

Multimodal Classification of Remote Sensing Images: A Review and Future Directions

This paper provides a review of current methodologies in seven challenging data fusion applications for remote sensing and highlights the promising methods that benefit from the synergy between machine learning and signal processing.

By LUIS GÓMEZ-CHOVA, Senior Member IEEE, DEVIS TUIA, Senior Member IEEE,
GABRIELE MOSER, Senior Member IEEE, AND GUSTAU CAMPS-VALLS, Senior Member IEEE

ABSTRACT | Earth observation through remote sensing images allows the accurate characterization and identification of materials on the surface from space and airborne platforms. Multiple and heterogeneous image sources can be available for the same geographical region: multispectral, hyperspectral, radar, multi-temporal, and multiangular images can today be acquired over a given scene. These sources can be combined/fused to improve classification of the materials on the surface. Even if this type of systems is generally accurate, the field is about to face new challenges: the upcoming constellations of satellite sensors will acquire large amounts of images of different spatial, spectral, angular, and temporal resolutions. **In this scenario, multimodal image fusion stands out as the appropriate framework to address these problems. In this paper, we provide a taxonomical view of the field and review the current methodologies for multimodal classification of remote sensing images. We also highlight the most recent advances, which exploit synergies with machine learning and signal processing: sparse methods, kernel-based fusion, Markov modeling, and manifold**

alignment. Then, we illustrate the different approaches in seven challenging remote sensing applications: **1) multiresolution fusion for multispectral image classification; 2) image downscaling as a form of multitemporal image fusion and multidimensional interpolation among sensors of different spatial, spectral, and temporal resolutions; 3) multiangular image classification; 4) multisensor image fusion exploiting physically-based feature extractions; 5) multitemporal image classification of land covers in incomplete, inconsistent, and vague image sources; 6) spatiotemporal multisensor fusion of optical and radar images for change detection; and 7) cross-sensor adaptation of classifiers.** The adoption of these techniques in operational settings will help to monitor our planet from space in the very near future.

KEYWORDS | Classification; fusion; multiangular; multimodal image analysis; multisource; multitemporal; remote sensing

I. Introduction

Earth observation through remote sensing techniques is a research field dealing with signals of various physical nature measured by instruments onboard space and airborne platforms. This type of technology can be exploited in different ways, focusing either on obtaining quantitative measurements and estimations of geo–bio–physical variables, or on the identification of materials by the analysis of the acquired signals. Meeting these objectives is possible because materials in a scene reflect, absorb, and emit electromagnetic radiation depending on their molecular composition and shape. Remote sensing exploits these physical facts and deals with the acquisition of information about a scene (or specific object) at a short, medium, or long

Manuscript received November 16, 2014; revised June 18, 2015; accepted June 20, 2015. Date of publication August 7, 2015; date of current version August 20, 2015. This work was supported in part by the Generalitat Valenciana under Project GV/2013/079, by the Swiss National Science Foundation under Grant PPOOP2_150593, by the Spanish Ministry of Economy and Competitiveness (MINECO) under Project LIFE-VISION TIN2012-38102-C03-01, and by the Italian Space Agency under Projects ID-2181 (COSMO-SkyMed: Announcement of Opportunity) and “OPERA-Civil protection from floods.”

L. Gómez-Choiva and G. Camps-Valls are with the Image Processing Laboratory (IPL), Universitat de València, E-46980 Paterna, València, Spain (e-mail: chovago@uv.es; gcamps@uv.es).

D. Tuia is with the Department of Geography, University of Zürich, CH-8057 Zürich, Switzerland (e-mail: devis.tuia@geo.uzh.ch).

G. Moser is with the Department of Electrical, Electronic, Telecommunications Engineering and Naval Architecture (DITEN), University of Genoa, I-16145 Genova, Italy (e-mail: gabriele.moser@unige.it).

Digital Object Identifier: 10.1109/JPROC.2015.2449668

0018-9219 © 2015 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.
See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

distance from the target [1]–[3]. Among all the products that can be obtained from the acquired images, classification maps¹ are perhaps the most relevant ones [4], [5]. The remote sensing image classification problem is very challenging because land-cover and land-use maps are mandatory in multitemporal studies and constitute useful inputs to processes such as the modeling of climate change, the study of oceanic currents, arctic studies, or postcatastrophe deployments.

Besides the variety of applications and the need for accuracy, remote sensing also faces constraints of temporal updating of maps, which forces the analyst to recur to multiple sensors at different spatial and spectral resolutions and coverage. In this context, remote sensing image classification requires methods capable of combining information from multiple sensors with different spatial, spectral, and temporal resolutions to obtain image products with improved overall characteristics [6], [7].

The need for multimodal remote sensing image fusion is further exacerbated by the upcoming constellations of satellite sensors that will acquire a large variety of heterogeneous images of different spatial, spectral, angular, and temporal resolutions. In fact, we are witnessing a tremendous increase in the amounts of data gathered with current and upcoming Earth observation (EO) satellite missions, such as the ESA Sentinels and the still growing NASA A-Train satellite constellation. With the superspectral Copernicus Sentinel-2 [8] and Sentinel-3 missions [9] as well as the planned EnMAP [10], HypsIRI [11], and PRISMA [12] imaging spectrometer missions, an unprecedented data stream for land monitoring will become available in the next years. In this context, two major data processing challenges are expected. On the one hand, big data problems, induced by processing several remote sensing modalities simultaneously, motivate the implementation of the algorithms within distributed computer resources, which should lead to a reasonable computing time. On the other hand, this data flux will require classification techniques, which should be accurate, robust, reliable, and fast. Besides, and very importantly, they should exploit all the heterogeneous and structured pieces of information jointly.

Multimodal remote sensing image fusion involves the terms “modalities” and “fusion.” In our context, the different modalities are mostly linked to different physical quantities describing the observed scenes. These quantities are obtained by sensors issued from multiple measuring principles and technologies, e.g., spectrometry or synthetic aperture radar (SAR). An increasing number of Earth observation satellites acquire information from the observed scenes at completely different spatial, spectral, and

temporal resolutions. Moreover, depending on whether the particular application is defined at a local scale or at a global scale, the environmental conditions (surroundings, environment, climatology, location) can be also considered as different modalities of the same problem. Therefore, even restricting ourselves to the classification of remote sensing images, the number and types of modalities are potentially very high. The second term is fusion, which describes how the different modalities are used together in the system: in remote sensing image classification, fusion can be performed at the feature level [which features are to be fed into the classifier(s) and how], at the classifier level (one per modality or multimodal classifiers), or at the output level (how to combine different model outputs, that can be either ensembles or reparametrizations of the multimodal model).

The remote sensing community has been very active in the last decade in proposing methods that combine different modalities [7], [13], [14]. In addition, to foster the research on this important topic, the Image Analysis and Data Fusion Technical Committee (IADFTC) of the IEEE Geoscience and Remote Sensing Society (GRSS)² has been annually proposing a Data Fusion Contest since 2006 (see dedicated paper in this issue [15]). More recently, data fusion challenges have also been proposed by the ISPRS society for photogrammetry and remote sensing.³ All these efforts reflect the high interest and timely relevance of the posed problems.

In this paper, we will review the current methodologies for multimodal classification of remote sensing images. In Section II, we will provide a taxonomical view of the field into **four main families of multimodal image fusion for classification: fusion at subpixel level, pixel level, feature level, and decision level**. As we show in Section III, many techniques and approaches are available for each one of them, but the choice ultimately depends on the application at hand and the desired end product [14], [16]. Besides the great many algorithms presented in the field, **we observed new emergent methodological avenues for multimodal remote sensing data fusion: research has gradually turned to advanced methods drawn from machine learning and signal processing, which we will review in detail in Section IV:** **1)** **multiple kernel learning that performs fusion in implicit high-dimensional feature representations;** **2)** **sparse dictionary learning, where fusion is conducted at the signal or information-theoretic level;** **3)** **Markov modeling that formalizes spatial and multimodal fusion through global minimum energy concepts;** **4)** **deep learning algorithms that easily accommodate to different sources and missing data;** and **5)** **domain adaptation and manifold alignment,** where sources are combined at a geometrical level in a latent space, regardless of the different nature and dimensionality of the sources. All these sound forms of multimodal

¹In the remote sensing community, the term “classification” is often preferred instead of the (perhaps more appropriate) term “semantic segmentation,” which is typically used in machine learning and computer vision. In this paper, we use the term “classification” to describe the process of attributing each pixel (or segment) to a single semantic class, corresponding to a type of land cover.

²<http://www.grss-ieee.org/community/technical-committees/data-fusion/>

³<http://www2.isprs.org/semantic-labeling.html>

fusion permit the combination and interaction between modalities with different levels of sophistication. Then, instantiations of the standard and advanced techniques will be illustrated in Section V through a sequence of case studies of growing sophistication. We end this article in Section VI with some concluding remarks and outline of the future directions in the field.

II. A TAXONOMY FOR MULTIMODAL REMOTE SENSING

In this section, we introduce key concepts related to multimodal fusion in the context of remote sensing image classification in order to provide a taxonomical view of the field. Multimodal image classification can be defined as a problem where one has to decide: which are the required preprocessing steps for each data source (modality); which features are fed into the classifier(s); which training samples are used to select model parameters; and how to combine the individual decisions of the set of trained classifiers in order to obtain the optimal ensemble of classifiers for our particular problem. The modalities to be included in the analysis, the most appropriate fusion techniques, and the assessment approach must be clearly identified for an efficient fusion [6], [17], [18]. Below, we detail these three elements.

A. Multimodal Data From Multiple Sources

In the field of remote sensing image analysis, the different modalities often represent completely different data entities carrying complimentary information about the surface observed [19]. A typical example is represented by the complementarity between multispectral (several-to-many observations along the spectrum, sensitive to sun illumination and cloud cover, easily photo-interpretable) and synthetic aperture radar (one or a few polarimetric observations, insensitive to sun illumination and almost insensitive to cloud cover, more difficult visual interpretation). Such a complementarity has proven to be very fruitful in many applications of land-cover classification, especially because of the improved discrimination capability obtained for several land-cover classes with both data sources [20].

In other words, the modalities can typically be interpreted as different views of the same scene and objects. In those cases, data fusion implies consistent data preprocessing to make the different data modalities “comparable.” Typical preprocessing steps applied to remote sensing data can be divided as follows. First, instrumental errors are corrected and acquired data (digital counts) are usually transformed to calibrated physical radiance units. In the case of optical data [21], at-sensor radiance is usually converted to top-of-atmosphere (TOA) reflectance. This allows us to remove the dependence on particular illumination conditions (day of the year, solar irradiance, and angular configuration) and illumination effects due to rough terrain (cosine correction that usually requires a di-

gital elevation model of the observed region). Moreover, TOA reflectance can be further corrected from atmospheric and bidirectional reflectance effects when enough atmospheric and angular information is available. In the case of thermal instruments, at-sensor radiance can be also converted to equivalent brightness temperature (in Kelvin degrees) to facilitate the interpretation of the observations. In the case of complex SAR images [22], we observe changes suffered by microwaves in amplitude and phase when they interact with the Earth’s surface. Therefore, SAR backscattering intensity provides information about dielectric and structural properties of the target depending on frequency and polarization; and, in repeat-pass SAR interferometry, the interferometric phase measures distance to the targets and temporal stability (coherence). However, SAR data is dominated by the presence of a multiplicative noise-like component called speckle, which comes from the coherent image formation process, that makes it often necessary to apply an additional filtering preprocessing stage [23]. Finally, all modalities to be combined must be coregistered and eventually resampled to the same grid taking into account the geolocation information available or using automatic coregistration techniques. During this phase, if a digital elevation model is available, geometric distortions and topographic effects over rugged terrains are also corrected. Once the data modalities have been normalized, geolocated and coregistered, data fusion can be undertaken, as detailed in the next section.

B. Standard Fusion Approaches

A variety of image fusion schemes have been proposed in the remote sensing literature, concerning multisource data combination and support decision making [19], [24]. Each fusion method is designed for a specific problem, with disparate inputs, processing approach, and outputs. Fusion approaches are usually named depending on the type of modalities that are involved in the remote sensing application [25]. Given a series of k data modalities, a simple taxonomy of fusion problems regarding the most common modalities in remote sensing would be as follows.

- 1) **Fusion at subpixel level.** The k data sets, which usually involve different spatial scales, are fused at subpixel level using appropriate transforms [26], [27].
- 2) **Fusion at pixel level.** Direct fusion of the k data sets at pixel level, which implies that a direct pixel correspondence between different modalities can be established [28].
- 3) **Fusion at feature level.** Feature extraction for all k data sets followed by a fusion at feature level, which can involve both the extraction and selection of more appropriate attributes [29].
- 4) **Fusion at decision level.** Different processing paths for each modality, followed by fusion at the decision level. This assumes that the k outputs can be combined to improve the achieved accuracy [30].

Depending on the peculiarities of the application and the multimodal inputs, these four scenarios might be used individually or combined in an ensemble of classifiers [31].

C. Performance and Quality Indicators

As stressed above, each type of fusion seeks a specific objective. For example, sometimes the only objective is to obtain a visually enhanced fused image with extended information content: in this case, one can rely on visual inspection to assess the quality of fusion. In other cases, fusion is only performed for classification: in this case one has to ensure that the extra effort required to combine complementary heterogeneous data is beneficial enough in terms of performance and accuracy [24], [32], [33]. Therefore, the correct definition of the fusion quality indicators ultimately depends on the application.

In the field of remote sensing data fusion, four main applications dominate: classification (of land cover, land use, etc.) or detection (of changes, targets, objects, etc.), retrieval of bio-geo-physical parameters (temperature, moisture, chlorophyll content, etc.), image pan-sharpening, and image retrieval from big databases. Depending on the application, different figures of merit are designed and optimized. Classification typically focuses on F -scores, area under the receiver operating characteristic (ROC) curve, the estimated Cohen's kappa statistic, and (overall, average, producer's, user's) accuracy scores. Regression typically evaluates accuracy, bias, and goodness-of-fit statistics, but also spatial homogeneity of the retrieved maps [20]. Pan-sharpening products [34] are typically evaluated through spatial-spectral metrics such as the spectral angle mapper (SAM) and the relative dimensionless global error in synthesis (ERGAS), or more perceptually meaningful criteria such as the structural similarity index (SSIM) or Q4 [33], [35], [36]. For product retrieval and search in big collections and repositories the goal is multifaceted: one aims simultaneously to maximize the success in retrieving the target class of images and the ranked list of best matches. The impact of the different modalities is essentially evaluated comparing the results with and without each modality.

III. REVIEW OF MULTIMODAL REMOTE SENSING METHODS

This section is devoted to a review of the vast literature of multimodal remote sensing methods. Five main standard scenarios are described: pan-sharpening, spatial-spectral unmixing, multiangular fusion, multitemporal fusion, and multisensor fusion.

A. Multiresolution Fusion and Pan-Sharpening

A major image fusion problem is to increase the spatial resolution of satellite images by merging their data with information collected at higher resolution. From geostationary sensors with kilometer-wide resolutions to near

polar sensors with up to 25-cm resolution, current missions offer a variety of observation scales. Many systems, such as IKONOS, QuickBird, SPOT 5 HRG, Pléiades, Landsat 7 ETM+, and Landsat 8 OLI, also exhibit multiresolution capabilities. A typical configuration, suggested by design issues about the signal-to-noise ratio of passive sensors [37], includes a set of coarser-resolution multispectral (MS) bands and one additional finer resolution panchromatic (PAN) band spanning approximately the same spectral range of all MS bands together. Since PAN and MS data are acquired by the same platform with the same viewing geometry, accurate image coregistration is expected.

A wide variety of image fusion methods have been proposed in the remote sensing literature [7], [17], [18], [38], [39], being pan-sharpening methods [24], [32], [35], [40] among the most widespread. Pan-sharpening uses the PAN band to increase the spatial resolution of the MS bands through image processing [34]. In general, there is not a big difference between the spatial resolutions of the PAN and MS bands, but this improvement is critical since the fused images are commonly used for photointerpretation. Hence, pan-sharpening is often especially focused on the injection of spatial details in the MS bands to visually enhance color representations [33], [36]. Then, a classical single-resolution classifier can be applied to the pan-sharpened image [41], [42]. Alternately, case-specific multiresolution methods, based, for example, on **Bayesian modeling** [43] or hierarchical **Markov random fields** (MRFs) [44], can be applied. An example of the latter approach is presented in Section V-A.

B. Spatial-Spectral Unmixing

In quantitative applications, in which relevant information is contained in the spectral domain, alternative fusion methods to pan-sharpening are usually selected. In these cases, the main objective is to preserve the valuable spectral information acquired by multispectral and hyperspectral sensors [33]. Recent proposals based on spectral unmixing have shown good properties providing spectrally consistent fused images and alleviating mixing problems in pixels composed of more than one land cover [27], [45]–[48]. These methods describe a mixed pixel as a combination of pure spectra called endmembers, which contribute into different proportions called fractional abundances [49]. These approaches are especially suited when the input images exhibit significantly different spatial resolutions or temporal revisit times [50]. This assumption was used by the spatial and temporal adaptive reflectance fusion model (STARFM) [46], [51] for combining information from Landsat (30-m resolution) and MODIS (250-m to 1-km resolution, more frequent overpass), and by a full family of methods based on [52] for increasing the spatial resolution of MERIS (300 m) by mapping fractional abundances from Landsat through classification [53]–[55]. These methods are particularly interesting examples of fusion as multiple

(spatial, spectral, temporal) information modes are jointly considered. A case study is presented in Section V-B.

