

Sentinel2 ToolBox Level2 Products  
Date Issued: 02.05.2016  
Issue: V1.1

---



# **S2ToolBox Level 2 products: LAI, FAPAR, FCOVER**

**Version 1.1**

Book Captain: Marie Weiss (INRA)

Contributing Author: Fred Baret (INRA)

## Change Record

<i>Issue/Rev</i>	<i>Date</i>	<i>Page(s)</i>	<i>Description of Change</i>
1.0	30.04.2016	All	Release of Draft 1
1.1	03.05.2016	§ 2.1.4 § 2.1.5 § 3.3.4.1 § 3.4.4 Table 4 Table 9	Specifications and theoretical performances of neural networks trained for Canopy Chlorophyll Content and Canopy Water Content retrieval were added.
1.1	03.05.2016		Include a sentence in the conclusion concerning the possibility to use the neural networks at lower spatial resolution than 20m.

## TABLE OF CONTENTS

<b>1</b>	<b>Background of the document</b> .....	<b>8</b>
1.1	Executive Summary .....	8
1.2	Scope and Objectives .....	8
1.3	Content of the document.....	9
1.4	Symbols and Acronyms .....	9
<b>2</b>	<b>Algorithm Justification and Overview</b> .....	<b>11</b>
2.1	The considered Products and their definitions.....	11
2.1.1	Leaf Area Index (LAI).....	11
2.1.2	FAPAR .....	11
2.1.3	Cover fraction (FVC).....	11
2.1.4	Canopy chlorophyll content (CCC).....	12
2.1.5	Canopy water content (CWC) .....	12
2.2	Instrument characteristics .....	12
2.3	Justification for the algorithm selection and design .....	13
2.4	Algorithm outline .....	15
2.4.1	Training the neural network.....	16
2.4.2	Operational use of the neural network .....	16
<b>3</b>	<b>Algorithm Description</b> .....	<b>18</b>
3.1	Inputs and outputs .....	18
3.1.1	Inputs.....	18
3.1.2	Outputs .....	20
3.2	Reflectance models .....	20
3.2.1	Leaf optical properties.....	21
3.2.2	Canopy radiative transfer models .....	22
3.2.3	Background reflectance model.....	22
3.3	Generation of the training database.....	25
3.3.1	Radiative transfer model .....	25
3.3.2	Distribution of the vegetation input variables .....	25
3.3.3	Actual generation of the inputs for the training data base.....	32
3.3.4	Consequence on the distribution of the output variables .....	38
3.4	Training the neural network.....	39
3.4.1	Normalization of the input and output values .....	39
3.4.2	Network architecture .....	39



3.4.3	Learning process .....	41
3.4.4	Theoretical Performances .....	41
3.5	Definition Domain .....	44
3.5.1	Input out of range flag .....	44
3.5.2	Output out of range flag .....	46
<b>4</b>	<b>Practical considerations for the algorithm implementation .....</b>	<b>47</b>
4.1	Requirements for easily updating the algorithm .....	47
4.2	Algorithm Implementation .....	47
<b>5</b>	<b>Conclusion.....</b>	<b>49</b>
<b>6</b>	<b>References.....</b>	<b>50</b>

## List of Figures

Figure 1. SENTINEL2A MSI spectral response (Top: Visible – near Infrared, Bottom: ShortWave Infrared) from ESA SENTINEL online web site. ....	13
Figure 2. Flow chart showing how the products ( $\hat{V}$ ) are generated operationally. ANN corresponds to an Artificial Neural Network characterized par its structure and its coefficients (synaptic weights and bias); $R_{TOC}$ corresponds to the SENTINEL2 Top Of Canopy reflectance used in the operational mode and $V$ correspond to the biophysical variable in the training data base and estimated by running the ANN over the simulated SENTINEL2 TOC reflectance and geometry.....	16
Figure 3. Comparison between the performances of a global neural network versus interpolated ones. Top: theoretical performances against a simulated test database, Bottom: scatterplot between interpolated and global neural network LAI estimates in the test database. Two interpolated neural networks are considered: 1°step or 5° step.....	20
Figure 4: The coupled PROSPECT+SAIL model to generate the training database made of TOC reflectances and corresponding biophysical variables.....	25
Figure 5: Distribution of leaf Chlorophyll content as described by Feret et al (2008). Values of Chlorophyll include chlorophylls a and b as well as carotenoids. Distribution in blue corresponds to data mainly from tropical forests (Hawaii data set in Feret’s paper). Distribution in red corresponds to other species. ....	26
Figure 6: On the left, distribution of dry matter contents (Cdm) derived from literature review: (Prado and De Moraes 1997; Poorter and Evans 1998; Reich, Ellsworth et al. 1998; Poorter and De Jong 1999). On the right, distributions reported by Feret et al (2008): distribution in blue corresponds to data mainly from tropical forests (Hawaii data set in Feret’s paper). Distribution in red corresponds to other species.....	27
Figure 7. RMSE values associated to the reconstruction of all the soils investigated for SENTINEL2 bands (left) and distribution of the residuals when considering 7 soils.....	27
Figure 8. The 7 standard soil reflectance spectra in SENTINEL2 spectral domain as measured by Liu et al (2002).....	28
Figure 9. Distribution of brightness values as observed over Liu’s data base (Liu et al. 2002) when using 7 reference spectra.....	28
Figure 10. Distribution of leaf area index from the 1008 ground LAI measurements provided by Scurlock et al (2002) .....	29
Figure 11. Effective LAI distribution over the VALERI sites .....	30
Figure 12. Distribution of the average leaf angle (ALA) over the VALERI sites.....	31
Figure 13. Relative leaf size (hot spot parameter) as a function of plant densities from (Baret et al. 2009). Typical values for 3D simulated scenes (crosses) and for actual canopies (circles): Wheat (Wh), Maize (Mz), Sorghum (Sg), Vineyards (Vy), Sunflower (Sf), Soybean (Sy), Tomatoe (To), Olive tree (Ol), Peach (Pe), Lettuce (Lt) and Beet (Bt). ....	31
Figure 14. Scheme explaining the way the distribution of variable $V$ is linked to that of LAI. ....	32
Figure 15: Distribution of values of observational characteristics used to simulate the learning data base.....	34
Figure 16. Distribution of the soil brightness ( $B_s$ , left) and soil type (right, $I_{Soil}$ ) used. ....	34



Figure 17. Distribution used for canopy variables LAI, ALA and Hot (HsD). ..... 35

Figure 18. Distribution of the leaf characteristics used. .... 35

Figure 19. Distribution and co-distributions of the input soil, canopy and leaf characteristics used to populate the training data base..... 36

Figure 20: Distribution and co-distribution of the training database simulated reflectances, in the 8 SENTINEL2 bands..... 37

