BIOGRAPHY

John H. Powell received a Bachelor of Science degree in Engineering Physics from the University of Colorado, Boulder in 1986. He received his Master of Science in Meteorology and Oceanography from the Naval Postgraduate School in Monterey, California in 1996. He is presently employed as a Principal Scientist at the National Geospatial-Intelligence Agency in Virginia.