C. Multiangular Image Fusion

Recent satellite platforms such as WorldView-2 or CHRIS Proba include onboard control systems to retarget the sensor to a wide range of viewing angles. Images can be acquired over the same region with higher temporal resolution from different observation angles during different satellite passes. However, since they are collected with significantly different acquisition geometries, the radiometry of each image is altered by angular effects and, while the images are acquired during different satellite passes, also by atmospheric and illumination conditions. Dealing with multiangular images has therefore recently become an alive topic of research in remote sensing.

Multiangular data applications consider two main multimodal settings. The first one is to consider the angular images as different views of the same objects, and exploit differences in the angular sequence as additional features [56] that convey information about, for example, the lateral sides of the objects (strong off-nadir angles), roof materials (nadir view) [57], bidirectional light scattering properties [58], [59], or shadowing [56], [60]. The second setting is to consider multiangular images as a series of observations of the same land covers that have undergone spectral distortions to be removed. This setting, closely related to domain adaptation [61], aims at building classification models that, after training on a specific angular view, must be able to predict the land-cover classes in all the remaining angular views [62]. This is relevant, as ground truth is not available for all acquisitions, so an angular-invariant model would improve land-cover maps issued from multitemporal mosaics. An example is shown in Section V-C.

D. Multitemporal Image Fusion

Multitemporal image classification algorithms classify pixels by exploiting the temporal evolution of their intensities. When a labeled image data set is available, supervised classifiers can yield improved performance over unsupervised approaches. Other advantages are the capabilities to detect land-cover transitions and to process multisensor/multisource images [63], and the robustness to atmospheric and light conditions. Many multitemporal supervised methods have been used, including evidence reasoning [64], generalized least squares [65], or neural networks [66]–[68]. Nevertheless, several problems can be identified. First, classifiers are in general sensitive to the high dimension of the input space (curse of dimensionality) when hyperspectral images are used or multisensor features from different times are put together (stacked approach). This issue has been lately alleviated by using support vector machines (SVMs) [69], [70]. Second, classifiers can suffer from false alarms especially when contextual or textural information is not considered. This is important as, in practice, the user is ultimately interested

in detecting precisely the position and spatial extent of the class(es) of interest. Spatial multitemporal SAR classification of urban areas has been addressed with statistical and neural approaches [71], and at a feature and pixel levels [72]. Third, and very importantly, most methods do not consider the (potentially nonlinear) cross information among pixels and features at different times. The learning paradigm is often violated, as the classifier is trained and tested with data coming from different distributions due to differences in atmospheric and light conditions, sensor drifts, etc. To address this problem, several strategies, including hidden MRFs [73] or fuzzy spatial–spectral fusion and transition probabilities [74], have been presented. Last, but not least, in most cases, only two dates are considered to illustrate method capabilities, and thus performance for long-term operational studies is unclear. In [75], a method aimed at long time series and based on kriging-integrated variograms and Gaussian maximum likelihood, was presented. A case study proposing an advanced kernel-based framework [76] for multitemporal image classification is discussed in Section V-E.

E. Multisensor Image Fusion

A well-known example of the complementarity among multiple EO data modalities is represented by optical and radar imagery. Through optical sensors, the incoming radiation, indirectly related to the reflectance, emittance, and temperature of the observed surface, is measured in the visible and near-to-thermal infrared wavelength ranges [37]. A radar imager for EO, most often implemented through the SAR technology to optimize spatial resolution, transmits a microwave pulse toward the target area and collects the backscattered return [22]. Unlike optical cameras, SAR provides information about roughness and soil moisture, operates regardless of sun illumination, and is almost insensitive to cloud cover [22]. However, visual photointerpretation is more difficult with SAR than with optical imagery, and automated analysis of SAR amplitude data is made complicated by the speckle phenomenon, which is due to radar backscattering and acts similarly to a multiplicative noise [22]. Joint optical-SAR data classification has been addressed for a couple of decades and many methodological approaches, including statistical pattern recognition [23], [77], neural networks [78], decision fusion [79], evidence theory [80], kernel-based learning [76], [81], and MRFs [82], have been proposed. An example is shown in Section V-F.

Light detection and ranging (LiDAR) instruments also constitute an interesting active sensing technology. LiDAR sensors provide relevant height information that can be used to derive digital surface models. The synergistic use of optical and LiDAR data has become a hot research topic in recent years [83]. In most applications, the topographical information of the scene extracted from LiDAR data is used to increase the discrimination of land-cover classes with similar spectral characteristics [20].

IV. ADVANCES IN MULTIMODAL REMOTE SENSING

Multimodal remote sensing data fusion is far from being a dead field of investigation. Dealing with data coming from different sensors, at different resolutions, and with different physical meanings is at the core of the future of remote sensing image processing, for which synergies with compressed sensing, machine learning, and computer vision are giving a new momentum to multimodal remote sensing. In this section, we briefly review these new research directions.

A. Feature Level Fusion Through Multiple Kernels

Kernel methods have been extensively studied in remote sensing [84], since they offer a general and modular framework to encode data into classification and prediction models, to reduce dimensionality or, as for the interest of this paper, to combine different data modalities. First attempts to use kernels for combining data of different nature are found in [81], where Camps-Valls *et al.* design a composite kernel by weighted summation of spectral and spatial features issued from the same, i.e., coregistered, geographical region. Each kernel is evaluated in the space of a specific modality and represents a measure of similarity between the same pixels in that space. Rooted on functional analysis, valid kernels can be summed or multiplied by positive scalars while still being a valid kernel function. Therefore, kernels of different data modalities can be easily combined, thus providing a multimodal data representation. Proceeding this way has the limitation of requiring coregistered data. This reasoning was brought into multitemporal scenarios, where difference kernels were applied for multitemporal processing or change detection. Extending this reasoning to ensembles of kernels (more than two sources) was possible through multiple kernel learning [85], which was used in multimodal studies for combining spectral and spatial information [86], to combine optical and radar data [76], [87], data from the same satellite but completely different locations [88], or optical data from different satellites [89] in change detection. In all these examples, each data source is used to generate a kernel matrix and all the source-specific kernel matrices are then combined by linear combination into a multimodal similarity matrix. In Section V-E, we show an example of how kernel machines allow us to fuse in a very natural way different sources of information in a multitemporal image classification problem [76].

B. Sparse Dictionary Learning

The field of learning sparse dictionaries has also captured the attention of remote sensing data fusion practitioners. In particular, algorithms for learning dictionaries have been applied recently to the problem of image pan-sharpening following restoration-based approaches [90]–[93]. Essentially, these methods rely on the concept of compressed sensing, by which pan-sharpening can be cast

as a signal restoration problem with sparsity regularization. The methods essentially reconstruct the fused multispectral image from the obtained sparse coefficients [e.g., using an orthogonal matching pursuit (OMP) algorithm] and the dictionary of the high-resolution multispectral image. The main directions investigate techniques to select the dictionary, improve the sparse codes, and reduce its dimensionality. In order to solve the problem of learning the dictionary for the high-resolution multispectral images, which is unknown, a recent work [94] proposes to train the dictionary using both the high- and low-resolution multispectral images and to constrain the learning of the K-SVD method to improve the representation. The work in [95] also introduces a compressed sensing method that significantly minimizes the spectral distortion in the pan-sharpened multispectral bands with respect to the original ones. Finally, it is worth mentioning the work in [96], in which pan-sharpening is conducted with a Bayesian nonparametric dictionary learning model that infers dictionary size, patch sparsity, and noise variances, and includes meaningful (piece-wise smoothness and sharpness) image constraints.

Beyond pan-sharpening, sparse learning techniques have been also used to fuse spatial, spectral, and temporal modalities. A technique based on sparse matrix factorization has been recently proposed to fuse remote sensing imagery with different spatial and spectral properties such as Landsat 7 Enhanced Thematic Mapper Plus (ETM+) and Terra Moderate Resolution Imaging Spectroradiometer (MODIS) [50]. Other works fuse spatio-temporal information with sparse techniques via dictionary learning generated from the high- and low-resolution difference image patches across times [97], [98].

C. Structured Prediction and Markov Modeling

Structured prediction and learning is also a relatively recent avenue for remote sensing image processing; relatively, since Markovian modeling has been used in multi-source remote sensing classification [82], [99] since the 1990s. At the same time, Markov random field (MRF) models are being used more and more for tasks such as image reconstruction [100] or super resolution using multiple satellite image acquisitions [101], which are basically estimation problems. MRFs encode naturally neighborhood relations and can be used efficiently to generate classification results with energy functions enforcing spatial smoothing [102], or directional structures [103] and so on. In recent works, MRFs with several connected layers have been used for change detection, either in bitemporal settings [104]–[106] or multiday [106], [107]. In MRF models, a generative model is used to estimate the posterior probability of class assignment and an interaction term is used to encode spatial relations between labels. A new breath for random field models also came from two recent technical advances: first, the availability of tractable, yet accurate methodologies to solve the related energy minimization problems (among them, graph cuts have been the

most successful ones [108]); and second, the introduction of **conditional random field (CRF) models, which are discriminative models where the data distribution interplays with the labels in the interaction term of the energy function** (with the MRFs, only label consistency is involved in the interaction term). By this difference, CRFs are a **flexible framework that allows incorporating all kinds of spatial priors about label structure as well as in terms of image gradients**. In [109], a CRF model is proposed to detect buildings with InSAR and optical images: both data sources are used in a common feature vector and bounded ratios are used in the interaction term the authors exploit. But the interesting part still lies ahead since: 1) semantic logic can be integrated in CRF models, as shown in [110], where the authors improve 3-D scene reconstruction by including class posterior probabilities and spatial logic in the interaction term; 2) by encoding higher order information in the energy function while letting the different modalities play in different terms [111], [112]; and 3) by learning the interaction potentials from data at different scales [113] and using structured prediction to regularize the results obtained with a deep learning model [114]. In Section V-F, we show an example of how a fusion approach based on MRF models allows fusing the information associated with optical data, SAR data, and spatial context at the same time.

D. Deep Learning for Multimodal Image Classification

Deep neural networks (DNNs) are a well-known family of methods for data processing. For example, convolutional neural networks were already proposed in 1989 by the work of Le Cun *et al.* [115], which has today reemerged as a powerful, effective, and flexible paradigm to extract knowledge from data [116]. **In itself, DNNs are not specific to data fusion, but in fact, they can accommodate data from different sources** [117] and, up to some extent, **function in absence of one of the modalities involved** [118], [119]. DNNs are recognized currently as among the most accurate methods in computer vision applications, but there is **little evidence of the good performance in remote sensing image classification**: Vaduva *et al.* [120] introduced a deep learning algorithm for supervised classification of (low-dimensional) very-high-resolution (VHR) images that combines spatial and spectral information, and Mnih and Hinton [121] explored the robustness of deep networks to noisy class labels for aerial image classification. Both contributions consider supervised settings though. Actually, deep nets typically excel when trained with huge amounts of data and in a supervised fashion. For this reason, their application to remote sensing data is still at its beginning [122], [123]. **Recently, some works have introduced unsupervised deep learning for multimodal remote sensing image classification** [124], [125], **hierarchical feature learning and selection** [126], and studied the potential for transferring general object features to the aerial domain [127]. These approaches bring new insights on the remote

sensing data processing problem, also leaving large room for investigation in terms of multimodal data fusion.

E. Alignment of Data Representations for Domain Adaptation

In the previous subsections, we were interested in using different views of the same data to solve a single processing task. Using different views allowed to have a richer description of the process observed. In other scenarios, we are interested in solving a series of similar tasks observed through time. For example, in multitemporal remote sensing, we may want to classify land cover at each time step in order to describe its status or its changes. In large scale classification, we may have a set of images acquired by different sensors observing different regions of Earth. One possible solution to these problems would be to build separate models for each time step (or sensor or location), but this is sometimes not possible, more often because of the lack of ground-truth label information. Moreover, comparing the results becomes difficult (if not impossible), since each map is issued from a separate classification system (with its own specificities and limitations). The general problem depicted here is also known as domain adaptation [61].

An alternative to attack the domain adaptation problem is to reduce all data sources to a common representation. In this common representation, all data sets behave similarly with respect to the criterion of interest: following on the example of classification, samples of the same class will lie close to each other and far from those of the other classes. This solution, also known as manifold alignment [128], consists of finding a set of mappings projecting each input space (each image) into the common space, also called the latent space. The properties of the latent space will vary according to the method, but generally they follow local consistency (neighbors in the input space remain neighbors in the latent space), cross correlation, or class discrimination (in the last case, some labeled examples must be present in all images). In remote sensing, alignment of manifolds of different sensors following these principles has been proposed in [129], where the alignment is performed using proximity graphs, and in [130], where a kernel measure of discrepancy between the image spaces is minimized. Both approaches reduce to solve a generalized eigenproblem expressed in terms of a data similarity and dissimilarity matrices that, depending on the method, can include class similarities, local similarities, and graph Laplacian affinity matrices. Then, the alignment is optimized using the best set of eigenvectors as projectors. For this reason, constraints of coregistration are avoided and the alignment can be found on a subset of the images and then applied to the entire data by simple matrix multiplication. In Section V-G, we show an example of how these techniques can be used to adapt a classifier trained on a given sensor to predict on an image acquired by another sensor with similar resolution, but fewer spectral bands.

V. MULTIMODAL REMOTE SENSING VIA CASE STUDIES

This section illustrates standard and novel methodologies of multimodal remote sensing image classification via seven case studies.

A. Multiresolution Image Fusion

As an example of multiresolution fusion [131] for classification purposes we consider a Bayesian approach based on MRFs and linear mixtures. We tackle the problem of fusing the information associated with PAN data, MS observations, and the spatial context of each pixel [43], [131]. The focus on this example case study is explained by the customary importance of PAN and MS imagery, which are made available by several space missions with VHR payloads, and by the opportunity to use their complementary spatial and spectral resolutions to capture the spatial distribution of land cover. Let the pixel lattices at the PAN and MS spatial resolutions be properly coregistered and aligned. Input MS and PAN data are modeled as realizations of 2-D continuous-valued stochastic processes sampled over the PAN and MS lattices, respectively. The output classification map is also modeled as a realization of a 2-D discrete-valued stochastic process over the PAN lattice. An example of Bayesian approach to multiresolution spatial–contextual fusion can be formalized on the basis of two major assumptions: 1) the stochastic process of MS observations can be obtained by mosaic averaging a 2-D Gaussian process of “virtual” feature vectors that is sampled over the PAN lattice (linear mixture assumption) [43]; and 2) the stochastic process of class labels is an MRF. The former assumption allows the relationship between data at the two spatial resolutions to be formalized in terms of an ideal stochastic process defined at the same resolution as the PAN image, but including the same spectral channels of the MS data.

The MRF assumption basically implies that the probability distribution of the label of each pixel in the PAN lattice, when conditioned to the labels of all other pixels in the same lattice, can be restricted to the distribution conditioned only to a subset of labels of neighboring pixels [132]–[134]. This Markovianity condition, which extends the analogous well-known assumption for 1-D Markov chains [135] to 2-D stochastic processes over discrete lattices, is aimed at formalizing spatial–contextual information (spatial memory) associated with the process. Markov modeling has a major impact on classification because it can be proven, through the so-called Hammersley–Clifford theorem, that, if Markovianity holds, then, under rather mild assumptions, the global Bayesian maximum *a posteriori* (MAP) image classification rule, which would be a computationally intractable problem, is equivalent to the tractable minimization of an “energy” function defined according to the neighborhood system [132]. This analytical result, along with remarkable flexibility in defining the

energy to incorporate multiple information sources [82] (see also Section V-F) and to favor desired label configurations, and with the availability of computationally effective global or near-global energy minimization algorithms (e.g., the aforementioned graph cuts [108]), explains the role of MRF models as mature methodological tools for multisource analysis problems [102].

Based on the aforementioned assumptions, the problem of mapping class labels at the PAN resolution on the basis of realizations of the MS and PAN stochastic processes can be addressed by iteratively combining: 1) a Bayesian MAP estimation of the Gaussian virtual process, given MS and PAN data; and 2) the estimation of the class label process, given the PAN observations and the estimated virtual process that are both available on the PAN lattice, through the minimization of an MRF energy. Details can be found in [131].

Fig. 1 shows an example of application to land-cover mapping from semi-simulated multiresolution data obtained from a 4-m resolution IKONOS image of the area of Itaipu (Brazil/Paraguay border), including seven classes ranging from urban and built-up areas, to vegetated covers and water bodies. Simulated MS and PAN data were obtained by subsampling the IKONOS spectral channels and by averaging them together, respectively. This semi-simulated example allows the results of multiresolution fusion to be compared with those of a benchmark single-resolution classification of the original (full-resolution) multispectral image, and the possible accuracy loss due to the degraded spatial resolution of the MS channels and spectral resolution of the PAN channel to be evaluated. A single-resolution MRF-based classifier was used for this comparison (see [131] for details). Fig. 1 shows details of the maps obtained through multiresolution fusion in the cases of resolution ratios (i.e., ratio between the linear spatial resolutions of the MS and PAN lattices) equal to 2 and 4, respectively. The comparison with the single-resolution benchmark points out remarkable similarity in terms of both visual analysis and classification accuracy on test samples located inside homogeneous image regions. The main difference may consist in slight blocky artifacts that can be noted in the multiresolution result obtained with resolution ratio equal to 4 (but not 2). These artifacts are consistent with the spatial degradation introduced in the larger ratio. Nevertheless, these results overall suggest the effectiveness of Bayesian MRF-based approaches to multiresolution fusion, and illustrate their potential to map land-cover classes at the finest available resolution with accuracies analogous to the ideal result associated with full spectral information at full panchromatic resolution.