Figure 21. Distribution and co-distribution of the output variables. .... 38

Figure 22: On the left, relationship between LAI and FAPAR as established over the ensemble of VALERI sites. Symbols correspond to the sites. The solid line corresponds to the best fit model. On the right, symbols correspond to MODIS C4 relationship between LAI and FAPAR as simulated by the 3D radiative transfer model for the several vegetation types considered. The several lines correspond to the distribution within the learning data base. The line with triangles represents a threshold under which values are not expected..... 39

Figure 23. Neural network architecture developed for the estimation of the biophysical variables considered from the 9 SENTINEL2 bands and the 3 angles defining the geometry of observation. The network is made of 1 hidden layer of 5 neurons and 1 linear output neuron. The 'Norm' symbols correspond to the normalization process as described by Equation 10. Symbols 'S' and 'L' correspond respectively to the sigmoid (tansig) and linear transfer functions of the neurons..... 40

Figure 24. Theoretical performances of the neural networks for LAI on the test database..... 42

Figure 25. Theoretical performances of the neural networks for FAPAR on the test database ..... 42

Figure 26. Theoretical performances of the neural networks for FCOVER on the test database ... 43

Figure 27. Theoretical performances of the neural networks for CCC on the test database ..... 43

Figure 28. Theoretical performances of the neural networks for CWC on the test database ..... 44

Figure 29. Schematic representation of the convex hull in the case of 2 dimensional inputs of the network. The convex hull is approximated by a regular grid. The cells with ones (the gray cells) correspond to cases of possible input range (i.e. inside the convex hull) while cells with zeros correspond to inputs out of range (outside the convex hull). ..... 45

## List of Tables

Table 1. SENTINEL2 orbit characteristics and maximum scan angle .....	12
Table 2. SENTINEL2 spectral characteristics: band centre and width, spatial resolution and use.	13
Table 3. SENTINEL2 spectral characteristics: band centre and width, spatial resolution of the 8 selected bands used .....	18
Table 4: Minimum, maximum values and associated resolution for the five products. ....	20
Table 5. Distribution of the input variables of the radiative transfer model used to generate the training data base. Truncated Gaussian, log-normal or uniform distribution laws are used, characterized by the mode, the standard deviation (s), and minimum and maximum values. The number of classes for each variable is presented (Nb_Class).....	33
Table 6. Values used for the co-distributions of the variables with LAI (see Figure 14). ....	33
Table 7. Characteristics of the uncertainties model used. ....	37
Table 8. Minimum and maximum input values (bounding box) of the definition domain for inputs from SENTINEL2 TOC reflectance. ....	45
Table 9. Tolerance, and Minimum ( $P_{min}$ ) and maximum ( $P_{max}$ ) values admitted for the products. ...	46

# 1 BACKGROUND OF THE DOCUMENT

## 1.1 EXECUTIVE SUMMARY

This ATBD (Algorithm Theoretical Based Document) describes the proposed algorithm for Level 2 products that will be implemented in the SENTINEL2 Toolbox. The level2 products are derived from SENTINEL2 top of canopy normalized reflectance data and correspond to the following set of biophysical variables: *LAI*, (Leaf Area Index) FAPAR (fraction of absorbed photosynthetically active Radiation) and FVC that are essential climate variables (ECVs) as recognized by international organizations such as GCOS and GTOS.

The proposed algorithm is based on methods that have already been proven to be efficient. They have been implemented to generate biophysical products from VEGETATION, MERIS, SPOT, and LANDSAT sensors. It mainly consists in generating a comprehensive data base of vegetation characteristics and the associated SENTINEL2 top of canopy (TOC) reflectances. Neural networks are then trained to estimate the canopy characteristics from the TOC reflectances along with set corresponding angles defining the observational configuration.

This ATBD is derived from the ones proposed in the frame of previous ESA projects:

- S2SPAD, Contract ESRIN #21450/08/I-EC: first ATBD delivered, theoretical performances only.
- VALSE2 (VALidation of SEntinel 2 – ESTEC AO/1-6958/11/NL/BJ): validation of the S2SPAD algorithm over experiment ground campaigns (airborne acquisitions to simulate S2-like data). Improvements of the algorithm (inclusion of acquisition geometry and illumination conditions)
- SL2P (Simplified L2 Product Prototype Processor – ESTEC AO/1-7455/13/NL/BJ): refinement of the optimal set of band inputs with regards to atmospheric noise, set up of a definition domain flag, and uncertainty estimates.

The modification with regards to SL2P latest version concern the fact that the actual SENTINEL2 spectral sensitivity (in replacement to the waveband centre values) were taken into account when generating the training database. We also compared two ways of taking into account the geometry of acquisition: either as inputs of the neural nets, or by interpolating values between multiple neural nets trained on specific geometries of acquisition. We concluded that the two methods led to similar results.

## 1.2 SCOPE AND OBJECTIVES

SENTINEL2 is part of the GMES space segment. Users include the scientific community as well as other stakeholders including policy makers, the proper information required for several applications, as detailed in the Mission Requirement Document (Gascon and Berger 2007). For the exploitation of MSI (Multi Spectral Instrument) data, ESA develops the Sentinel-2 Toolbox which consists of a rich set of visualisation, analysis and processing tools products. The objective of this document is to provide a detailed description and justification of the algorithm proposed for the SENTINEL2 Toolbox level2 biophysical variables algorithm.



## 1.3 CONTENT OF THE DOCUMENT

This ATBD document is split in 3 main sections:

1. **Algorithm Justification and overview.** This section contains:
  - A definition of the proposed products.
  - A brief description of SENTINEL2 sensor from which the products will be derived
  - A justification of the algorithm selected
  - The outline of the algorithm.
2. **Description of the algorithm.** This section contains:
  - The required inputs and outputs provided by the algorithm.
  - The retrieval technique used: neural network techniques constitute the core of the operational algorithm. Quality indicators are also provided.
3. **Recommendations for the algorithm implementation.** This section contains details on how the algorithm could be implemented in the processing chain.

## 1.4 SYMBOLS AND ACRONYMS

µg	Microgram
ALA	Average Leaf Angle
ANN	Artificial Neural Network
ATBD	Algorithm theoretical based Document
BHR	Bi-Hemispherical Reflectance
BRF	Bidirectional reflectance factor
BRDF	Bidirectional Reflectance Distribution Function
Bs	Soil Brightness
Cab	Chlorophyll content in the leaf ( $\mu\text{g}\cdot\text{cm}^{-2}$ )
Cbp	Content of brown pigments in the leaf (no units)
CCC	Canopy Chlorophyll Content
Cdm	Content of dry matter in the leaf ( $\text{g}\cdot\text{cm}^{-2}$ )
CEOS	Committee for Earth Observation Satellite
ECV	Essential Climate Variable
FAPAR	Fraction of Absorbed Photosynthetically Active Radiation
CWC	Canopy Water Content
FVC	Fraction of vegetation cover
GCOS	Global Climate Observation System
GMES	Global Monitoring of Environment and Security
GTOS	Global Terrestrial Observation System

L2	Level 2 product
L3	Level 3 product
LAI	Leaf Area Index
MODIS	Moderate Imaging Spectrometer
N	Structural parameter of the leaf (unitless)
NIR	Near Infrared
NNT	Neural Network Technique
RMSE	Root Mean Square Error
RTM	Radiative Transfer Model
SPOT	Satellite Pour l'Observation de la Terre
SWIR	Short Wave Infra red
TOA	Top of Atmosphere
TOC	Top of Canopy
VI	Vegetation Index

## 2 ALGORITHM JUSTIFICATION AND OVERVIEW

### 2.1 THE CONSIDERED PRODUCTS AND THEIR DEFINITIONS

The considered products correspond to the actual vegetation biophysical variables defined below.