The second example is shown in Fig. 2 with regard to a real IKONOS data set acquired over Alessandria, Italy, and composed of a 1-m resolution PAN [Fig. 2(b)] image and four 4-m resolution MS channels [Fig. 2(a)]. The imaged scene is mostly associated with urban and agricultural land covers. In this case, no benchmark comparison with the

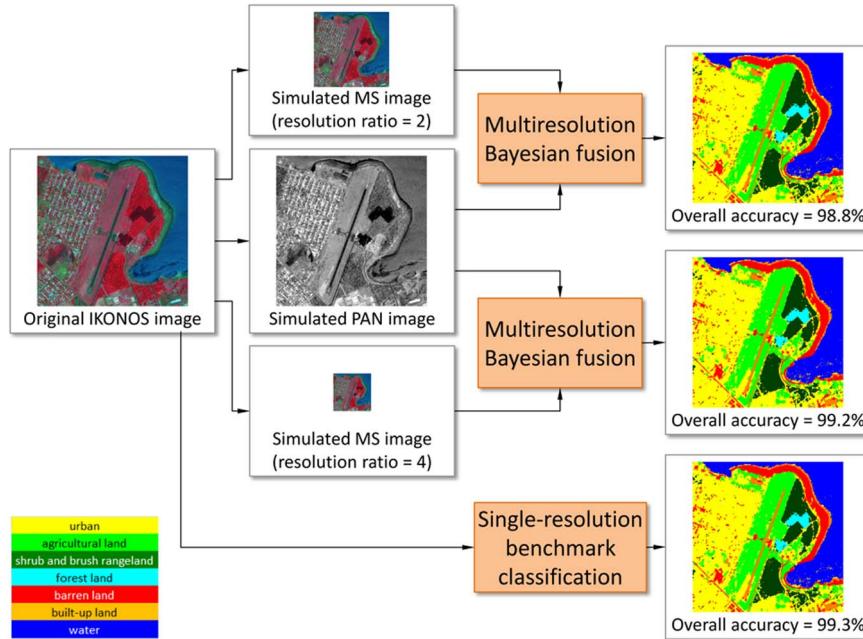


Fig. 1. Example of multiresolution fusion, Itaipu, Brazil/Paraguay border, semi-simulated data set: detail of the IKONOS image (false-color composite) and of the classification maps obtained through multiresolution fusion (resolution ratios equal to 2 and 4) of simulated PAN and MS channels, and through single-resolution classification of the original full-resolution image.

aforementioned ideal single-resolution result is obviously feasible. However, through multiresolution fusion, accurate classification was obtained again on the PAN lattice, both quantitatively (overall accuracy equal to 98.6% on test samples located in homogeneous regions) and visually. It is also worth noting that the resolution ratio was 4;

nonetheless no blocky effect was obtained through multiresolution analysis probably due to the spatial smoothing associated with the MS camera.

B. Multitemporal and Multisensor Fusion and Interpolation

The usefulness of data fusion approaches based on spatial and spectral unmixing (see Section III-B for more details) can be exploited in multitemporal settings involving time series of different sensors. We consider here time series of MERIS (15 bands at 300 m in full resolution) and Landsat TM (30 m) [48]. The final goal is to complete or fill gaps in the Landsat time series, which have a revisit time of 16 days, by using MERIS data that exhibit more frequent temporal coverage of up to three days, hence obtaining consistent time series at high spatial resolution.

The analyzed fusion approach is illustrated in Fig. 3. It assumes a linear mixing model for the lower spatial resolution observations, i.e., each MERIS pixel is unmixed using information about its composition in terms of land-cover class proportions. However, rather than estimating them from the MERIS spectra, class proportions and their spatial distributions are obtained from the high spatial resolution Landsat time series. A high-resolution soft clustering of the Landsat time series is used to estimate land-cover class proportions inside each coarser resolution MERIS pixel by also taking into account the acquisition characteristics of MERIS. Class endmembers are derived from MERIS data through a sliding-window spatial unmixing algorithm, and

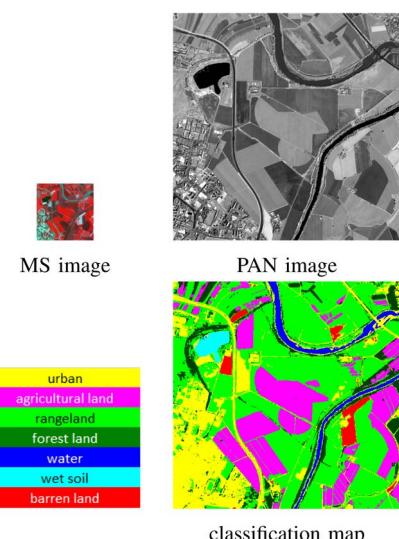


Fig. 2. Example of multiresolution fusion, Alessandria, Italy: finer resolution PAN channel, false-color composite of coarser resolution MS channels, and classification map.

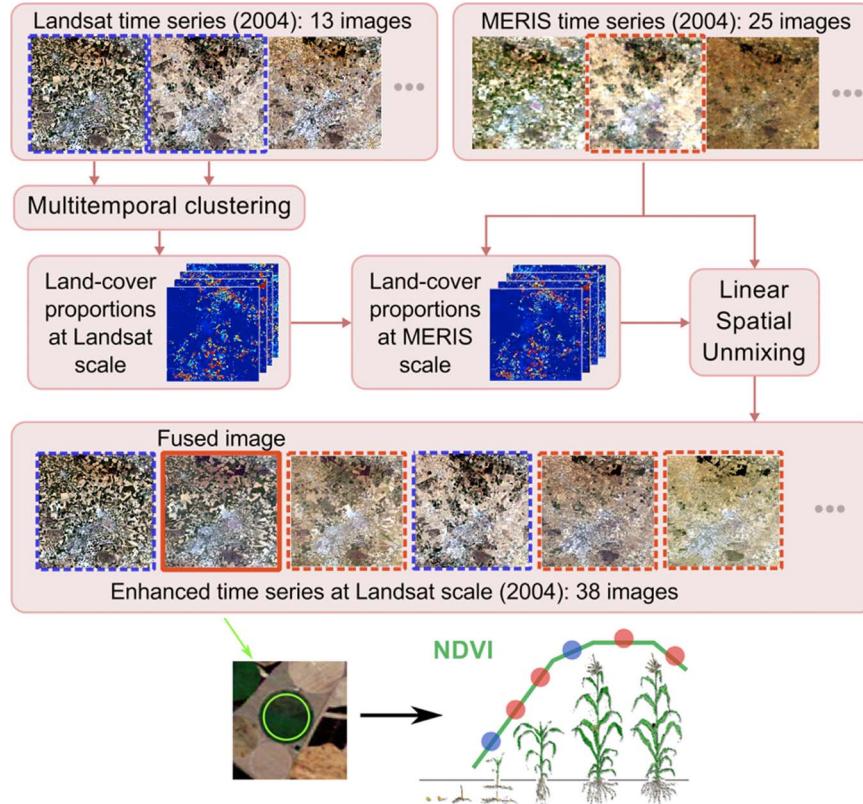


Fig. 3. Multitemporal fusion scheme for medium (MERIS) and high (Landsat TM) spatial resolution image time series [48].

high-resolution fused pixel intensities are obtained as linear combinations of the estimated MERIS endmembers weighted by the Landsat soft memberships. Details can be found in [48]. The enhanced time series allow developing operational applications that require monitoring rapidly varying phenomena at high spatial resolution, such as precision agriculture, irrigation advisory services, and near-real-time change detection. For example, the temporal profile of the vegetation parameters is enhanced by including the downsampled products and allows a more accurate determination of the crop type and phenology.

The availability of long time series allows some of the images to be used for quantitative performance analysis. One of the Landsat images is held out from the fusion algorithm, used as a reference, and compared to the downsampled image for the same date. Here, we show examples of downscaling results in the spatial, spectral, and temporal domains to illustrate the potential and benefits of this multimodal fusion approach.

Fig. 4 shows image products (namely, false-color composites and the normalized difference vegetation index (NDVI) [37]) for MERIS and Landsat images acquired on the same date together with the fused image for this date. The fusion result is visually very similar to the reference Landsat image. In addition, a quantitative comparison of the fused reflectance to the original Landsat and MERIS

images is also shown in terms of root mean square error (RMSE). The largest errors correspond to urban areas, where the MERIS spatial resolution cannot provide the fine spatial detail or the pure pixels that are necessary to capture the different land covers. Similarly, the retrieved vegetation index map shows a Pearson's correlation coefficient of 0.87 and an average RMSE of 0.091 (results not shown in the figure).

We note that a fine temporal sampling is essential in agricultural monitoring for an accurate determination of crop type and phenology. Fig. 5 shows the NDVI temporal profile for three different agricultural fields. The profiles are drastically improved when the downsampled MERIS images are included in the time series, thus allowing a deeper analysis of the related phenological cycles. Crop "A" corresponds to a summer cereal crop, whose cycle is characterized by growth and development stages in spring followed by a drying phase. Crop "B" corresponds to an early spring crop that is harvested in June. A sudden change in its NDVI profile cannot be properly predicted using only the Landsat time series. Crop "C" shows a double-cropped irrigated field, in which a second crop is planted after the first one is harvested. As noted for crop "B," when the downsampled MERIS images are used, a more accurate temporal NDVI profile is obtained, which allows a better monitoring and classification of the agricultural crops.

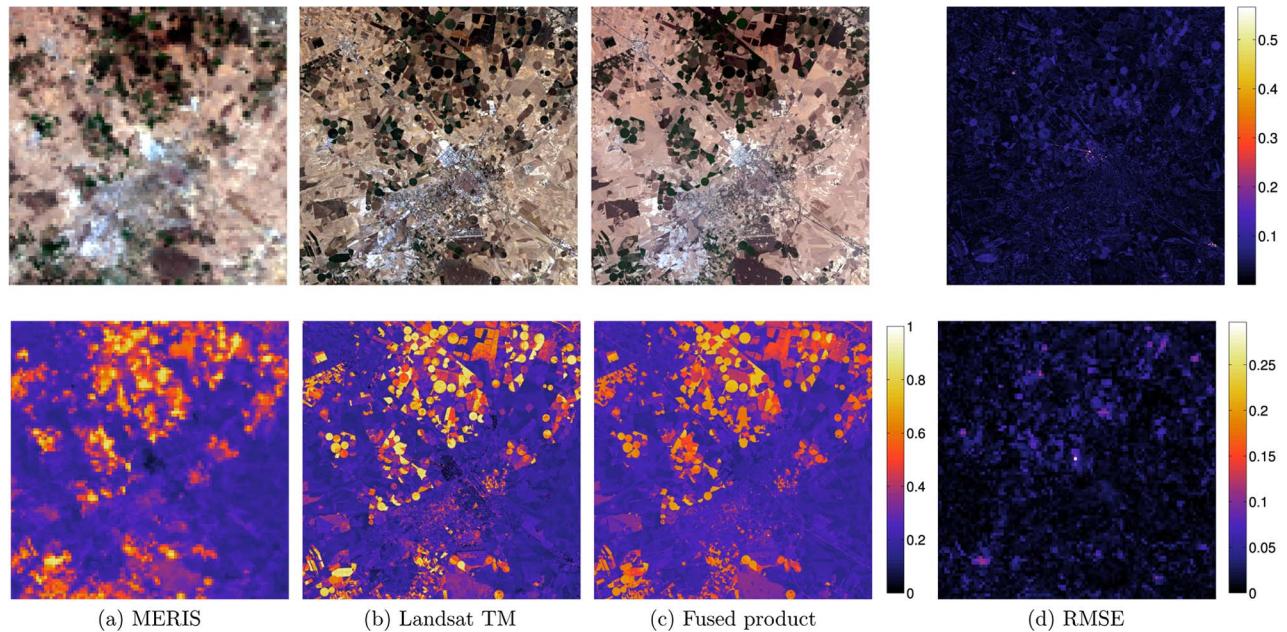


Fig. 4. False-color composition (top row) and NDVI (bottom row) for the original (a) MERIS and (b) Landsat TM acquisitions, (c) the downscaled product, and (d) the RMSE values between the Landsat and the fused product at Landsat (top rightmost panel) and MERIS (bottom rightmost panel) scales.

C. Multiangular Image Fusion and the Spatial Distortions Problem

In this example, we investigate the accuracy of a supervised classifier when adapting on a multiangular sequence acquired over Rio de Janeiro, Brazil, by the WorldView-2 satellite in January 2010 [136]. The sequence depicts a 5-min-long acquisition over a single pass of the satellite and consists of 20 images with off-nadir angles θ going from -47.3° to $+47.5^\circ$ (looking southward and northward, respectively, from the satellite to the imaged area). Each image contains eight spectral bands covering the visible and near-infrared (NIR) spectral ranges (see the right column of Table 2 for the wavelengths of each band). The flight

path is perpendicular to the sun (see Fig. 6), which causes a symmetric scattering/shadowing behavior along the sequence. A ground truth of five classes has been derived by photointerpretation on each image separately. The most nadir image and the corresponding ground truth are shown in Fig. 7.

The experiment is designed as follows. Three images are considered as those providing training samples; each one corresponds to a different column of Fig. 8. Each image in the data set is used in turn as the testing image. The x-axis in the panels of Fig. 8 represents the angular distribution, so that each x-location in these plots corresponds to one of the 20 testing images. Accuracy is computed on the entire

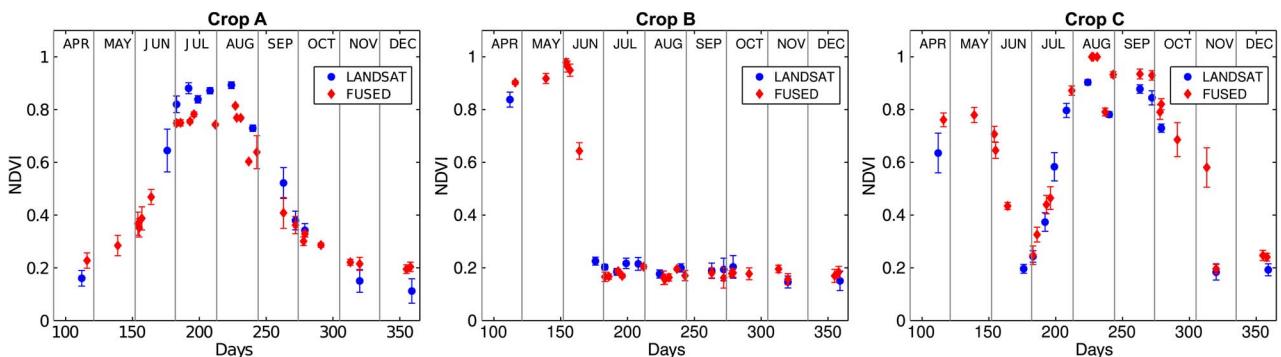


Fig. 5. Temporal profiles of NDVI values (mean and standard deviations within the crop fields) derived from the Landsat-5 TM (blue dots) and the downscaled images (red diamonds).

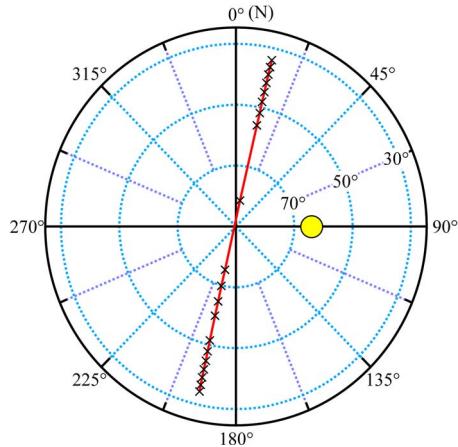


Fig. 6. Viewing geometry of the 20 considered multiangular acquisitions. Each one is represented by a cross along the flight line (in red) on an azimuth (North = 0°) versus the elevation angle (increasing radially from the center, so the angle is 90° minus the value in the graph) plot. The yellow circle represents the position of the sun [136].

ground truth of the testing image. Two types of compensation of angular effects are considered: one is physical (the DG-Acomp proprietary compensation algorithm) and the other one is statistical (based on nonlinear data transformation through kernel principal component analysis (KPCA); details can be found in [137]). The portabilities obtained through linear and nonlinear methods are compared by using both a linear discriminant analysis (LDA) and the SVM with Gaussian kernel function, respectively.

Fig. 8 illustrates the portability results. From the plots, we can draw a series of interesting observations: first, the domain adaptation problem described above is real and visible especially when using a linear classifier that cannot cope with nonlinear distortions of the signals due to the angular conditions. When training the classifiers on the nadir image (central panel), we observe a deterioration of the performances at large viewing angles and an almost symmetrical loss (the sun is perpendicular to the acquisition path, so it impacts the images taken from the leftmost

and rightmost positions in similar ways). When training on off-nadir images (left and right columns), we observe a deterioration of the performances in the central part of the plot, corresponding to the prediction of the nadir images with models trained under strongly off-nadir conditions. Since the angular distortions are symmetric in this case, better predictions are observed at the other end of the angular sequence: this behavior is specific to the angular conditions of this data set, whereas, in more complex scenarios, the classifier performance will decrease even more when predicting data from extremely different angular conditions (see [136] for an example over the city of Atlanta).

We also observe that the tested compensation strategies are very effective in compensating for changes in a viewing angle. Atmospheric compensation is very effective in retrieving an almost flat classification curve, i.e., a classifier trained on a specific image effectively generalizes to data issued from other acquisitions. In the case of the linear classifier, the best results are obtained by the joint use of both statistical and physical adaptation, as KPCA provides the nonlinearity necessary to solve the classification problem, for which the signatures have been adapted by the atmospheric correction. In the case of the nonlinear classifier, the KPCA step seems unnecessary, most probably because the classifier is already sufficient to describe the nonlinearities in the data regardless of the domain adaptation problem.

D. Multisensor Image Fusion Exploiting Physically-Based Feature Extraction

This section presents an example of multisensor cloud screening [138] that exploits combined information from the MERIS [139] and AATSR [140] instruments on board the ENVISAT satellite (2002–2012). They had similar spatial resolutions and swath but were complementary in terms of spectral domains and viewing geometries. The main objective is to explore, through information fusion at different levels, their synergistic use to increase cloud detection accuracy. The considered cloud-screening approach is based on a supervised classification methodology that makes use of ensemble classifiers and relies on both real observations and simulated data. The latter are generated using coupled surface and atmospheric radiative transfer models (RTMs). The main idea is to generate, under known conditions, enough situations to cover a wide range of natural scenarios and then use them for developing robust supervised cloud classifiers.

First, meaningful physically-based features (e.g., brightness, whiteness, temperature, atmospheric absorption features) optimizing the separability between clouds and surface are extracted from suitably calibrated data of both sensors [141]. Then, to solve the complex classification problem of cloud screening, these added-value features and the original spectral bands are used as inputs to advanced classifiers, consisting of several artificial neural networks (ANNs) trained with different sets of input features and training samples. An ensemble classifier is

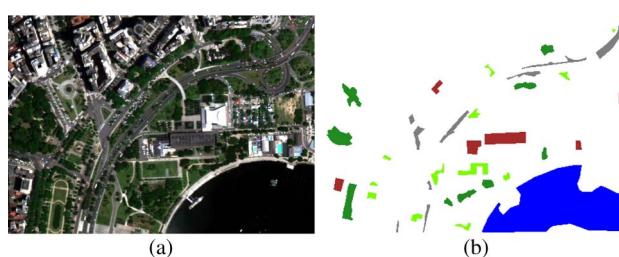


Fig. 7. Example of one among the 20 multiangular images acquired over Rio de Janeiro (most nadir acquisition). Legend: water, roads, meadow, trees, buildings [136]. (a) Image. (b) Ground truth.

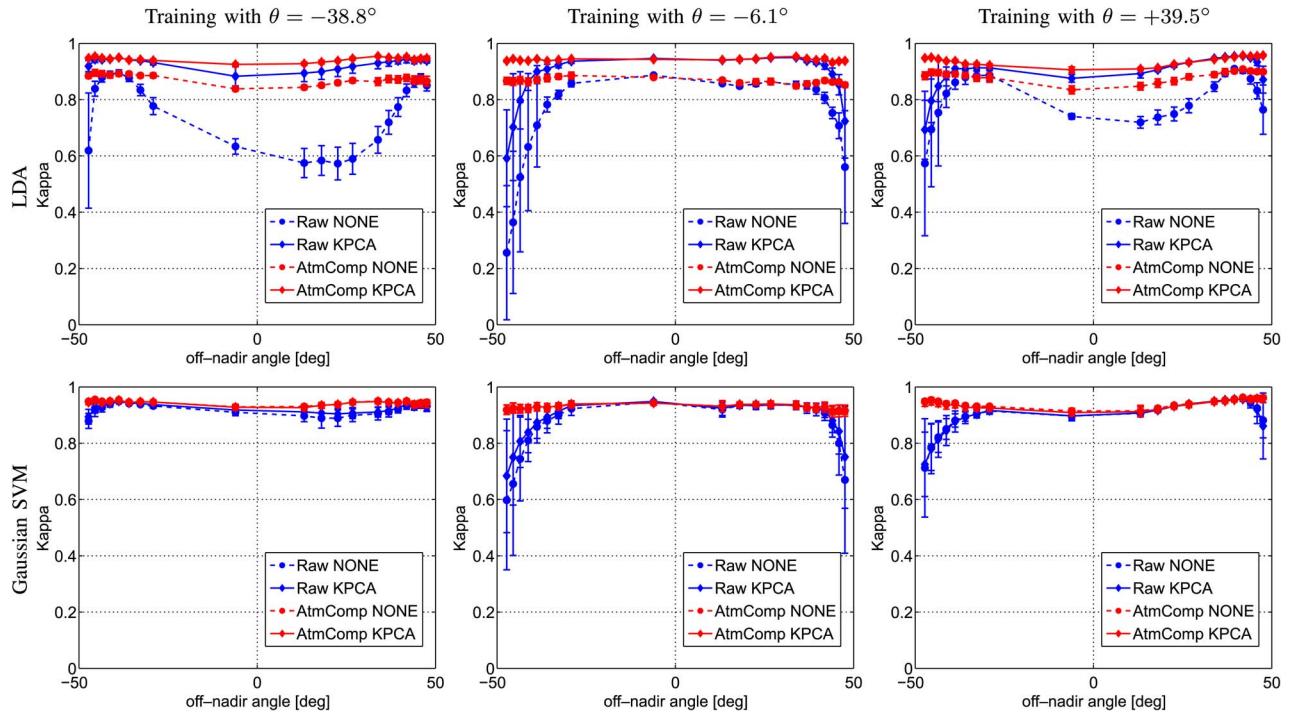


Fig. 8. Model portability results when training classifiers on the image at nadir angles (left) $\theta = -38.8^\circ$, (center) -6.1° , and (right) $+39.5^\circ$. Top and bottom rows illustrate results obtained with LDA and Gaussian SVM, respectively. Accuracy is expressed through the kappa statistic agreement score [37].

obtained by combining the ANN outputs at the decision level. Therefore, fusion at the pixel, feature, and decision levels as well as multiple sources, sensors, and spectral ranges are involved in the overall process.

The synergistic use of AATSR and MERIS data is of paramount importance to improve cloud screening accuracy [138] because: 1) MERIS provides an accurate spectral characterization of the target reflectance; 2) AATSR provides angular information (nadir, forward) and additional spectral bands in the infrared regions; and 3) AATSR and MERIS channels measuring atmospheric absorptions provide information about cloud height. Fig. 9 shows the locations of the MERIS and AATSR channels, as compared to standard spectral curves of healthy vegetation, bare soil, and atmospheric transmittance. The spectral bands free

from atmospheric absorption contain information about surface reflectance, while others are mainly affected by the atmosphere. The above mentioned features together with the MERIS and AATSR spectral bands are fed to the ANNs.

Obtaining training samples with a true label (“cloudy” or “cloud-free”) is not an easy task since coincident simultaneous cloud data at the same resolution are not available for MERIS and AATSR. Therefore, coupled surface and atmospheric RTMs were used to simulate satellite observations associated with a wide range of real “cloudy” and “cloud-free” scenarios. However, models relying only on simulated data can provide poor results when applied to real data depending on the quality and representativeness of the simulated information and of the related noise models. Therefore, for training the models, the simulated samples

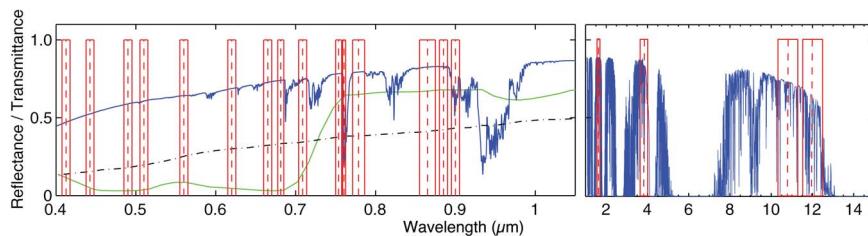


Fig. 9. MERIS (optical) and AATSR (thermal) channels (red bars) superimposed to the atmospheric transmittance (blue solid line), and the reflectance spectra of vegetation (green thin solid line) and bare soil (dashed-dotted line).

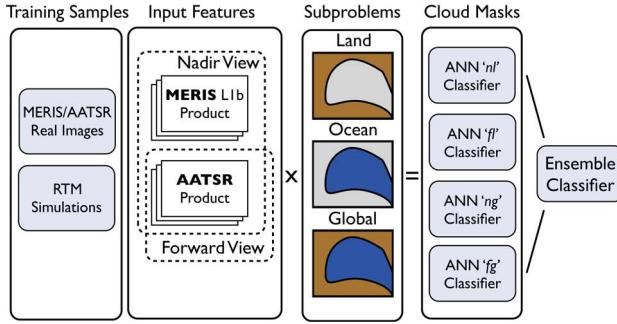


Fig. 10. Scheme of the different fusion levels considered to implement the ensemble of classifiers for cloud screening.

were complemented with real MERIS/AATSR spectra that were manually labeled as cloud contaminated or cloud free.

Finally, the basic idea of ensemble methods is to construct a set of classifiers and combine their individual decisions (e.g., by means of voting or weighted average) with the aim of providing a more accurate classifier [30], [31]. In this case study, a simple ensemble of six ANNs was constructed on the basis of the different modalities involved in the problem by taking benefit from the looking-angle characteristics of AATSR (forward-looking f or nadir n) [140] and from the opportunity to use separate training sets for cloud screening over “land,” “ocean,” or “both.” Indeed, only four independent cloud masks are obtained for a given image since the local “land” and “ocean” networks are selectively applied depending on the pixel location. Therefore, two “local” (N_{nl} and N_{fl}) and two “global” (N_{ng} and N_{fg}) classifiers are obtained for the “nadir” and “forward” views, respectively. Fig. 10 shows the different multimodal information and fusion levels for all processing steps, which perfectly illustrate the number of sources, sensors, data, views, and fusion levels that can be present in a remote sensing application.

A set of coregistered MERIS and AATSR products associated with a number of representative sites worldwide was used to train this multimodal classification algorithm and test its accuracy (38 training and 46 test images, respectively). The performances are compared to those of a single-instrument approach that only uses MERIS data. The classification accuracy is measured by the kappa statistic score (κ) obtained individually for each one of the 84 MERIS/AATSR images. Fig. 11 shows this comparison as a scatter plot of κ values, in which the multimodal fusion (y-axis) clearly outperforms the single-sensor method (x-axis) for almost all the images (points over the 1 : 1 line).

E. Multitemporal Spectral–Spatial Image Classification With Kernels

As discussed in Section IV-A, kernel machines allow to fuse in a very natural way different sources of information. The idea is based on well-known properties of functional

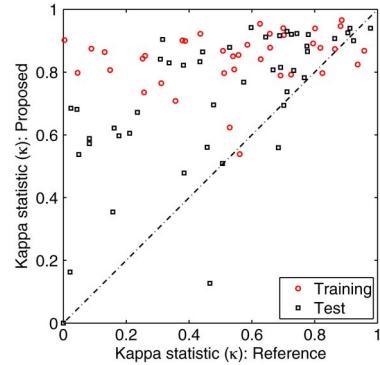


Fig. 11. Cloud classification accuracy (κ) for the 84 MERIS/AATSR (training and test) images obtained using a single-instrument method (cloud probability processor) and the presented multimodality ensemble method.

analysis [142], such as the direct sum of Hilbert spaces, by which a (weighted) sum of kernel functions is still a valid kernel function. The idea was originally introduced for combining spatial and spectral image information in [81], and latter extended to the multitemporal scenario in [76], which is reviewed in this section.

Let us first fix the notation. Let there be a process characterized with P different information sources (e.g., a pixel at T different instants), each of which is given by a (possibly different dimensional) vector $\mathbf{x}^{(t)}$, $t = 1, \dots, T$. The standard approach to combine information considers a stacked vector given by $\mathbf{x}_i = [\mathbf{x}_i^{(1)}, \dots, \mathbf{x}_i^{(t)}]$, which represents the concatenation of all “pixels views” into a single feature vector. Then, a kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$ is built for use in SVM classification. An alternative presented in [76] and [81] considers first mapping the vectors $\mathbf{x}_i^{(t)}$ and $\mathbf{x}_j^{(t)}$ and then stacking them in Hilbert spaces. This operation yields a summation kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{t=1}^T K_t(\mathbf{x}_i^{(t)}, \mathbf{x}_j^{(t)})$. Moreover, multiplying a kernel function times positive scalars $\mu_t > 0$ returns a positive definite kernel, giving rise to the weighted summation kernel $K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{t=1}^T \mu_t K_t(\mathbf{x}_i^{(t)}, \mathbf{x}_j^{(t)})$. More sophisticated concatenations can be formulated (see details in [76] and [81]).

In multitemporal image classification, one tries to classify pixels of an image at the observation time t_0 by using all available (instantaneous and/or previous) information $t \leq t_0$. Two scenarios can be found. The most common one is when labeled samples are available only for $t < t_0$ so we have to learn to discriminate classes at $t < t_0$ and then do inference at $t = t_0$. A more advantageous setting is when some (typically little) labeled information is available for $t \leq t_0$. In both scenarios, we will use the SVM and the one-class SVM (OC-SVM) classifiers, according to whether there is full information on the class labels or only on the class(es) of interest, respectively.

Currently, there are very few operating sensors providing times series of hyperspectral images, but this will be

Table 1 Introduced Changes at the Corresponding Instants

Time	Change	#Pixels
$t_2 \rightarrow t_3$	Grassland (C2) to Bare soil (C4)	1361
$t_4 \rightarrow t_5$	Bare soil (C4) to Rural urban area (C5)	608
$t_6 \rightarrow t_7$	Sallow water (C3) to Sand (C6)	229
$t_8 \rightarrow t_9$	Bare soil (C4) to Winter crops (C7)	1470
$t_{10} \rightarrow t_{11}$	Sand (C6) to Rural urban area (C5)	305

soon a common scenario. We simulated this scenario by simulation of a synthetic time series of hyperspectral images. We used data from the compact high-resolution imaging spectrometer (CHRIS), which is mounted onboard the European Space Agency (ESA) small satellite platform called the Project for on Board Autonomy (PROBA). The CHRIS sensor provides hyperspectral images in the spectral range from 400 to 1050 nm (62 spectral channels for acquisition Mode 1) [143]. The selected image was acquired in the 2006 Agricultural Bio-/Geophysical Retrievals from Frequent Repeat SAR and Optical Imaging (AgriSAR 2006) campaign over the Demmin site (Germany) [144]. This image was selected for the study in order to take into account different surface types, patterns, and spatial textures (soil, vegetation, water, urban areas, etc.). A time series of 12 images was generated by modifying the class covariance matrices in a time fashion, and incorporating natural variability and artificial changes in the scenes along time: 1) natural spectral variability of the class accounted by the covariance matrix and the random generation of the samples for the different dates; 2) changes of the class distributions between dates (e.g., due to illumination or atmospheric effects) simulated with a multiplicative factor over the distribution parameters ($\mu_t = \delta_t \mu$ and $\Sigma_t = \delta_t^2 \Sigma$, where $\delta_t = 0.01t + 0.94$, $t = 1, \dots, 11$ for all classes); and 3) artificially generated changes in the ground truth only included at even instants ($t_3, t_5, t_7, t_9, t_{11}$). The latter allows

us to study the adaptation capabilities of the time-varying kernel classifiers. See Table 1 for details on the introduced changes (more details in [76]).

In order to analyze the performance of the proposed methods under realistic ill-posed situations, we varied the number of training samples per class ($n = \{5, 10, 15, 20, 30, 40, 50\}$) and measured the overall accuracy (OA[%]), the estimated kappa statistic (κ), and the complexity of the machines using the rate of support vectors (in [%]). Classifiers were trained following a tenfold cross-validation method, and the best composite kernels were selected according to the κ score. Average results for a number of ten realizations and over all time prediction instants are shown in Fig. 12. In all cases, the RBF kernel outperformed the linear kernel.

Several conclusions can be derived. First, as we increase the number of training samples, accuracy and sparsity increase. Second, the best kernel classifier is constituted by the weighted kernel, yet results are very similar in accuracy and sparsity to the summation kernel. The stacked kernel approach produces the worst results in all the domains, probably due to the extremely high input space dimensionality generated. It is also important to note that the proposed classifiers obtain close results to the MAP classifier,⁴ even when working with a reduced number of training samples. This can be explained since the composite kernel classifiers consider the temporal information in addition to the static spectral information.

For further comparing the classifiers, we focused on the case of 50 training samples for each class (the last point in the curves of Fig. 12). Fig. 13 shows the sequence of synthetic images, their corresponding true maps, and the classification maps obtained with the optimal MAP

⁴Note that the MAP classification constitutes an upper bound of model performance since the true distribution that generated the synthetic data is used to generate the classifier.