#### 2.1.1 Leaf Area Index (LAI)

LAI is defined as half the developed area of photosynthetically active elements of the vegetation per unit horizontal ground area. It determines the size of the interface for exchange of energy (including radiation) and mass between the canopy and the atmosphere. This is an intrinsic canopy primary variable that should not depend on observation conditions. *LAI* is strongly non linearly related to reflectance. Therefore, its estimation from remote sensing observations will be strongly scale dependent (Garrigues et al. 2006a; Weiss et al. 2000). Note that vegetation *LAI* as estimated from remote sensing will include all the green contributors, i.e. including understory when existing under forests canopies. However, except when using directional observations (Chen et al. 2005), *LAI* is not directly accessible from remote sensing observations due to the possible heterogeneity in leaf distribution within the canopy volume. Therefore, remote sensing observations are rather sensitive to the 'effective' leaf area index, i.e. the value that would produce the same remote sensing signal as that actually recorded, while assuming a random distribution of leaves. The difference between the actual *LAI* and effective *LAI* may be quantified by the clumping index (Chen et al. 2005) that roughly varies between 0.5 (very clumped canopies) and 1.0 (randomly distributed leaves).

#### 2.1.2 FAPAR

FAPAR corresponds to the fraction of photosynthetically active radiation absorbed by the canopy. The FAPAR value results directly from the radiative transfer model in the canopy which is computed instantaneously. It depends on canopy structure, vegetation element optical properties and illumination conditions. FAPAR is very useful as input to a number of primary productivity models based on simple efficiency considerations (Prince 1991). Most of the primary productivity models using this efficiency concept are running at the daily time step. Consequently, the product definition should correspond to the daily integrated FAPAR value that can be approached by computation of the clear sky daily integrated FAPAR values as well as the FAPAR value computed for diffuse conditions. To improve the consistency with other FAPAR products that are sometimes considering the instantaneous FAPAR value at the time of the satellite overpass under clear sky conditions (e.g. MODIS), a study was proposed to investigate the differences between alternative FAPAR definitions (Baret et al. 2003). Results show that the instantaneous FAPAR value at 10:00 (or 14:00) solar time is very close to the daily integrated value under clear sky conditions.

FAPAR is relatively linearly related to reflectance values, and will be little sensitive to scaling issues. Note also that the FAPAR refers only to the green parts (leaf chlorophyll content higher than  $15\mu\text{g}\cdot\text{cm}^{-2}$ ) of the canopy.

#### 2.1.3 Cover fraction (FVC)

It corresponds to the gap fraction for nadir direction. FVC is used to separate vegetation and soil in energy balance processes, including temperature and evapotranspiration. It is computed from the leaf area index and other canopy structural variables and does not depend on variables such as the geometry of illumination as compared to *FAPAR*. For this reason, it is a very good candidate for the replacement of classical vegetation indices for the monitoring of green vegetation. Because of its quasi-linear relationship with reflectances, *FVC* will be only marginally scale dependent

(Weiss et al. 2000). Note that similarly to *LAI* and *FAPAR*, only the green elements (leaf chlorophyll content higher than  $15\mu\text{g}\cdot\text{cm}^{-2}$ ) will be considered.

### 2.1.4 Canopy chlorophyll content (CCC)

The chlorophyll content is a very good indicator of stresses including nitrogen deficiencies. It is strongly related to leaf nitrogen content (Houlès et al. 2001). This quantity can be calculated both at the leaf level and at the canopy level by multiplication of the leaf level chlorophyll content by the leaf area index. In this case it is obviously an intrinsic secondary variable. Recent studies tend to prove that this product could be of very high interest in primary production models because it partly determines the photosynthetic efficiency (Green et al. 2003). In addition, studies have demonstrated that a direct estimation of *CCC* is more robust and accurate than an estimation based on the product of the individual estimation of *LAI* and  $C_{ab}$  (Weiss et al. 2000). Therefore, the estimation of *CCC* has been preferred to that of the leaf chlorophyll content.

### 2.1.5 Canopy water content (CWC)

Since radiation is absorbed significantly by water in the near and middle infrared, the spectral configuration of SENTINEL2 allows accessing this variable. Water represents between 60 % and 80% of the living plant mass. The variable that is the best related to the remote sensing signal is defined as the mass of water per unit ground area ( $\text{g}\cdot\text{m}^{-2}$ ). One of the difficulties in retrieving this variable is the possible confusion with soil moisture effects. Canopy water content (CWC) is proposed here as a possible candidate in the list of SENTINEL2 products.

## 2.2 INSTRUMENT CHARACTERISTICS

SENTINEL2 is an optical sensor aboard a polar platform providing helio-synchronous observations. The characteristics of the orbit are presented in Table 1.

Acronym	Values
Orbit altitude (km)	786
Repeat cycle (days)	10
Period (min)	100.7
Inclination (°)	98.62°
Equatorial descending node crossing time (hr)	10:30
Maximum scan angle	20.6°

**Table 1. SENTINEL2 orbit characteristics and maximum scan angle**

The 13 characteristics of the 13 bands are presented in Table 2. Bands B1, B2, B9 and B10 are more dedicated to atmosphere or cloud. The 9 remaining bands may be used for vegetation characterization. They are all with 20 m spatial resolution, except B4 and B8 that have a 10 m resolution. This will thus need binning these bands to provide the 20 m spatial resolution. The detailed spectral characteristics are presented in Figure 1.

Acronym	Central (nm)	Width (nm)	Spatial resolution (m)	Potential Applications
B1	443	20	60	Atmosphere
B2	490	65	10	Atmosphere
B3	560	35	10	Vegetation
B4	665	30	10	Vegetation
B5	705	15	20	Vegetation

B6	740	15	20	Vegetation
B7	783	20	20	Vegetation
B8	842	115	10	Vegetation
B8a	865	20	20	Vegetation
B9	945	20	60	Atmosphere
B10	1375	30	60	Atmosphere
B11	1610	90	20	Vegetation
B12	2190	180	20	Vegetation

Table 2. SENTINEL2 spectral characteristics: band centre and width, spatial resolution and use.