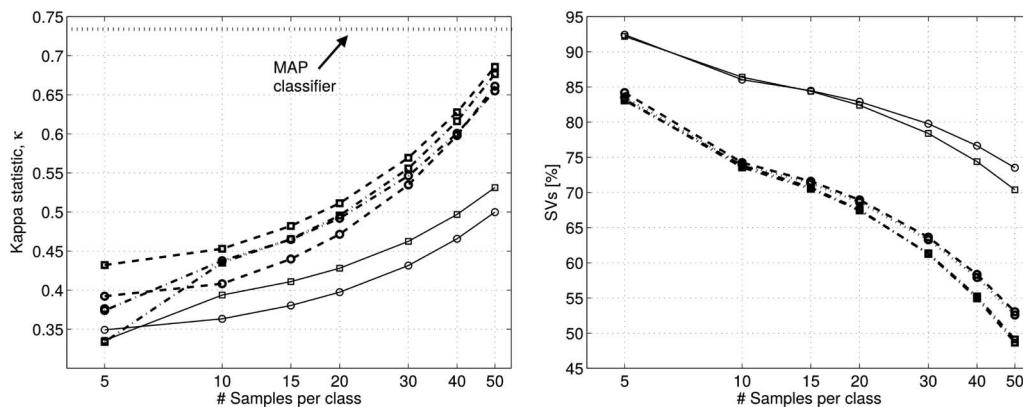


Fig. 12. Results for the multitemporal hyperspectral image classification problem. (left) Kappa statistic and (right) rate of support vectors as a function of the number of training samples per class. Several kernel classifiers are shown using SVMs (squares) and OC-SVMs (circles) with the RBF kernel function: summation kernel (dashed lines), weighted kernel (dashed-dotted lines), and stacked approach (thin solid lines). Average results for all instants and ten realizations are shown for all methods and the MAP classifier.

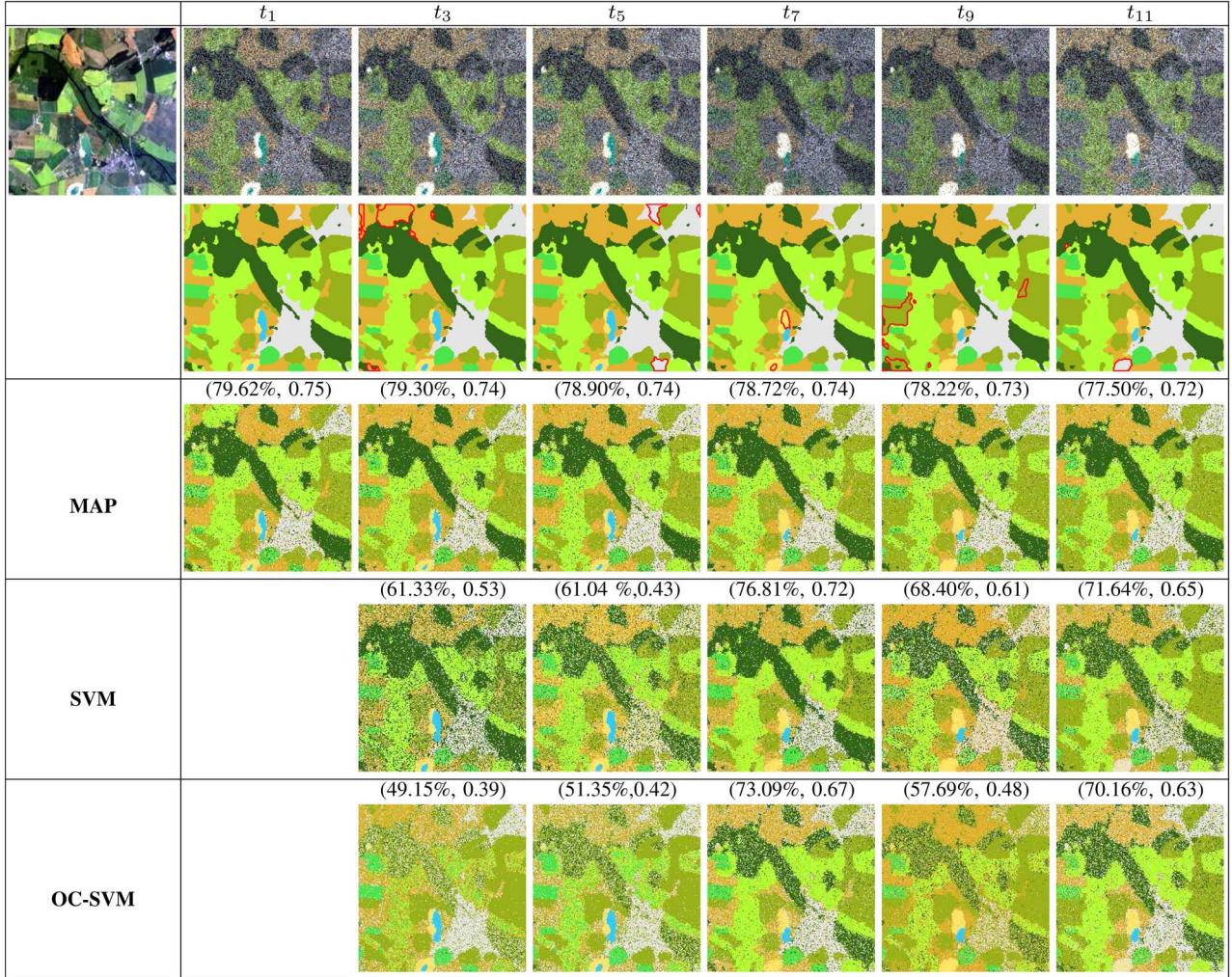


Fig. 13. Sequence of synthetic images (RGB composition, bands [19 12 2]) and their corresponding true classification maps. Only 50 training samples were used for training. Accuracies are indicated on the form (OA[%], κ).

classifier, and the best composite kernel SVM and OC-SVM classifiers at each instant. It can be noted that, again, the SVM works slightly better than the OC-SVM in almost all the cases. This can be due to the fact that distributions vary along time quite smoothly, and thus, pure inductive classifiers can yield good results in this experiment. One can also note that, in general, results are improved when enough temporal data are available to allow classifiers to follow the dynamic changes; accuracy reaches the maximum values starting at time t_7 , getting closer to the MAP solution with relatively little training data. Also note that all the classification maps appear quite noisy, which is a direct consequence of not considering the contextual or textural information. We should stress however that following the same strategy for data fusion with kernels, one can incorporate very easily contextual, textural, and ancillary multisource information to improve the classification maps (cf., Section IV-A).

F. Multimodal Fusion of Optical and Radar Data

In the framework of the analysis of multisensor optical and SAR imagery, the Markovian approach to data fusion, which is based on the minimization of an energy function combining the input information sources, is exemplified by considering the joint use of VHR optical and SAR data for detecting the ground changes that occurred in between multitemporal observations [145], [146]. This case study has been selected as a challenging problem of VHR mapping to stress the benefits of exploiting the complementary properties of SAR and optical data through multisensor fusion as compared to the processing of individual data sources. A multitemporal pair of coregistered multisensor images, each composed of optical multispectral channels and of a SAR amplitude channel acquired over the same geographical area before and after ground changes, is assumed available. Change detection methods that were found successful at coarser resolution [147], [148] can be

often ineffective when applied at VHR because they do not properly capture the strong spatial information associated with VHR imagery. Moreover, change-detection methods for current VHR satellite data usually focus separately on either VHR SAR or VHR optical data and aim at detecting changes that affect the radar backscattering coefficient or the reflectance in the visible and near-infrared ranges, but not both [149], [150]. The approach to fusion based on MRF models allows both issues to be addressed by fusing the information associated with the optical data, the SAR data, and the spatial context [82], [102], [151].

As a simple illustration, let us assume that two real-valued features that emphasize the discrimination between changed and unchanged areas are extracted from the optical and SAR temporal pairs, respectively. To this end, techniques such as change vector analysis, image differencing and ratioing [146], information-theoretic distances [152], likelihood ratios [153], or copula functions [154] can be applied. Let \mathcal{I} be the pixel lattice, r_i and s_i the features computed for the i th pixel from the SAR and optical data, $\mathcal{X} = \{(r_i, s_i)\}_{i \in \mathcal{I}}$ the 2-D stochastic process of all available features, and $\mathcal{Y} = \{y_i\}_{i \in \mathcal{I}}$ the 2-D stochastic process of class labels. \mathcal{X} is continuous valued whereas \mathcal{Y} is a binary process such that $y_i = 1$ and $y_i = 0$ indicate that the i th pixel belongs to the “change” and “no-change” classes, respectively. Similar to Section V-A, \mathcal{Y} is modeled as an MRF.

In particular, according to the Markovian approach to multisource fusion, the energy function is defined as a linear combination of contributions associated with the input individual information sources [82], [102]. A simple example is described by [145]

$$U(\mathcal{Y}|\mathcal{X}) = -\alpha \sum_{i \in \mathcal{I}} \ln f(r_i|y_i) - \beta \sum_{i \in \mathcal{I}} \ln g(s_i|y_i) - \gamma \sum_{j \sim i} \delta(y_i, y_j) \quad (1)$$

where $f(\cdot)$ and $g(\cdot)$ denote suitable parametric models for the probability density functions of the SAR and optical features, respectively, when conditioned to the class label, $\delta(\cdot)$ is the Kronecker symbol (i.e., $\delta(a, b) = 1$ for $a = b$ and $\delta(a, b) = 0$ otherwise), α, β and γ are positive weight parameters, and $i \sim j$ indicates that the i th and j th pixels are neighbors ($i, j \in \mathcal{I}$; see Section V-A). The first two energy contributions are associated with the pixel-wise statistics of each data source while the third term is related to spatial context and favors the same labeling for all pixels in a homogeneous image area; α, β , and γ tune the tradeoff among the three energy terms. Among the parametric models that have been proposed to characterize the class-conditional statistics $f(\cdot)$ and $g(\cdot)$ for change detection purposes, we recall Nakagami-ratio, Weibull-ratio, log-normal [148], and generalized Gaussian distributions [147], and dictionary-based approaches [145] in the case of SAR

features, as well as Gaussian [155], [156] and generalized Gaussian distributions [145] in the case of optical features. Compared to the basic formulation in (1), much more sophisticated prior models for the spatial information can also be incorporated in the energy to characterize, for example, regions or objects through segmentation [157], spatial edges through line processes [132] or discontinuity-adaptive models [134], texture through continuous-valued random fields [134] or adaptive moving-window algorithms [102], and multiscale and multiresolution information through hierarchical trees [133]. Regardless of the case-specific definition of the energy, the minimization of $U(\mathcal{Y}|\mathcal{X})$ with respect to \mathcal{Y} determines the output binary classification result and intrinsically takes into account both optical and SAR input sources.

Fig. 14(a)–(d) shows a pair of COSMO-SkyMed SAR stripmap images (single-polarization, 5-m spatial resolution) and a pair of QuickBird optical images (four spectral channels), acquired over an area of Port-au-Prince, Haiti, before and after the 2010 earthquake. The optical images have been resampled at the same resolution of the SAR data, and all input images have been registered. Multiple types of change may be noted due, for example, to damages caused by the earthquake or seasonal variations. The aforementioned Markovian multisensor fusion approach has been applied in an unsupervised way (i.e., without involving training samples) by using change vector analysis to extract a feature from the optical pair; generalized Gaussian distributions to model the statistics of this feature [i.e., to determine $g(\cdot)$] [147]; image rationing to extract a feature from the SAR pair; a dictionary of SAR-specific parametric families to characterize the statistics of this feature [i.e., to determine $f(\cdot)$] [148]; higher order moments [147], log-cumulants, the expectation–maximization algorithm, and the mode-field approximation to estimate the parameters of these models as well as α, β , and γ [148]; and a graph cut algorithm to minimize the energy [108] (details can be found in [145]).

Fig. 14(e) shows the resulting change map. No ground truth was available, so a qualitative assessment was conducted on the basis of visual photointerpretation and comparison with ancillary information. This analysis suggested that the result in Fig. 14(e) accurately captured the main ground changes, which could be appreciated in the optical and/or SAR multitemporal data set, and allowed the meanings of the major detected “change” areas to be identified. Details on this discussion can be found in [145]. Here, we focus on comparing these results with those in Fig. 14(f) and (g) that were obtained by applying the same MRF-based binary classification to only the optical or the SAR data. On the one hand, a visual analysis points out that several major changed areas were detected through multisensor fusion while they were missed or poorly detected when only one of the two individual single-sensor sources was used. For example, both changes pointed out by the SAR pair because of their impacts on strong radar scatterers

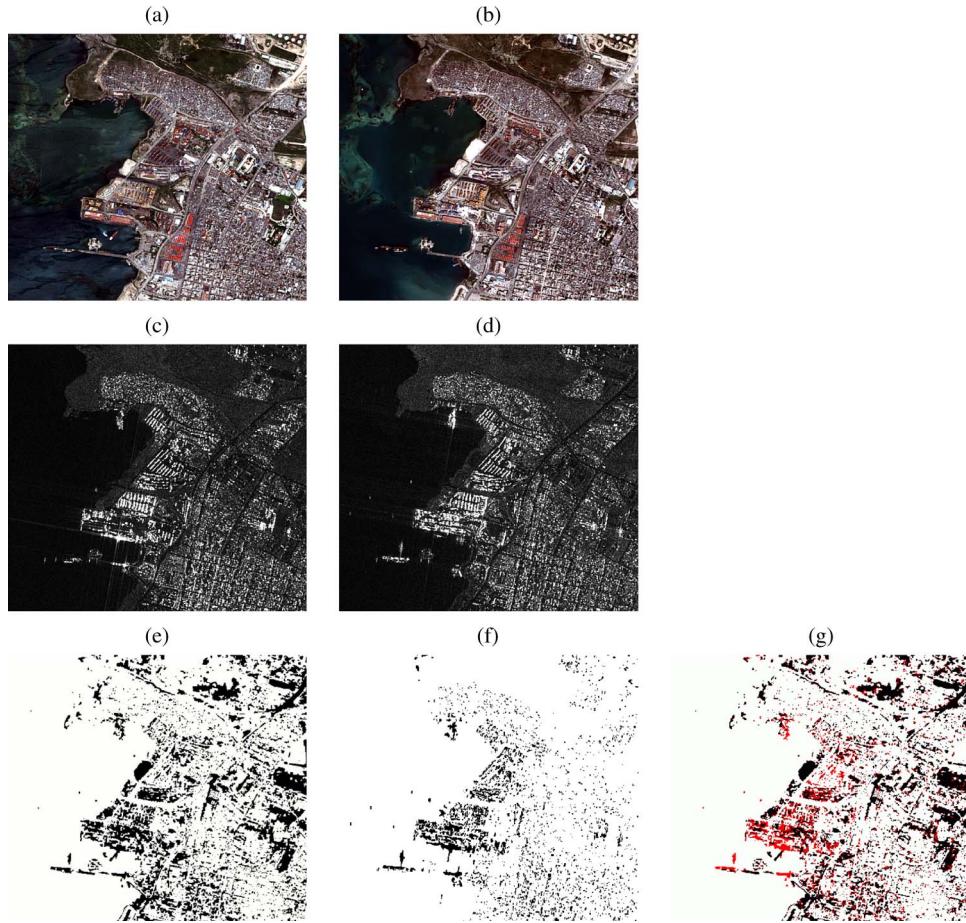


Fig. 14. Example of multisensor fusion of optical and SAR data, Port-au-Prince, Haiti: QuickBird images, September 16, 2009 (a) and, January 15, 2010 (b); COSMO-SkyMed images, April 28, 2009 (c) and January 15, 2010 (d); and change maps obtained through multisensor fusion (e) and through the separate analyses of the SAR (f) and optical (g) temporal pairs. Areas classified as changed and unchanged are marked in black and white, respectively. Differences between the results obtained through multisensor fusion and through the analysis of only optical data are shown in red in (g).

and less relevant in the optical pair, and seasonal vegetation changes, which are apparent in the optical pair but have little to no impact on the SAR pair, are correctly detected in the multisensor result. On the other hand, no visually appreciable increase in the number of false alarm was noted when using multisensor fusion rather than when dealing separately with the two individual sources.

From an application-oriented viewpoint, this scenario is consistent with the presence of types of change that are visually appreciable in either optical VHR data but not in SAR VHR data or vice versa, because of the different natures of the observed physical variables. Furthermore, there are different time lags between the SAR and optical multitemporal acquisitions. From a methodological standpoint, these results confirm the capabilities of MRF models as data fusion tools that can benefit from the information conveyed by both input EO data modalities to generate a thematic map overall detecting the information emphasized by, at least, one of the two modalities.

G. Cross-Sensor Adaptation Via Manifold Alignment

In this case study, we aim at adapting a classifier trained on a multispectral sensor to predict land-cover classes of another multispectral sensor with similar resolution but different spectral bands. We also want this process to be free from coregistration requirements. This is of particular interest for new satellite sensors such as WorldView-2 and 3, which have been built on the basis of the spectral configuration of the previous QuickBird sensor, while including additional spectral bands (see Table 2). In this case, we could think of reusing models built on QuickBird images to upgrade them to make use of the higher spectral resolution of WorldView-2 or to use classifiers built on WorldView-2 images to predict on images acquired by QuickBird (that is still flying and acquiring images): in this case, QuickBird could be used as if it was part of a WorldView-2 constellation. Such setting cannot be handled by standard classification pipelines, unless we proceed by downgrading the WorldView-2 images (e.g., by removing the four

TABLE 2 Comparison of the Spectral Bands Provided by the QuickBird and WorldView-2 Sensors

QuickBird	Band Name	WorldView-2
—	coastal	400-450 nm
450-520 nm	blue	450-510 nm
520-600 nm	green	510-580 nm
—	yellow	585-625 nm
630-690 nm	red	630-690 nm
—	red edge	705-745 nm
760-900	near IR-1	770-895 nm
—	near IR-2	860-900 nm

nonmatching bands). Recent systems based on multiview analysis [89], [129] can answer this call and be used to transfer model regardless of the dimensionality of the data space, i.e., of the sensor used in the first place.

The setting considered in this case study is summarized in Fig. 15. The images are acquired by either QuickBird or by WorldView-2. They are not coregistered and provide $d_{QB} = 4$ and $d_{WV2} = 8$ spectral bands, respectively. We use the semi-supervised manifold alignment algorithm (SS-MA [129], [158]) to align the probability density functions of the images in a latent space of dimension $d_{QB} + 2 * d_{WV2} = 20$. The basic concept of SS-MA is to align the data manifolds by projecting them to this common space, which has two desirable properties: 1) it maintains the local geometry of each manifold unchanged; and 2) it brings regions belonging to the same class closer together while it pushes those belonging to different classes apart. To enforce these properties, it uses graph Laplacians [158]. The advantage of this method over canonical correlation-based methods [89] is that SS-MA does not require the different views to be coregistered to each other, thus opening possibilities for aligning images of different locations (as in the example below). On the downside, SS-MA requires few labeled pixels in all the domains aligned (studies on how to relax this assumption can be found in [159] and [160]).

We consider three images (Fig. 16): one acquired by QuickBird over Zurich, Switzerland, in October 2006 and two acquired over Lausanne, Switzerland, in September 2010 and August 2011, respectively. We study the capability of SS-MA of exploiting one single classifier to classify the three acquisitions, even though they are neither synchronous, coregistered, or acquired by the same sensor. We use only spectral information, to study the benefits of

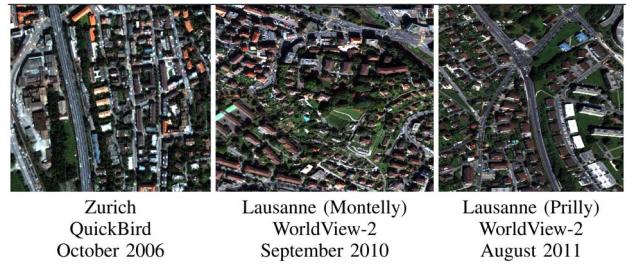


Fig. 16. Images considered in the cross-sensor classification experiment [129].

the algorithm in aligning the three original data spaces, but adding contextual information would boost the results for all cases.

We start by the assumption that one image comes with several labels (the leading training image hereafter). The other acquisitions come with a smaller number of labels. This assumption is reasonable, since one wants to transfer classifiers with minimal effort and avoid many labeled samples from the newly acquired images. We take 100 labeled pixels per class in the leading training image and then align it to the two others using an increasing number of labels from those ([10, ..., 90] labels per class). The results are illustrated in Fig. 17. We take each image in turn as the leading one (corresponding to each row of Fig. 17), align the three images, and predict all of them at once.

We compare SS-MA (blue bars) with a classification without alignment (“original” hereafter, red bars) where, in order to perform classification, we removed the four bands missing in QuickBird from the WorldView-2 acquisitions. We also report a baseline (green bars), which is a classifier trained on 100 labeled pixels per class from the test image only. For all train/test pairs, SS-MA leads to higher κ scores than the method without alignment. When testing on the same image as the leading training one (diagonal blocks of Fig. 17), adding pixels from the other images results in a degradation of the performance if no alignment is performed. But when using SS-MA, κ remains stable at a level generally higher than the baseline. Since it aligns domains in a discriminative sense, SS-MA allows to use many images without affecting the quality of the classifier in the original domain. When predicting other images than the leading training one (off-diagonal blocks of Fig. 17) the alignment procedure boosts the performance of the classifier and leads to κ scores that are between 0.1 and 0.3 κ points higher than the case where only pixels from the leading domain are used for training.

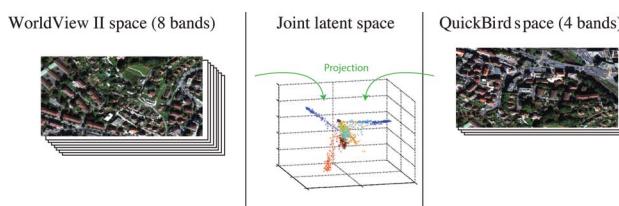


Fig. 15. Setting of the cross-sensor case study.

VI. CONCLUDING REMARKS

This paper summarized our view on the multimodal classification of Earth observation remote sensing images. We reviewed the main distinctive aspects of the typical tasks

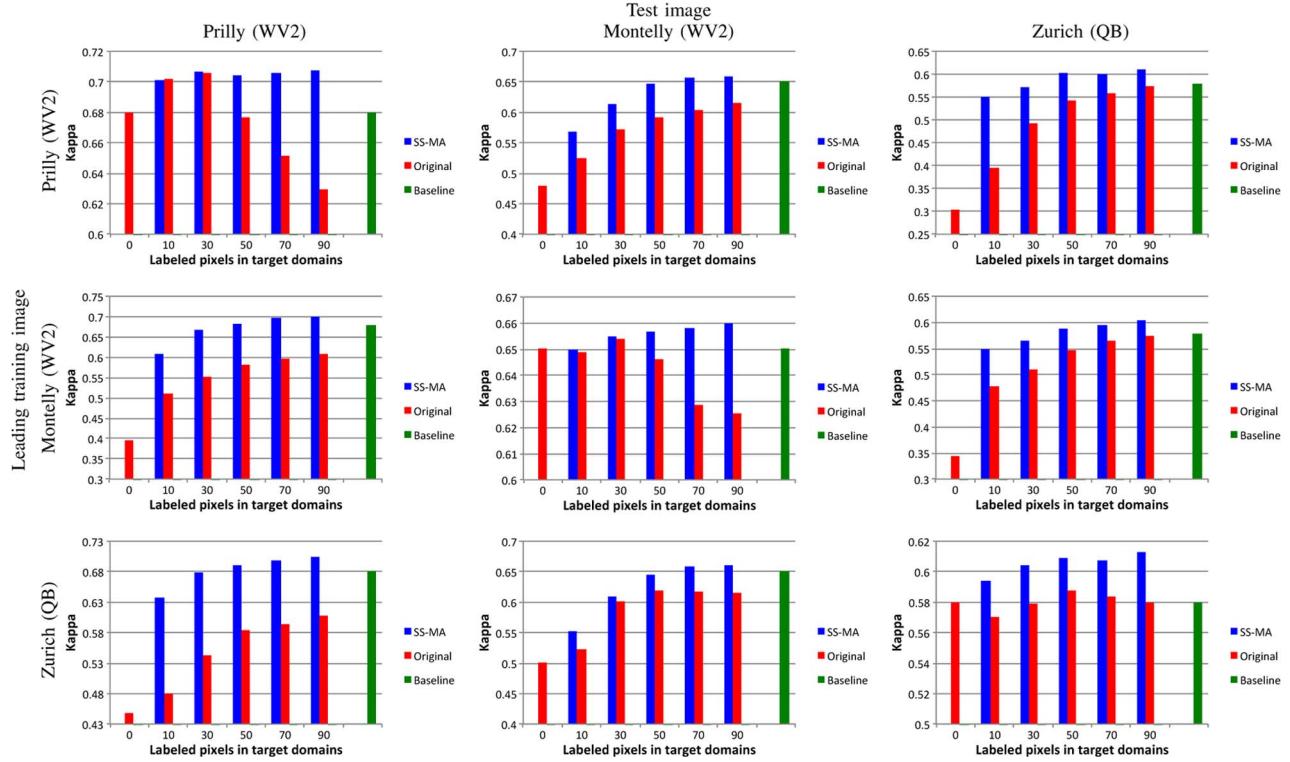


Fig. 17. Cross-sensor classification with manifold alignment. Rows indicate the image from which 100 labeled pixels per class are used. κ performances for increasing number of labeled pixels in the two other images (with [10, ..., 90] labels per class) are reported. Columns correspond to the image that has been used for testing. The baseline is the model obtained using 100 pixels per class from the test image only [129].

involved, the recent literature, and some new methodological approaches, which are bringing new perspectives to the field. In particular, we paid attention to the problems of combining multiple and heterogeneous image sources (e.g., multispectral, hyperspectral, radar, multitemporal, and multiangular images) for data classification. Multimodal image fusion is becoming the most prominent and important field in remote sensing data processing, and its importance will increase in the near future with the advent of new sensors and satellite constellations.

We have identified standard approaches that fuse data at subpixel level, pixel level, feature level, and decision level. Many techniques and approaches are available, but the choice ultimately depends on the specific application and on the desired end product. Among the emergent fields for multimodal data fusion in the remote sensing community, we highlighted multiple kernel learning that performs fusion in implicit high-dimensional feature representations; sparse dictionary learning, where fusion is conducted at the signal or information-theoretic level; structured prediction and Markov modeling, which formalize spatial and multimodal fusion through global minimum energy concepts; deep learning, probably the fastest moving area of modern computer vision; and manifold alignment where sources of different nature and dimensionality are combined at a geometrical level in a latent space. All these sound forms

of multimodal fusion allow the combination and interaction between modalities with different levels of sophistication.

We illustrated the different approaches in seven real challenging remote sensing applications:

- 1) multiresolution fusion of medium resolution multispectral images and very high spatial resolution panchromatic images;
- 2) image downscaling as a form of multitemporal image fusion and multidimensional interpolation between sensors of different spatial, spectral, and temporal resolutions;
- 3) multiangular image classification of very high spatial and angular resolutions for urban monitoring;
- 4) multisensor image fusion exploiting physically-based feature extracted from intrinsically different sensors like MERIS and AATSR for cloud screening problems;
- 5) multitemporal online image classification of land use in incomplete, inconsistent, and vague image sources;
- 6) spatial–spectral information fusion for classification of optical and radar images;
- 7) cross-sensor adaptation of classifiers via manifold alignment.

Every technique has shown desirable properties in the scenarios considered, although it is obvious that, for example,

image fusion will fail when combining images of completely different spatial resolutions, or domain adaptation techniques will fail when adapting too different scenes that have nothing in common. We feel that the adoption of these techniques in operational settings could help to monitor our planet from space in the near future. The inputs from machine learning, computer vision, and signal processing, as well as those from other applied disciplines such as bio-

medical engineering, are crucial for advancing this reborn and exciting field, whose expansion is, in our view, only beginning. ■

Acknowledgment

The authors would like to thank DigitalGlobe Inc. for the optical data on Rio and Haiti, and the Italian Space Agency for the SAR data on Haiti.

REFERENCES

- [1] S. Liang, *Quantitative Remote Sensing of Land Surfaces*. New York, NY, USA: Wiley, 2004.
- [2] T. M. Lillesand, R. W. Kiefer, and J. Chipman, *Remote Sensing and Image Interpretation*. New York, NY, USA: Wiley, 2008.
- [3] G. Shaw and D. Manolakis, "Signal processing for hyperspectral image exploitation," *IEEE Signal Process. Mag.*, vol. 50, no. 1, pp. 12–16, Jan. 2002.
- [4] J. M. Bioucas-Dias et al., "Hyperspectral remote sensing data analysis and future challenges," *IEEE Geosci. Remote Sens. Mag.*, vol. 1, no. 2, pp. 6–36, Jun. 2013.
- [5] G. Camps-Valls, D. Tuia, L. Bruzzone, and J. A. Benediktsson, "Advances in hyperspectral image classification: Earth monitoring with statistical learning methods," *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 45–54, Jan. 2014.
- [6] T. Stathaki, *Image Fusion: Algorithms and Applications*. New York, NY, USA: Academic, 2008.
- [7] C. Pohl and J. L. Van Genderen, "Multisensor image fusion in remote sensing: Concepts, methods and applications," *Int. J. Remote Sens.*, vol. 19, no. 5, pp. 823–854, Mar. 1998.
- [8] M. Drusch et al., "Sentinel-2: ESA's optical high-resolution mission for GMES operational services," *Remote Sens. Environ.*, vol. 120, pp. 25–36, 2012.
- [9] C. Donlon et al., "The Global Monitoring for Environment and Security (GMES) Sentinel-3 mission," *Remote Sens. Environ.*, vol. 120, pp. 37–57, 2012.
- [10] T. Stoffler et al., "The EnMAP hyperspectral imager—An advanced optical payload for future applications in Earth observation programmes," *Acta Astronautica*, vol. 61, no. 1–6, pp. 115–120, 2007.
- [11] D. A. Roberts, D. A. Quattrochi, G. C. Hulley, S. J. Hook, and R. O. Green, "Synergies between VSWIR and TIR data for the urban environment: An evaluation of the potential for the hyperspectral infrared imager (HyspIRI) decadal survey mission," *Remote Sens. Environ.*, vol. 117, pp. 83–101, 2012.
- [12] D. Labate et al., "The PRISMA payload optomechanical design, a high performance instrument for a new hyperspectral mission," *Acta Astronautica*, vol. 65, no. 9–10, pp. 1429–1436, 2009.
- [13] L. Wald, "Some terms of reference in data fusion," *IEEE Trans. Geosci. Remote Sens.*, vol. 37, no. 3, pp. 1190–1193, May 1999.
- [14] I. R. Farah, "A Multi Views Approach for Remote Sensing Fusion Based on Spectral, Spatial and Temporal Information Image Fusion," Rijeka, Croatia: InTech, 2011, pp. 43–70.
- [15] J. Chanussot et al., "Challenges and opportunities of multimodality and data fusion in remote sensing," *Proc. IEEE*, 2015.
- [16] P. S. Mahler Ronald, *Statistical Multisource-Multi-Target Information Fusion*. Norwood, MA, USA: Artech House, 2007.
- [17] Y. Zhang, "Understanding image fusion," *Photogramm. Eng. Remote Sens.*, vol. 70, pp. 657–661, Jun. 2004.
- [18] A. A. Goshtasby and S. Nikolov, "Image fusion: Advances in the state of the art," *Inf. Fusion*, vol. 8, no. 2, pp. 114–118, 2007.
- [19] I. R. Farah, W. Boulila, K. Saheb Ettabaa, and M. Ben Hamed, "Multiapproach system based on fusion of multispectral images for land-cover classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 12, pp. 4153–4161, Dec. 2008.
- [20] M. Alonso, B. Bookhagen, and D. A. Roberts, "Urban tree species mapping using hyperspectral and LiDAR data fusion," *Remote Sens. Environ.*, vol. 148, pp. 70–83, 2014.
- [21] R. N. Clark, *Spectroscopy and Principles of Spectroscopy, Manual of Remote Sensing*. New York, NY, USA: Wiley, 1999.
- [22] C. Oliver and S. Quegan, *Understanding Synthetic Aperture Radar Images*. Raleigh, NC, USA: SciTech Publishing, 2004.
- [23] L. Gómez-Chova et al., "Urban monitoring using multitemporal SAR and multispectral data," *Pattern Recognit. Lett.*, vol. 27, no. 4, pp. 234–243, Mar. 2006.
- [24] Z. Wang, D. Ziou, C. Armenakis, D. Li, and Q. Li, "A comparative analysis of image fusion methods," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 6, pp. 1391–1402, Jun. 2005.
- [25] P. Gamba, F. Dell'Acqua, and B. V. Dasarathy, "Urban remote sensing using multiple data sets: Past, present, future," *Inf. Fusion*, vol. 6, no. 4, pp. 319–326, 2005.
- [26] I. R. Farah and M. Ben Ahmed, "Towards an intelligent multi-sensor satellite image analysis based on blind source separation using multi-source image fusion," *Int. J. Remote Sens.*, vol. 31, no. 1, pp. 13–38, Jan. 2010.
- [27] S. Delalieux et al., "Unmixing-based fusion of hyperspatial and hyperspectral airborne imagery for early detection of vegetation stress," *IEEE J. Sel. Top. Appl. Earth Observat. Remote Sens.*, vol. 7, no. 6, pp. 2571–2582, Jun. 2014.
- [28] J. Marcello, A. Medina, and F. Eugenio, "Evaluation of spatial and spectral effectiveness of pixel-level fusion techniques," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 3, pp. 432–436, May 2013.
- [29] B. Luo, M. M. Khan, T. Bienvenu, J. Chanussot, and L. Zhang, "Decision-based fusion for pansharpening of remote sensing images," *IEEE Geosci. Remote Sens. Lett.*, vol. 10, no. 1, pp. 19–23, Jan. 2013.
- [30] W. Li, S. Prasad, and J. E. Fowler, "Decision fusion in kernel-induced spaces for hyperspectral image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 6, pp. 3399–3411, Jun. 2014.
- [31] T. G. Dietterich, "Ensemble methods in machine learning," in *Multiple Classifier Systems*, vol. 1857, Berlin, Germany: Springer-Verlag, 2000, pp. 1–15.
- [32] L. Alparone et al., "Comparison of pansharpening algorithms: Outcome of the 2006 GRS-S data-fusion contest," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 10, pp. 3012–3021, Oct. 2007.
- [33] H. R. Shahdousti and H. Ghassemian, "Fusion of MS and PAN images preserving spectral quality," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 3, pp. 611–615, Mar. 2015.
- [34] Z. Chen, H. Pu, B. Wang, and G.-M. Jiang, "Fusion of hyperspectral and multispectral images: A novel framework based on generalization of pan-sharpening methods," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 8, pp. 1418–1422, Aug. 2014.
- [35] M. M. Khan, L. Alparone, and J. Chanussot, "Pansharpening quality assessment using the modulation transfer functions of instruments," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 11, pp. 3880–3891, Nov. 2009.
- [36] B. Aiazzi, L. Alparone, S. Baronti, and A. Garzelli, "Quality assessment of Pansharpening methods and products," *IEEE Geosci. Remote Sens. Soc. Newslett.*, vol. 1, no. 161, pp. 10–18, 2011.
- [37] D. A. Landgrebe, *Signal Theory Methods in Multispectral Remote Sensing*. New York, NY, USA: Wiley, 2003.
- [38] D. L. Hall and J. Llinas, "An introduction to multisensor data fusion," *Proc. IEEE*, vol. 85, no. 1, pp. 6–23, Jan. 1997.
- [39] C. Thomas, T. Ranchin, L. Wald, and J. Chanussot, "Synthesis of multispectral images to high spatial resolution: A critical review of fusion methods based on remote sensing physics," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 5, pp. 1301–1312, May 2008.
- [40] X. He, L. Condat, J. M. Bioucas-Dias, J. Chanussot, and J. Xia, "A new pansharpening method based on spatial and spectral sparsity priors," *IEEE Trans. Image Process.*, vol. 23, no. 9, pp. 4160–4174, Sep. 2014.
- [41] K. G. Nikolakopoulos, "Comparison of nine fusion techniques for very high resolution data," *Photogramm. Eng. Remote Sens.*, vol. 74, no. 5, pp. 647–659, 2008.
- [42] R. Colditz et al., "Influence of image fusion approaches on classification accuracy: A case study," *Int. J. Remote Sens.*, vol. 27, no. 15, pp. 3311–3335, 2006.
- [43] G. Storvik, R. Fjortoft, and A. H. S. Solberg, "A Bayesian approach to classification of multiresolution remote sensing data," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 3, pp. 539–547, Mar. 2005.