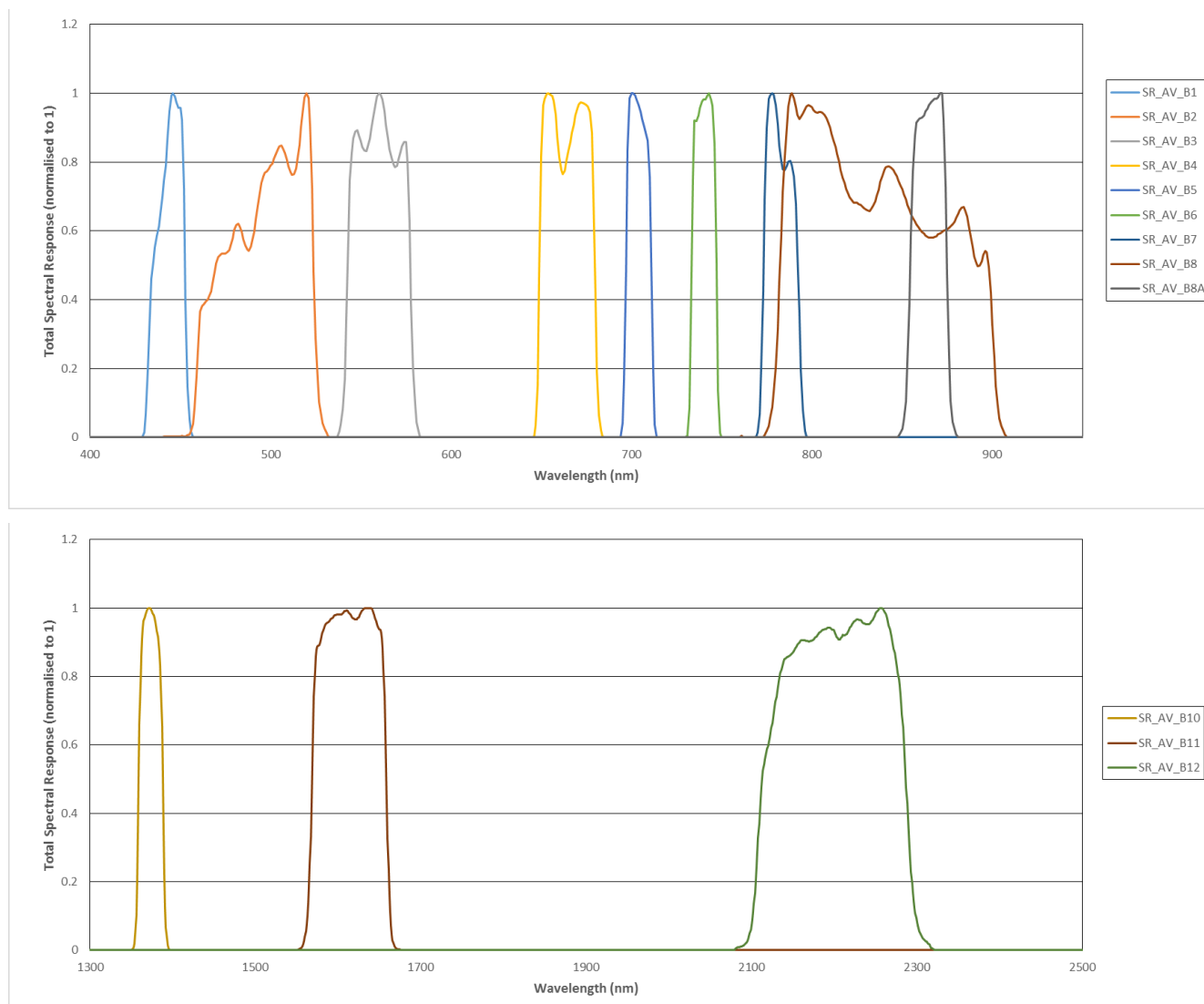


Figure 1. SENTINEL2A MSI spectral response (Top: Visible – near Infrared, Bottom: ShortWave Infrared) from ESA SENTINEL [online web site](#).

## 2.3 JUSTIFICATION FOR THE ALGORITHM SELECTION AND DESIGN

The requirements for the selection and design of the algorithm proposed in this study for SENTINEL2 level 2 products are presented below:

- **Explicit use of all the SENTINEL2 pertinent spectral information.** The spectral sampling of SENTINEL2 provides potentially a higher level of information on canopy structure and optical properties of its elements as compared to the simple use of the classical red and near infrared bands implemented in most other retrieval approaches. The exploitation of the adequate SENTINEL2 spectral information should hopefully allow to restrain the solution space and lead to a more robust and accurate retrieval as compared to high resolution sensors such as SPOT or LANDSAT.
- **Accuracy of the retrieval and computational efficiency.** A review of current state of the art for the estimation of biophysical variables from remote sensing data (Baret and Buis 2007) was proposed. It shows that among the several retrieval algorithms, those based on the minimisation of the distance in the space of canopy variables provides the best accuracy on retrievals while being very efficient computationally wise. Therefore, techniques based on neural networks are selected in this study. In addition, their limitation mainly driven by the necessity to have a fixed number of input variables would not constitute any problem to process SENTINEL2 data up to level 2, if the geometrical configuration is input explicitly. Note that such techniques have already been implemented and lead to good performances of the retrieval (Weiss et al. 2002); (Baret et al. 1997); (Combal et al. 2002); (Kimes et al. 2002) (Bacour et al. 2006). It requires a learning process achieved over a training data base.
- **No use of ancillary data difficult to derive at high spatial resolution.** As a matter of fact, retrieval of surface characteristics from remote sensing observations is an ill posed problem, leading to uncertainties in the solution (Combal et al. 2002). The only way to regularize the problem is to use ancillary information or additional constraints on the solution to reduce its domain of variation. Spatial and temporal constraints may be used to regularize the solution as demonstrated by (Lauvernet et al. 2008) and (Kötz et al. 2005; Kötz et al. 2007) because we are mainly considering LEVEL2 products. Further, because of the generic (global) nature of the proposed vegetation products, no specific ancillary information can be efficiently used: the absence of reliable and regularly updated land cover map at a spatial resolution similar to that of SENTINEL2 from which some vegetation architecture features could be derived prevents from tailoring specific algorithm for each of the land cover classes. The proposed algorithm will therefore rely only on SENTINEL2 instantaneous observations. As a consequence, the algorithm may provide reasonably good estimates over all the cases, but certainly poorer performances as compared to algorithm specific to a given surface type for this particular surface type, although this specific algorithm would fail over most of other cases. Nevertheless, it could be possible to develop such specific algorithms and use them when the land cover is known. It would be also possible to develop corrections of the generic algorithm that will be specific to a surface type. Note also that the differences between a specific and a generic algorithm will depend on the considered variables: large differences are expected for LAI, while less sensitivity is expected for FAPAR and FVC. The genericity of the algorithm will obviously have important consequences on the generation of the training data base.
- **Generation of the training data base.** The training data base should sample all the vegetation types and conditions that can be observed from SENTINEL2 over land surfaces. In addition it should reflect the uncertainties in the reflectance values as observed by SENTINEL2. Ideally, the training data base should therefore be made of SENTINEL2 observations that are paired with accurate ground measurements of the considered biophysical variables. However, because of the uncertainties attached to the ground measurements and the difficulty associated to the collection of such measurements within a large range of vegetation types and conditions, this simple 'experimental' approach is currently not feasible. Therefore, the use of simulations by radiative transfer models would be preferable. The radiative transfer model should simulate within a good accuracy the



canopy reflectance as observed within SENTINEL2 bands and geometry over most vegetation types and conditions that can be observed over the Earth. A particular attention should be brought on:

- the leaf optical properties, particularly regarding the effect of the chlorophyll and water content on reflectance and transmittance,
  - the background reflectance that should include in addition to a large variety of soils, litter and senescent vegetation.
- **Quality assessment.** Quantitative and qualitative indicators should be attached to the product so that the user could properly ‘weigh’ the data within its application according to the confidence he puts on. This could be achieved within several ways:
  - Quality of the TOC reflectance used as input to the algorithm. This would simply correspond to the replication of indicators produced previously such as cloud occurrence, sensor problem and atmospheric correction.
  - Additional indicators based on:
    - The input out of range. This should indicate that either the input reflectances have problems (cloud contamination, poor atmospheric correction, shadow) or anyway that the application of the algorithm could result in unreliable results. This would be common to all the derived products.
    - Output out of range: flags raised when the product appears to be out of the nominal range of variation.
    - Product uncertainty. The algorithm provides a quantitative estimation of the uncertainty associated to the product.

## 2.4 ALGORITHM OUTLINE

For the reasons exposed above, ***it is proposed to use neural networks for the SENTINEL2 biophysical variables estimation.*** Once the neural network inputs and outputs are defined, it is very easy to change the network coefficients to take advantage of recent advances in the generation of the training data base. The processing chain will therefore be very easily upgraded if it designed for, i.e. if coefficients are stored as parameters in a well identified and documented manner.

For each product, one particular network will be calibrated. Two main steps are foreseen (Figure 2):

- Training the neural network.
- Operational use of the neural network.

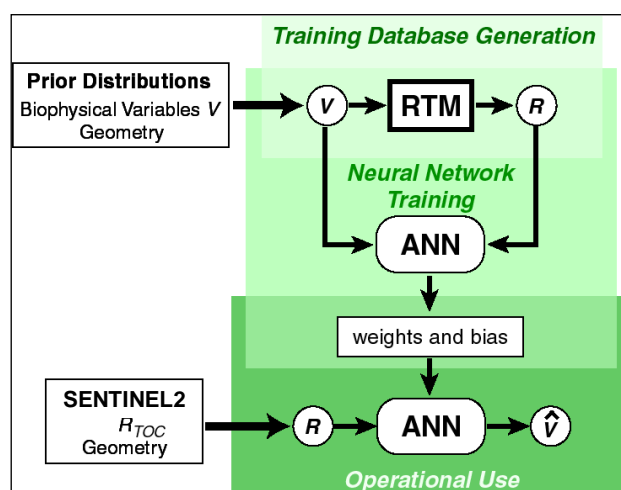


Figure 2. Flow chart showing how the products ( $\hat{V}$ ) are generated operationally. ANN corresponds to an Artificial Neural Network characterized par its structure and its coefficients (synaptic weights and bias);  $R_{TOC}$  corresponds to the SENTINEL2 Top Of Canopy reflectance used in the operational mode and  $V$  correspond to the biophysical variable in the training data base and estimated by running the ANN over the simulated SENTINEL2 TOC reflectance and geometry.

### 2.4.1 Training the neural network

This process consists mainly in three steps:

- Generating training data base
- Defining the neural network architecture
- Calibrating the network.

#### 2.4.1.1 Generation of the training data base

The generation of the training data base corresponds to the most critical issue to be solved. It should be constituted of a representative set of top of canopy reflectances and incorporate prior information on the distribution of the input variables. The model used to simulate the learning database should represent a good compromise between the level of accuracy and the complexity of setting up the simulations: number and distribution of the input variables to describe the canopy as well as computer ressources and time.

#### 2.4.1.2 Designing network architecture

It consists in defining the optimal structure (typically the number of layers and the number of neurons per layers) as well as possible transformations of the inputs and outputs such as normalization.

#### 2.4.1.3 Calibrating the network

This last step corresponds to the actual training, i.e. tuning the coefficients (synaptic weights and bias) that provide the best estimates of the biophysical variables. Dedicated tools are available to achieve this training, and this issue will be detailed later on.

### 2.4.2 Operational use of the neural network

Once the neural network is trained, it will be run in operational mode. 3 networks will produce in parallel estimates of the considered biophysical variables: LAI, FAPAR, Fcover. Additionally, quality assessment indicators will also be generated:



- **Input consistency with the training data base.** This represents the consistency of the measured SENTINEL2 input reflectances with those used in the training data base. The training definition domain of the inputs is therefore identified, and a flag will be raised when observations are outside the training definition domain.
- **Output consistency with expected range.** This represents the consistency of the actual network outputs (the biophysical variables) with those used in the training data base.
- **Quality indicators:** These are a replication of the previously computed quality indicators, including those related to the atmospheric correction and cloud filtering.

## 3 ALGORITHM DESCRIPTION

In this section, the algorithmic elements are described, including:

- The definition of the inputs and outputs,
- The radiative transfer model used to generate the learning data base
- The inversion technique
- The quality assessment

### 3.1 INPUTS AND OUTPUTS

#### 3.1.1 Inputs

The neural networks will apply on instantaneous top of canopy reflectance data (level 2). All the following inputs are required for each considered pixel.

##### 3.1.1.1 *Top of canopy reflectance.*

Individual daily observations will be used. Reflectances should be expressed in terms of reflectance factor, mainly varying between 0 and 0.7 for most land surfaces outside hot-spot or specular directions and snow or ice cover. Only nine bands are used: B3, B4, B5, B6, B7, B8a, B11 and B12. This band selection is derived from the VALSE2 and SL2P project results.