- [44] A. Voisin, V. A. Krylov, G. Moser, S. B. Serpico, and J. Zerubia, "Supervised classification of multisensor and multiresolution remote sensing images with a hierarchical copula-based approach," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 6, pp. 3346–3358, Jun. 2014.
- [45] A. H. J. M. Pellemans, R. W. L. Jordans, and R. Allewijn, "Merging multispectral and panchromatic SPOT images with respect to the radiometric properties of the sensor," *Photogramm. Eng. Remote Sens.*, vol. 59, no. 1, pp. 81–87, 1993.
- [46] F. Gao, J. Masek, M. Schwaller, and F. Hall, "On the blending of the Landsat and MODIS surface reflectance: Predicting daily Landsat surface reflectance," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 8, pp. 2207–2218, Aug. 2006.
- [47] R. Zurita-Milla, J. G. P. W. Clevers, J. A. E. Van Gijzel, and M. E. Schaepman, "Using MERIS fused images for land-cover mapping and vegetation status assessment in heterogeneous landscapes," *Int. J. Remote Sens.*, vol. 32, no. 4, pp. 973–991, 2011.
- [48] J. Amorós-López et al., "Multitemporal fusion of Landsat/TM and ENVISAT/MERIS for crop monitoring," *Int. J. Appl. Earth Observat. Geoinf.*, vol. 23, no. 0, pp. 132–141, Aug. 2013.
- [49] N. Keshava and J. F. Mustard, "Spectral unmixing," *IEEE Signal Process. Mag.*, vol. 19, no. 1, pp. 44–57, Jan. 2002.
- [50] B. Huang, H. Song, H. Cui, J. Peng, and Z. Xu, "Spatial and spectral image fusion using sparse matrix factorization," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 3, pp. 1693–1704, Mar. 2014.
- [51] C. M. Gevaert and F. J. García-Haro, "A comparison of STARFM and an unmixing-based algorithm for Landsat and MODIS data fusion," *Remote Sens. Environ.*, vol. 156, pp. 34–44, 2015.
- [52] B. Zhukov, D. Oertel, F. Lanzl, and G. Reinhackel, "Unmixing-based multisensor multiresolution image fusion," *IEEE Trans. Geosci. Remote Sens.*, vol. 37, no. 3, pp. 1212–1226, May 1999.
- [53] A. Minghelli-Roman, L. Polidori, S. Mathieu-Blanc, L. Loubersac, and F. Cauneau, "Spatial resolution improvement by merging MERIS-ETM images for coastal water monitoring," *IEEE Geosci. Remote Sens. Lett.*, vol. 3, no. 2, pp. 227–231, Apr. 2006.
- [54] R. Zurita-Milla, J. Clevers, and M. E. Schaepman, "Unmixing-based Landsat TM and MERIS FR data fusion," *IEEE Geosci. Remote Sens. Lett.*, vol. 5, no. 3, pp. 453–457, Jul. 2008.
- [55] J. Amorós-López et al., "Regularized multiresolution spatial unmixing for ENVISAT/MERIS and Landsat/TM image fusion," *IEEE Geosci. Remote Sens. Lett.*, vol. 8, no. 5, pp. 844–848, Sep. 2011.
- [56] N. Longbotham et al., "Very high resolution multiangular urban classification analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 4, pp. 1155–1170, 2012.
- [57] S. Stagakis, N. Markos, O. Skyoti, and A. Kirparissis, "Monitoring canopy biophysical and biochemical parameters in ecosystem scale using satellite hyperspectral imagery: An application on a Phlomis fruticosa Mediterranean ecosystem using multisangular CHRIS/PROBA observations," *Remote Sens. Environ.*, vol. 114, pp. 977–994, 2010.
- [58] M. E. Schaepman et al., "The future of imaging spectroscopy—Prospective technologies and applications," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Denver, CO, USA, 2006, pp. 2005–2009.
- [59] E. Puttonen, J. Suomalainen, T. Hakala, and J. Peltoniemi, "Measurement of reflectance properties of asphalt surfaces and their usability as reference targets for aerial photos," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 7, pp. 2330–2339, Jul. 2009.
- [60] J. Verrelst, M. E. Schaepmann, B. Koetz, and M. Kneubühler, "Angular sensitivity analysis of vegetation indices derived from CHRIS/PROBA data," *Remote Sens. Environ.*, vol. 112, pp. 2341–2353, 2008.
- [61] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Trans. Knowl. Data Eng.*, vol. 22, no. 10, pp. 1345–1359, Oct. 2010.
- [62] F. Pacifici, N. Longbotham, and W. Emery, "The importance of physical quantities for the analysis of multitemporal and multiangular optical very high spatial resolution images," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 10, pp. 6241–6256, Oct. 2014.
- [63] L. Bruzzone and S. Serpico, "An iterative technique for the detection of land-cover transitions in multitemporal remote-sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 35, no. 4, pp. 858–867, Jul. 1997.
- [64] F. Wang, "A knowledge-based vision system for detecting land changes at urban fringes," *IEEE Trans. Geosci. Remote Sens.*, vol. 31, no. 1, pp. 136–145, Mar. 1993.
- [65] J. T. Morisette and S. Khorram, "An introduction to using generalized linear models to enhance satellite-based change detection," in *Proc. Int. Conf. Geosci. Remote Sens.*, 1997, pp. 1282–1284.
- [66] D. L. Civco, "Artificial neural networks for land cover classification and mapping," *Int. J. Geograph. Inf. Syst.*, vol. 7, no. 2, pp. 173–186, 1993.
- [67] D. Kushardono, K. Fukue, H. Shimoda, and T. Sakata, "Comparison of multi-temporal image classification methods," in *Proc. Int. Conf. Geosci. Remote Sens.*, 1995, pp. 1282–1284.
- [68] S. Gopal and C. Woodcock, "Remote sensing of forest change using artificial neural networks," *IEEE Trans. Geosci. Remote Sens.*, vol. 34, no. 2, pp. 189–202, Mar. 1996.
- [69] J. Li and R. M. Narayanan, "A shape-based approach to change detection of lakes using time series remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 41, no. 11, pp. 2466–2477, Nov. 2003.
- [70] D. Liu, M. Kelly, and P. Gong, "Classifying multi-temporal Landsat TM imagery using Markov random fields and support vector machines," in *Proc. 3rd Int. Workshop Anal. Multi-temporal Remote Sens. Images*, 2005, pp. 225–228.
- [71] T. M. Pellizzetti, P. Gamba, P. Lombardo, and F. Dell'Acqua, "Multitemporal/multiband SAR classification of urban areas using spatial analysis: Statistical versus neural kernel-based approach," *IEEE Trans. Geosci. Remote Sens.*, vol. 41, no. 10, pp. 2338–2353, Oct. 2003.
- [72] P. Gamba, F. Dell'Acqua, and G. Lisini, "Change detection of multitemporal SAR data in urban areas combining feature-based and pixel-based techniques," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 11, pp. 2820–2827, Oct. 2006.
- [73] F. Melgani and S. B. Serpico, "A Markov random field approach to spatiotemporal contextual image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 41, no. 11, pp. 2478–2487, Nov. 2003.
- [74] F. Melgani, "Classification of multitemporal remote-sensing images by a fuzzy fusion of spectral and spatio-temporal contextual information," *Int. J. Pattern Recognit. Artif. Intell.*, vol. 18, no. 2, pp. 143–156, Mar. 2004.
- [75] A. Boucher, K. C. Seto, and A. G. Journeel, "A novel method for mapping land cover changes: Incorporating time and space with geostatistics," *IEEE Trans. Geosci. Remote Sens.*, vol. 44, no. 11, pp. 3427–3435, Nov. 2006.
- [76] G. Camps-Valls, L. Gómez-Chova, J. Muñoz Marí, J. L. Rojo-Álvarez, and M. Martínez-Ramón, "Kernel-based framework for multitemporal and multisource remote sensing data classification and change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 6, pp. 1822–1835, Jun. 2008.
- [77] M. Datcu, F. Melgani, A. Piardi, and S. B. Serpico, "Multisource data classification with dependence trees," *IEEE Trans. Geosci. Remote Sensing*, vol. 40, no. 3, pp. 609–617, Mar. 2002.
- [78] J. A. Benediktsson, P. H. Swain, and O. K. Ersoy, "Neural network approaches versus statistical methods in classification of multisource remote sensing data," *IEEE Trans. Geosci. Remote Sens.*, vol. 28, no. 4, pp. 540–552, Jul. 1990.
- [79] B. Waske and J. A. Benediktsson, "Fusion of support vector machines for classification of multisensor data," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 12, pp. 781–790, Dec. 2007.
- [80] S. L. Hegaré-Mascle, I. Bloch, and D. Vidal-Madjar, "Application of Dempster-Shafer evidence theory to unsupervised classification in multisource remote sensing," *IEEE Trans. Geosci. Remote Sens.*, vol. 35, no. 4, pp. 1018–1021, Jul. 1997.
- [81] G. Camps-Valls, L. Gómez-Chova, J. Muñoz Marí, J. Vila-Francés, and J. Calpe-Maravilla, "Composite kernels for hyperspectral image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 3, no. 1, pp. 93–97, Jan. 2006.
- [82] A. H. S. Solberg, T. Taxt, and A. K. Jain, "A Markov random field model for classification of multisource satellite imagery," *IEEE Trans. Geosci. Remote Sens.*, vol. 34, no. 1, pp. 100–113, Jan. 1996.
- [83] C. Debes et al., "Hyperspectral and LiDAR data fusion: Outcome of the 2013 GRSS data fusion contest," *IEEE J. Sel. Top. Appl. Earth Observat. Remote Sens.*, vol. 7, no. 6, pp. 2405–2418, Jun. 2014.
- [84] G. Camps-Valls and L. Bruzzone, *Kernel Methods for Remote Sensing Data Analysis*. Hoboken, NJ, USA: Wiley, 2009.
- [85] A. Rakotomamonjy, F. R. Bach, S. Canu, and Y. Grandvalet, "SimpleMKL," *J. Mach. Learn. Res.*, vol. 9, pp. 2491–2521, 2008.
- [86] D. Tuia, F. Ratle, A. Pozdnoukhov, and G. Camps-Valls, "Multi-source composite kernels for urban image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 1, pp. 88–92, Jan. 2010.
- [87] D. Tuia, G. Camps-Valls, G. Matasci, and M. Kanevski, "Learning relevant image features with multiple kernel classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 10, pp. 3780–3791, Oct. 2010.
- [88] L. Gómez-Chova, G. Camps-Valls, L. Bruzzone, and J. Calpe-Maravilla, "Mean map kernel methods for semisupervised cloud classification,"

- [89] M. Volpi, G. Camps-Valls, and D. Tuia, "Spectral alignment of cross-sensor images with automated kernel canonical correlation analysis," *ISPRS Photogramm. Remote Sens.*, doi:10.1016/j.isprsjprs.2015.02.005.
- [90] Z. Li and H. Leung, "Fusion of multispectral and panchromatic images using a restoration-based method," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 5, pp. 1482–1491, May 2009.
- [91] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE Trans. Image Process.*, vol. 19, no. 11, pp. 2861–2873, Nov. 2010.
- [92] S. Li and B. Yang, "A new pan-sharpening method using a compressed sensing technique," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 2, pp. 738–746, Feb. 2011.
- [93] S. Li, H. Yin, and L. Fang, "Remote sensing image fusion via sparse representations over learned dictionaries," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 9, pp. 4779–4789, Sep. 2013.
- [94] M. Cheng, C. Wang, and J. Li, "Sparse representation based pansharpening using trained dictionary," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 1, pp. 293–297, Jan. 2014.
- [95] M. Ghahremani and H. Ghassemian, "Remote sensing image fusion using ripplet transform and compressed sensing," *IEEE Geosci. Remote Sens. Lett.*, vol. 12, no. 3, pp. 502–506, Mar. 2015.
- [96] X. Ding, Y. Jiang, Y. Huang, and J. Paisley, "Pan-sharpening with a Bayesian nonparametric dictionary learning model," in *Proc. 17th Int. Conf. Artif. Intell. Stat.*, Reykjavik, Iceland, 2014, pp. 176–184.
- [97] B. Huang and H. Song, "Spatiotemporal reflectance fusion via sparse representation," *IEEE Trans. Geosci. Remote Sens.*, vol. 50, no. 10, pp. 3707–3716, Oct. 2012.
- [98] H. Song and B. Huang, "Spatiotemporal satellite image fusion through one-pair image learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 51, no. 4, pp. 1883–1896, Apr. 2013.
- [99] B. C. K. Tso and P. M. Mather, "Classification of multisource remote sensing imagery using a genetic algorithm and Markov random fields," *IEEE Trans. Geosci. Remote Sens.*, vol. 37, no. 3, pp. 1255–1260, May 1999.
- [100] M. Xu, H. Chen, and P. K. Varshney, "An image fusion approach based on Markov random fields," *IEEE Trans. Geosci. Remote Sens.*, vol. 49, no. 12, pp. 5116–5127, Dec. 2011.
- [101] F. Li, X. Jia, D. Fraser, and A. Lambert, "Super resolution for remote sensing images based on a universal hidden Markov tree model," *IEEE Trans. Geosci. Remote Sens.*, vol. 48, no. 3, pp. 1270–1278, Mar. 2010.
- [102] G. Moser, S. B. Serpico, and J. A. Benediktsson, "Land-cover mapping by Markov modeling of spatial-contextual information," *Proc. IEEE*, vol. 101, no. 3, pp. 631–651, Mar. 2013.
- [103] P. C. Smits and S. G. Dellepiane, "Synthetic aperture radar image segmentation by a detail preserving Markov random field approach," *IEEE Trans. Geosci. Remote Sens.*, vol. 35, no. 4, pp. 844–857, Jul. 1997.
- [104] C. Benedek and T. Sziranyi, "Change detection in optical aerial images by a multilayer conditional mixed Markov model," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 10, pp. 3416–3430, Oct. 2009.
- [105] M. Vakalopoulou, K. Karantzalos, N. Komodakis, and N. Paragios, "Simultaneous registration and change detection in multitemporal, very high resolution remote sensing data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2015, pp. 61–69.
- [106] C. Benedek, M. Shadayeh, Z. Kato, T. Sziranyi, and J. Zerubia, "Multilayer Markov random field models for change detection in optical remote sensing images," *ISPRS Photogramm. Remote Sens.*, doi:10.1016/j.isprsjprs.2015.02.006.
- [107] T. Sziranyi and M. Shadayeh, "Segmentation of remote sensing images using similarity-measure-based fusion-MRF model," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 9, pp. 1544–1548, Sep. 2014.
- [108] V. Kolmogorov and R. Zabih, "What energy function can be minimized via graph cuts?" *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 2, pp. 147–159, Feb. 2004.
- [109] J. D. Wegner, R. Hänsch, A. Thiele, and U. Soergel, "Building detection from one orthophoto and high-resolution InSAR data using conditional random fields," *IEEE J. Sel. Top. Appl. Earth Observat. Remote Sens.*, vol. 4, no. 1, pp. 83–91, Mar. 2011.
- [110] C. Haeme, C. Zach, A. Cohen, R. Angst, and M. Pollefeys, "Joint 3D scene reconstruction and class segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013, pp. 97–104.
- [111] S. Kadoury, H. Labelle, and N. Paragios, "Automatic inference of articulated spine models in CT images using high-order Markov random fields," *Med. Image Anal.*, vol. 15, no. 4, pp. 426–437, 2011.
- [112] L. Xu et al., "Oil spill candidate detection from SAR imagery using a thresholding-guided stochastic fully-connected conditional random field model," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2015, pp. 79–86.
- [113] M. Volpi and V. Ferrari, "Semantic segmentation of urban scenes by learning local class interactions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2015, pp. 1–9.
- [114] S. Paisitkriangkrai, J. Sherrah, P. Janney, and A. Van-Den Hengel, "Effective semantic pixel labelling with convolutional networks and conditional random fields," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Jun. 2015, pp. 36–43.
- [115] Y. LeCun et al., "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, no. 4, pp. 541–551, 1989.
- [116] G. Hinton et al., "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Process. Mag.*, vol. 29, no. 6, pp. 82–97, Nov. 2012.
- [117] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proc. Int. Conf. Mach. Learn.*, Helsinki, Finland, 2008, pp. 160–167.
- [118] J. Ngiam et al., "Multimodal deep learning," in *Proc. Int. Conf. Mach. Learn.*, Bellevue, WA, USA, 2011, pp. 689–696.
- [119] N. Srivastava and R. Salakhutdinov, "Multimodal learning with deep Boltzmann machines," in *Proc. NIPS*, 2012, pp. 2231–2239.
- [120] C. Vaduva, I. Gavat, and M. Datcu, "Deep learning in very high resolution remote sensing image information mining communication concept," in *Proc. 20th Eur. Signal Process. Conf.*, Aug. 2012, pp. 2506–2510.
- [121] V. Mnih and G. Hinton, "Learning to Label Aerial Images From Noisy Data," in *Proc. International Conference on Machine Learning, ICML*, pp. 567–574, Jun. 2012.
- [122] Y. Chen, Z. Lin, X. Zhao, G. Wang, and Y. Gu, "Deep learning-based classification of hyperspectral data," *IEEE J. Sel. Top. Appl. Earth Observat. Remote Sens.*, vol. 7, no. 6, pp. 2094–2107, Jun. 2014.
- [123] X. Chen, S. Xiang, C.-L. Liu, and C.-H. Pan, "Vehicle detection in satellite images by hybrid deep convolutional neural networks," *IEEE Geosci. Remote Sens. Lett.*, vol. 11, no. 10, pp. 1797–1801, Oct. 2014.
- [124] A. Romero, C. Gatta, and G. Camps-Valls, "Unsupervised deep feature extraction of hyperspectral images," in *Proc. IEEE Workshop Hyperspectral Image Signal Process. Whispers*, 24–27 June, Lausanne, Switzerland, 2014.
- [125] A. Romero, C. Gatta, and G. Camps-Valls, "Unsupervised deep feature extraction for remote sensing image classification," *IEEE Trans. Geosci. Remote Sens.*, 2015.
- [126] D. Tuia, N. Courty, and R. Flamary, "Multiclass feature learning for hyperspectral image classification: Sparse and hierarchical solutions," *ISPRS J. Int. Soc. Photogramm. Remote Sens.*, vol. 105, pp. 272–285, 2015.
- [127] O. Penatti, K. Nogueira, and J. A. dos Santos, "Do deep features generalize from everyday objects to remote sensing and aerial scenes domains?" in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, pp. 44–51, Jun. 2015.
- [128] C. Wang, P. Kraft, and S. Mahadevan, "Manifold alignment," in *Manifold Learning: Theory and Applications*, Y. Ma and Y. Fu, Eds. Boca Raton, FL, USA: CRC Press, 2011.
- [129] D. Tuia, M. Volpi, M. Trolliet, and G. Camps-Valls, "Semisupervised manifold alignment of multimodal remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 52, no. 12, pp. 7708–7720, Dec. 2014.
- [130] G. Matasci, M. Volpi, M. Kanevski, L. Bruzzone, and D. Tuia, "Semisupervised transfer component analysis for domain adaptation in remote sensing image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 7, pp. 3550–3564, Jul. 2015.
- [131] G. Moser and S. B. Serpico, "Joint classification of panchromatic and multispectral images by multiresolution fusion through Markov random fields and graph cuts," in *Proc. 17th Int. Conf. Digital Signal Process.*, Corfu, Greece, 2011, doi:10.1109/ICDSP.2011.6005014.
- [132] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, the Bayesian restoration of images," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-6, no. 6, pp. 721–741, Dec. 1984.
- [133] Z. Kato and J. Zerubia, "Markov random fields in image segmentation," *Found. Trends Signal Process.*, vol. 5, no. 1–2, pp. 1–155, 2011.
- [134] S. Z. Li, *Markov Random Field Modeling in Image Analysis*. New York, NY, USA: Springer-Verlag, 2010.
- [135] L. R. Rabiner, "A tutorial on hidden Markov models and selected applications in speech recognition," *Proc. IEEE*, vol. 77, no. 2, pp. 257–286, Feb. 1989.

- [136] G. Matasci, N. Longbotham, F. Pacifici, M. Kanevski, and D. Tuia, "Understanding angular effects and their consequences on urban land-cover model portability in VHR in-track multiangle image sequences," *ISPRS Photogramm. Remote Sens.*, submitted for publication.
- [137] G. Matasci, L. Bruzzone, M. Volpi, D. Tuia, and M. Kanevski, "Investigating the feature extraction framework for domain adaptation in remote sensing image classification," in *Int. Conf. Pattern Recognit. Appl. Methods*, Barcelona, Spain, pp. 419–424, 2013.
- [138] L. Gómez-Chova, J. Muñoz-Marí, J. Amorós-López, E. Izquierdo-Verdiguier, and G. Camps-Valls, "Advances in synergy of AATSR-MERIS sensors for cloud detection," in *Proc. IEEE Int. Geosci. Remote Sens. Symp.*, Jul. 2013, pp. 4391–4394.
- [139] M. Rast, J. L. Bézy, and S. Bruzzi, "The ESA medium resolution imaging spectrometer MERIS: A review of the instrument and its mission," *Int. J. Remote Sens.*, vol. 20, no. 9, pp. 1681–1702, Jun. 1999.
- [140] D. Llewellyn-Jones et al., "AATSR: Global-change and surface-temperature measurements from Envisat," *ESA Bull.*, vol. 105, pp. 11–21, Feb. 2001.
- [141] L. Gómez-Chova, G. Camps-Valls, J. Calpe, L. Guanter, and J. Moreno, "Cloud-screening algorithm for ENVISAT/MERIS multispectral images," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 12, pt. 2, pp. 4105–4118, Dec. 2007.
- [142] M. C. Reed and B. Simon, *Functional Analysis*, vol. 1. New York, NY, USA: Academic, 1980, ser. Methods of Modern Mathematical Physics.
- [143] M. J. Barnsley, J. J. Settle, M. Cutter, D. Lobb, and F. Teston, "The PROBA/CHRIS mission: A low-cost smallsat for hyperspectral, multi-angle, observations of the Earth surface and atmosphere," *IEEE Trans. Geosci. Remote Sens.*, vol. 42, no. 7, pp. 1512–1520, Jul. 2004.
- [144] I. Hajnsek, R. Bianchi, M. Davidson, and M. Wooging, "AgriSAR 2006—Airborne SAR and Optics campaigns for an improved monitoring of agricultural processes and practices," in *Proc. 4th Int. Workshop Anal. Multitemporal Remote Sens. Images*, AgriSAR 2006 Team, Leuven, Belgium, pp. 1–8, Jul. 2007.
- [145] L. Cianci, G. Moser, and S. B. Serpico, "Change detection from very high-resolution multisensor remote-sensing images by a Markovian approach," in *Proc. IEEE GOLD Remote Sens. Conf.*, Rome, Italy, 2012, pp. 46–48.
- [146] A. Singh, "Digital change detection techniques using remotely-sensed data," *Int. J. Remote Sens.*, vol. 10, pp. 989–1003, 1989.
- [147] Y. Bazi, L. Bruzzone, and F. Melgani, "An unsupervised approach based on the generalized Gaussian model to automatic change detection in multitemporal SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 4, pp. 874–887, Apr. 2005.
- [148] G. Moser and S. B. Serpico, "Unsupervised change detection from multichannel SAR data by Markovian data fusion," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 7, pp. 2114–2128, Jul. 2009.
- [149] F. Bovolo and L. Bruzzone, "A detail-preserving scale-driven approach to change detection in multitemporal SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 12, pp. 2963–2972, Dec. 2005.
- [150] G. Moser and S. B. Serpico, "Unsupervised change detection with high-resolution SAR images by edge-preserving Markov random fields and graph-cuts," in *Proc. Int. Geosci. Remote Sens. Symp.*, Munich, Germany, 2012, pp. 1984–1987.
- [151] A. Bendjebbour, Y. Delignon, L. Fouque, V. Samson, and W. Pieczynski, "Multisensor image segmentation using Dempster-Shafer fusion in Markov fields context," *IEEE Trans. Geosci. Remote Sens.*, vol. 39, no. 8, pp. 1789–1798, Aug. 2001.
- [152] J. Ingla and G. Mercier, "A new statistical similarity measure for change detection in multitemporal SAR images and its extension to multiscale change analysis," *IEEE Trans. Geosci. Remote Sens.*, vol. 45, no. 5, pp. 1432–1445, May 2007.
- [153] K. Conradsen, A. Aasbjerg Nielsen, J. Schou, and H. Skriver, "A test statistic in the complex Wishart distribution and its application to change detection in polarimetric SAR data," *IEEE Trans. Geosci. Remote Sens.*, vol. 41, no. 1, pp. 4–19, Jan. 2003.
- [154] G. Mercier, G. Moser, and S. B. Serpico, "Conditional copulas for change detection in heterogeneous remote sensing images," *IEEE Trans. Geosci. Remote Sens.*, vol. 46, no. 5, pp. 1428–1441, May 2008.
- [155] L. Bruzzone and D. F. Prieto, "Automatic analysis of the difference image for unsupervised change detection," *IEEE Trans. Geosci. Remote Sens.*, vol. 38, no. 3, pp. 1171–1182, May 2000.
- [156] F. Melgani, G. Moser, and S. B. Serpico, "Unsupervised change detection methods for remote sensing images," *Opt. Eng.*, vol. 41, no. 12, pp. 3288–3297, 2002.
- [157] T. Blaschke, S. Lang, and G. Hay, *Object-Based Image Analysis*. New York, NY, USA: Springer-Verlag, 2008.
- [158] C. Wang and S. Mahadevan, "Heterogeneous domain adaptation using manifold alignment," in *Proc. Int. Joint Conf. Artif. Intell.*, 2011, pp. 1541–1546.
- [159] C. Wang and S. Mahadevan, "Manifold alignment without correspondence," in *Proc. Int. Joint Conf. Artif. Intell.*, Pasadena, CA, USA, 2009, pp. 1273–1278.
- [160] D. Tuia, M. Volpi, and G. Camps-Valls, "Unsupervised alignment of image manifolds with centrality measures," in *Proc. Int. Conf. Pattern Recognit.*, Stockholm, Sweden, 2014, pp. 912–917.

ABOUT THE AUTHORS

Luis Gómez-Chova (Senior Member, IEEE) received the B.Sc. (with first class honors), M.Sc., and Ph.D. degrees in electronics engineering from the University of Valencia, Valencia, Spain, in 2000, 2002, and 2008, respectively.

Since 2000, he has been with the Department of Electronics Engineering, University of Valencia, first enjoying a research scholarship from the Spanish Ministry of Education and currently as an Associate Professor. He is also a Researcher at the Image Processing Laboratory (IPL), where his work is mainly related to pattern recognition and machine learning applied to remote sensing multispectral images and cloud screening. He conducts and supervises research on these topics within the frameworks of several national and international projects. He is the author (or coauthor) of 30 international journal papers, more than 90 international conference papers, and several international book chapters.

Dr. Gómez-Chova was awarded the National Award for Electronics Engineering by the Spanish Ministry of Education.



Devis Tuia (Senior Member, IEEE) was born in Mendrisio, Switzerland, in 1980. He received the Diploma degree in geography from the University of Lausanne (UNIL), Lausanne, Switzerland, in 2004, the Master of Advanced Studies in environmental engineering from the Federal Institute of Technology of Lausanne (EPFL), Lausanne, Switzerland, in 2005, and the Ph.D. degree in environmental sciences from UNIL in 2009.

He was a Visiting Postdoctoral Researcher at the University of Valéncia, Valéncia, Spain, and the University of Colorado, Boulder, CO, USA. He then worked as a Senior Research Associate with EPFL under a Swiss National Foundation program. Since 2014, he has been an Assistant Professor with the Department of Geography, University of Zurich, Zurich, Switzerland. His research interests include the development of algorithms for information extraction and data fusion of remote sensing images using machine learning algorithms with focus on multimodal remote sensing involving multisensor, large data, and multimedia integration issues.



Dr. Tuia is an Associate Editor of the IEEE JOURNAL OF SELECTED TOPICS IN APPLIED EARTH OBSERVATION AND REMOTE SENSING. In 2014–2015, he was a Guest Editor for the Special Issue about Data Fusion of the IEEE GEOSCIENCE AND REMOTE SENSING MAGAZINE. Since 2013, has been serving as a Co-Chair of the Image Analysis and Data Fusion Technical Committee of the IEEE Geoscience and Remote Sensing Society.

Gabriele Moser (Senior Member, IEEE) received the Laurea degree (M.Sc. equivalent; *summa cum laude*) in telecommunications engineering and the Ph.D. degree in space sciences and engineering from the University of Genoa, Genoa, Italy, in 2001 and 2005, respectively.

Since 2014, he has been an Associate Professor of Telecommunications at the University of Genoa. Since 2001, he has cooperated with the Image Processing and Pattern Recognition for Remote Sensing Laboratory, University of Genoa. Since 2013, he has been the Head of the Remote Sensing for Environment and Sustainability Laboratory, Savona Campus of the University of Genoa. From January to March 2004, he was a Visiting Student at the Institut National de Recherche en Informatique et en Automatique (INRIA), Sophia Antipolis, France. His research activity is focused on the development of image processing and pattern recognition methodologies for remote sensing data interpretation. His current research interests include multitemporal, multisensor, and multiresolution data fusion; contextual classification; kernel-based methods; and geo/biophysical parameter estimation.

Dr. Moser has been a reviewer for several international journals. He has been an Associate Editor of the IEEE GEOSCIENCE AND REMOTE SENSING LETTERS and *Pattern Recognition Letters* since 2008 and 2011, respectively. He has been a Guest Editor of a Special Issue on Data Fusion in Remote Sensing of the IEEE GEOSCIENCE AND REMOTE SENSING MAGAZINE (September 2015). Since 2013, he has been the Chairman of the Image Analysis and Data Fusion Technical Committee of the IEEE Geoscience and Remote Sensing Society. He was the recipient of the Best Paper Award at the 2010 IEEE Workshop on Hyperspectral Image and Signal Processing.



Gustau Camps-Valls (Senior Member, IEEE) received the B.Sc. degree in physics, the B.Sc. degree in electronics engineering, and the Ph.D. degree in physics from the Universitat de València, València, Spain, in 1996, 1998, and 2002, respectively.

He is currently an Associate Professor (Hab. Full Professor) in the Department of Electronics Engineering, Universitat de València. His research is conducted in the Image and Signal Processing (ISP) group. He was a Visiting Researcher at the Remote Sensing Laboratory, University of Trento, Trento, Italy, in 2002 and the Max Planck Institute for Biological Cybernetics, Tübingen, Germany, in 2009; and an Invited Professor at the Laboratory of Geographic Information Systems, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland, in 2013. He is interested in the development of machine learning algorithms for geoscience and remote sensing data analysis. He is an author of 120 journal papers, more than 150 conference papers, 20 international book chapters, and an editor of the books *Kernel Methods in Bioengineering, Signal and Image Processing* (Hershey, PA, USA: IGI, 2007), *Kernel Methods for Remote Sensing Data Analysis* (New York, NY, USA: Wiley, 2009), and *Remote Sensing Image Processing* (California, Morgan & Claypool, MC, 2011). He's a coeditor of the forthcoming book *Digital Signal Processing with Kernel Methods* (New York, NY, USA: Wiley, 2015).



Dr. Camps-Valls holds a Hirsch's h index $h = 40$, entered the ISI list of Highly Cited Researchers in 2011, and Thomson Reuters ScienceWatch identified one of his papers on kernel-based analysis of hyperspectral images as a Fast Moving Front research. In 2015, he received an ERC consolidator grant on statistical learning for Earth observation data analysis. He is a referee and Program Committee member of many international journals and conferences. Since 2007, he has been a member of the Data Fusion Technical Committee of the IEEE Geoscience and Remote Sensing Society, and since 2009, he has been a member of the Machine Learning for Signal Processing Technical Committee of the IEEE Signal Processing Society. He is member of the MTG-IRS Science Team (MIST) of EUMETSAT. He is an Associate Editor of the IEEE TRANSACTIONS ON SIGNAL PROCESSING, IEEE SIGNAL PROCESSING LETTERS, and IEEE GEOSCIENCE AND REMOTE SENSING LETTERS.