Acronym	Central (nm)	Width (nm)	Spatial resolution (m)
B3	560	35	10
B4	665	30	10
B5	705	15	20
B6	740	15	20
B7	783	20	20
B8a	865	20	20
B11	1610	90	20
B12	2190	180	20

**Table 3. SENTINEL2 spectral characteristics: band centre and width, spatial resolution of the 8 selected bands used**

##### 3.1.1.2 *Geometry of acquisition*

Similarly to the spectral information content of the reflectance signal, the directional information must also be taken into account when training the neural networks. The cosine of the sun zenith angle ( $\theta_s$ ), view zenith angle ( $\theta_v$ ) and relative azimuth angle ( $\varphi$ ) at the time of the image acquisition are required. The geometry of acquisition can be taken into account in two ways when using the neural nets:

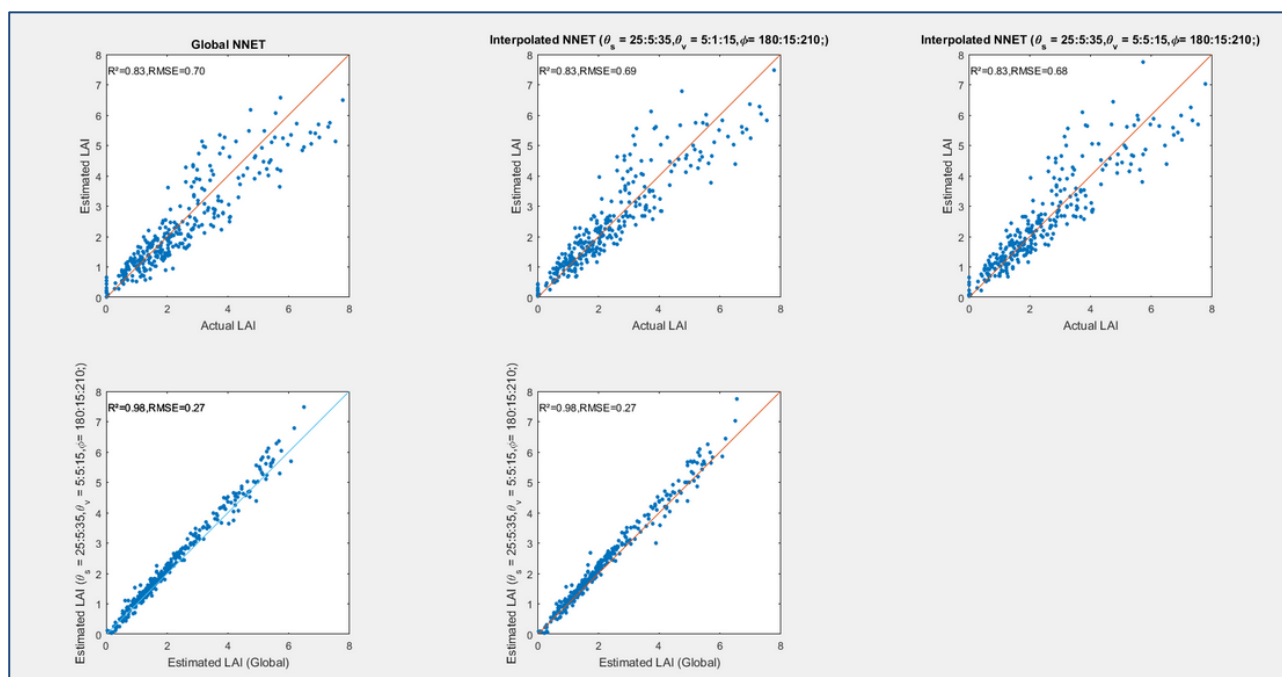
1. Global Neural Network: the cosine of the sun zenith angle ( $\theta_s$ ), view zenith angle ( $\theta_v$ ) and relative azimuth angle ( $\varphi$ ) are used as inputs to the neural networks. The simulations of the training data base must take into account the orbit characteristics of SENTINEL2 to simulate the satellite angular sampling.

2. Interpolated Neural Networks: the geometry of acquisition is not used as input to the neural networks: it is taken into account by training neural networks with simulations run for fixed values of the triplet  $(\cos(\theta_s), \cos(\theta_v), \cos(\varphi))$ : for each variable, a number N of neural networks is calibrated, each of them trained with a single geometric condition. When deriving the biophysical variable for a given sensor acquisition, the different neural networks are run, and the final product is estimated by interpolating the different outputs corresponding to fixed geometries into the actual sensor view angles.

We compared the results between these two implementations for LAI estimation. All the neural networks were trained with the same drawing of PROSAIL input variables (Table 5 and Table 6), except for the angles. For each simulation, the global neural net was trained with a uniform random drawing of the date and location, allowing to derive, thanks to the orbit characteristics of SENTINEL2, the corresponding geometry. For the interpolated Neural Networks, we trained 99 NNets (respectively, 27) with a single geometry configuration corresponding to sun zenith angles between  $25^\circ$  and  $35^\circ$  by step of  $5^\circ$ , view zenith angles between  $5^\circ$  and  $15^\circ$  by steps of  $1^\circ$  (respectively  $5^\circ$ ), and relative azimuth angle between  $180^\circ$  and  $210^\circ$  by steps of  $5^\circ$ , which is in agreement with S2 geometry (see section 3.3).

Then, we simulated an independent test database using the same radiative transfer model (PROSAIL) with a uniform random drawing of the date and location. From one hand, we applied the global neural network, by using the actual geometry as input. From the other hand, we applied the 99 (respectively 27) NNets and interpolated the obtained values at the actual geometry value (Figure 3). Results show that the performances of the 3 methods are comparable (Correlation coefficients of 0.83 for the 3, RMSE between 0.68 and 0.70). Moreover, the scatterplots between the global NNets and the interpolated ones shows a very good agreement ( $R^2=0.98$ ), with few points a little more scattered when degrading the resolution of the view angle. This result also show that the accuracy required on the view angle is not critical. As the results are very similar, we conclude that it is much simpler to use the global NNet for the implementation in the SENTINEL2 toolbox.

However, the determination of the viewing angles for a given pixel may not be currently an easy task when using current SENTINEL2 data. Indeed, there is one viewing configuration associated to each of the 12 detectors for a given band (making  $8 \times 12$  view angles per pixel). The current S2 L1C data format does not provide the information on the detector associated to the reflectance in case of overlapping between two detectors. Therefore, it is required to use an approximation of the view angle, which will not degrade the results as shown in (Figure 3, bottom). We propose to use the same approximation as in the sen2cor processor that is also implemented in the S2Toolbox to perform atmospheric corrections.



**Figure 3. Comparison between the performances of a global neural network versus interpolated ones. Top: theoretical performances against a simulated test database, Bottom: scatterplot between interpolated and global neural network LAI estimates in the test database. Two interpolated neural networks are considered: 1°step or 5° step.**

### 3.1.2 Outputs

The outputs will be provided by application of the algorithm over each pixel. They include the neural network derived LAI, FAPAR, FVC, CCC and CWC values as described previously. The proposed range of variation and resolution are presented in Table 4. In addition to the product values, quality flags are also generated.

Product	Unit	Minimum	Maximum	resolution
LAI	$mP^{2P} \cdot mP^{-2P}$	0	8.0	0.01
FAPAR	-	0	1.0	0.01
FVC	-	0	1.0	0.01
CCC	g/cm <sup>2</sup>	0	600	1
CWC	µg/cm <sup>2</sup>	0	0.55	0.0025

**Table 4: Minimum, maximum values and associated resolution for the five products.**

## 3.2 REFLECTANCE MODELS

As previously described, the algorithm is based on the training of neural networks for each product. Emphasis will therefore be put on the generation on the training data base with attention on the assumptions associated to the used models. Physically based radiative transfer models are considering 3 main components that will be described separately in the following:

- The leaf optical properties
- The canopy structure
- The background reflectance

### 3.2.1 Leaf optical properties

To estimate the chlorophyll content from canopy reflectance, chlorophyll content has to be explicitly introduced into the radiative transfer model to be used. Because of its versatility and performances, the PROSPECT model (Jacquemoud and Baret 1990) with the updated absorption coefficients proposed by (Fourty and Baret 1997) appears therefore to be a good candidate.

Note that the PROSPECT model considers the leaf as a lambertian surface. Sanz et al (1997) showed that leaves were mainly characterized by a specular behaviour in addition to an important diffuse scattering process that takes place within the leaf. These authors demonstrated that, except in the specular direction, the lambertian approximation was valid in all other viewing directions. In addition, PROSPECT assumes that the optical properties of both leaf faces are equal.

Several authors (Fourty and Baret 1997; Jacquemoud and Baret 1990; Newnham and Burt 2001) have successfully validated the model over broadleaf types. In addition, the PROSPECT model provides a reasonable description of the optical properties of the needles, even though the basic assumptions associated to the plate model are obviously violated (Zarco-Tejada et al. 2001). The following variables are required as input to the PROSPECT model:

- $N$  leaf mesophyll structure index. It varies between 1.0 for the most compact leaves (such as young cereal leaves) up to 3.5 for thick leaves with well developed spongy mesophyll or event senescent leaves having disorganized mesophyll with large amount of air spaces.
- $C_{ab}$  Leaf Chlorophyll content ( $\mu\text{g}\cdot\text{cm}^{-2}$ ). It actually corresponds to the content of chlorophyll a, chlorophyll b and carotenoids (Fourty and Baret 1997). Note that chlorophyll a and b are generally strongly correlated. The same is observed between chlorophyll a and b and carotenoids, particularly for medium to large chlorophyll content values. It basically varies between 0 to  $90 \mu\text{g}\cdot\text{cm}^{-2}$ .
- $C_{dm}$  Leaf dry matter content ( $\text{g}\cdot\text{cm}^{-2}$ ). Dry matter absorbs over the whole spectral domain, and its effect is maximal in the near infrared region. The leaf dry matter content is also called the specific leaf weight ( $SLW$ ) which is also the inverse of the specific leaf area ( $SLA$ ) used by physiologists.
- $C_w$  Leaf water content ( $\text{g}\cdot\text{cm}^{-2}$ ). Several studies showed that the relative water content could be approximated to a value close to 75 % for the green leaves. This allows linking the water ( $C_w$ ) and the dry matter ( $C_{dm}$ ) contents together.
- $C_{bp}$ . Leaf brown pigment content (relative units). Baret et al (2002) reported that chlorophyll and brown pigments are exclusive, i.e. green and non green elements (senescent leaves, branches, stems) are spatially dissociated. The canopy structure model should therefore include at least green and non green elements. Green leaves will have no brown pigments and senescent leaves will have no chlorophyll pigments.

Bacour et al (2002) and Le Maire (2002) have analysed the sensitivity of the radiometric response both at the leaf and canopy level. They showed that the chlorophyll content, the dry matter and the structure index are the main drivers of the optical properties in the visible to near infrared spectral domain.

### 3.2.2 Canopy radiative transfer models

The use of pure 3D models such as DART (Gastellu-Etchegorry et al. 1996) or DISORD (Myneni et al. 1992) for simulating a very large range of situations appears very appealing. Even though, the use of detailed 3D models that mimics actual canopy architecture and combined with ray tracing (Govaerts and Verstraete 1998), (España et al. 1999) or radiosity (Gerstl and Borel 1992), (Borel et al. 1991), (Chelle et al. 1997), (Soler et al. 2001) radiative transfer description and applied to a representative sample of biomes and conditions would be ideal. However, it might be difficult to implement for two practical reasons:

- The necessity to describe a very large range of realistic canopy architectures. This requires a huge effort in canopy architecture and optical properties measurements at ground level.
- The time associated to the model computation.

It is thus proposed to use a reflectance model that is computer efficient and uses a small number of input variables. The SAIL radiative transfer model (Verhoef 1984, 1985) is widespread in the remote sensing community for the estimation of vegetation biophysical variables. The canopy is described as a homogeneous medium where leaves are randomly distributed. The SAIL model uses a limited number of structural variables in addition to leaf reflectance and transmittance and soil back ground reflectance.

- Leaf area index (*LAI*),
- the average leaf angle (*ALA*), characterizing the leaf angle distribution that will be described by an ellipsoidal distribution (Campbell 1986). Note that a spherical distribution corresponds to an average angle close to 57°,
- the hot spot parameter (*HOT*) (Kuusk 1991).

### 3.2.3 Background reflectance model

The background reflectance corresponds to all the non green materials that constitute the last bottom layer in the canopy. Following the definition of the *products*, all the green vegetation layers have to be accounted for in the computation of these variables. Therefore, if the understory is green (including lichens and moss), it will not be considered as the background here and will be included within the green vegetation layer. The background reflectance may thus correspond to soil, litter, water and snow. However, because of the particularities of water and snow backgrounds, these will not be included in this study. Indeed, when these particular cases will be encountered when applying the algorithm on SENTINEL2 data, the 'input out of range' quality indicator should be flagged.

#### 3.2.3.1 *The background brightness concept*

The background reflectance, for a given wavelength, will depend on the background type (snow, soil type, litter, water), geometrical illumination and view conditions ( $\Omega$ ), roughness ( $z$ ) or moisture ( $H$ ). Note also that there is a continuum between soil background and water (which is always above soil!).

The approach used here to describe the background reflectance properties is based on the brightness concept allowing confounding the effect of geometrical conditions, roughness and moisture within a single parameter that will be assumed not to depend on wavelength.

The background reflectance  $\rho_b(\lambda, \Omega_i, H_j, z_k)$  for any wavelength  $\lambda$ , observation geometrical configuration  $\Omega_i$ , moisture  $H_j$  and roughness  $z_k$  is assumed proportional to the reflectance



background for the same wavelength  $\lambda$  but different observation geometrical configuration  $\Omega_l$ , moisture  $H_m$  and roughness  $z_n$ :

$$\rho_b(\lambda, \Omega_i, H_j, z_k) = Bs \cdot \rho_b(\lambda, \Omega_l, H_m, z_n) \quad \text{Equation 1}$$

where  $Bs$  is a brightness parameter that does not depend on wavelength  $\lambda$ , but depends on all the other factors ( $\Omega$ ,  $H$ ,  $z$ ). This convenient property is a consequence of the well known soil line concept (Baret et al. 1993) stating that a linear relationship exists between the reflectance of soils (and litter) in two wavelengths  $\lambda_1$  and  $\lambda_2$  when either the roughness, moisture or illumination or view directions vary:

$$\rho_b(\lambda_1, \Omega_i, H_j, z_k) = a(\lambda_1, \lambda_2) \cdot \rho_b(\lambda_2, \Omega_i, H_j, z_k) + b(\lambda_1, \lambda_2) \quad \text{Equation 2}$$

This property could be written for another set of sun and view directions:

$$\rho_b(\lambda_1, \Omega_l, H_m, z_n) = a(\lambda_1, \lambda_2) \cdot \rho_b(\lambda_2, \Omega_l, H_m, z_n) + b(\lambda_1, \lambda_2) \quad \text{Equation 3}$$

Replacing in **Equation 3**  $\rho_b(\lambda_1, \Omega_i, H_j, z_k)$  and  $\rho_b(\lambda_2, \Omega_i, H_j, z_k)$  by their expression derived from Equation 1 and Equation 2:

$$Bs \cdot \rho_b(\lambda_1, \Omega_l, H_m, z_n) = a(\lambda_1, \lambda_2) \cdot Bs \cdot \rho_b(\lambda_2, \Omega_l, H_m, z_n) + b(\lambda_1, \lambda_2) \quad \text{Equation 4}$$

Identifying Equation 4 to Equation 3 provides the condition under which Equation 2 is valid:

$$Bs \cdot b(\lambda_1, \lambda_2) \approx b(\lambda_1, \lambda_2) \quad \text{Equation 5}$$

which is true either for  $Bs \approx 1$  or  $b(\lambda_1, \lambda_2) \approx 0$ . Experimental and theoretical results (Baret et al. 1993) show that the soil line intercept,  $b(\lambda_1, \lambda_2)$ , is generally very small in comparison to the background reflectance value. For example, in the red and near infrared bands,  $0 < b(\lambda_{red}, \lambda_{nir}) < 0.1$ . Similarly, experimental evidences, when referring to a standard situation (dry soil, medium roughness, no hot-spot configuration),  $0.3 < Bs < 1.3$ . Therefore, the brightness concept is generally valid and has already been used extensively in past studies (Weiss et al. 2002); (Bacour, Jacquemoud et al. 2002).

The brightness concept allows describing the spectral variation of a given background when the geometrical configuration, moisture or roughness varies with two inputs:

- The brightness parameter ( $Bs$ ) that is independent on wavelength
- A reference soil reflectance spectrum.

### 3.2.3.2 **Background reference spectral variation**

The background spectral database must represent the different background types with special attention to soils and litter.

Note that the litter corresponds to an important background, particularly over forest areas. The spectral signature of litter is very close to that of the soil as noticed by several studies (Asner et al. 1998). Crop residues and natural vegetation residues may have also important contribution to the reflected signal during specific seasons. Similarly to litter, the reflectance of vegetation residues is also very comparable to that of soil background (Gausman et al. 1975), (Biard and Baret 1997), (Chen and McKyes 1993). Few litter and vegetation residues were measured at the laboratory to get some reasonable representation of these background types. Because of the similarity between litter, residues and soil reflectance, these are finally aggregated within the 'soil' background category.





### 3.3 GENERATION OF THE TRAINING DATABASE

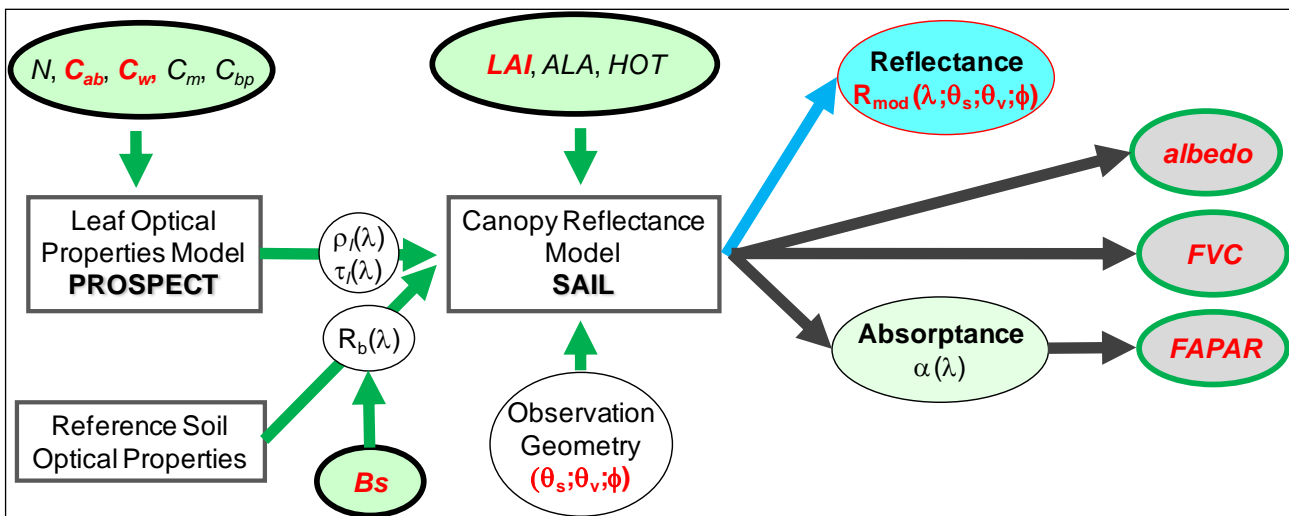
The training data base is generated in three steps:

- Generation of the data base containing the input radiative transfer model variables
- Generation of the corresponding top of canopy reflectance for the 8 SENTINEL2 bands considered
- Addition of uncertainties to the simulated top of canopy reflectance values previously simulated.

#### 3.3.1 Radiative transfer model

The top of canopy reflectances are simulated with the PROSPECT+SAIL model described earlier (Figure 4). The coupled model allows as well the computation of the secondary biophysical variables such as *FVC*, *FAPAR*. The model inputs are:

- the geometrical configuration of illumination and observation (i.e. the solar and view zenith angles,  $\theta_s$  and  $\theta_v$ , and the relative azimuth), derived from the SENTINEL2 orbit characteristics and swath,
- the background reflectance spectrum, as described earlier,
- the primary biophysical variables related to leaf optical properties ( $N$ ,  $C_{ab}$ ,  $C_w$ ,  $C_m$ , and  $C_{bp}$ ) and to the canopy structure ( $LAI$ ,  $ALA$ ,  $HOT$ , and  $B_s$ ). Their associated distribution law will be specified hereafter.



**Figure 4:** The coupled PROSPECT+SAIL model to generate the training database made of TOC reflectances and corresponding biophysical variables.

#### 3.3.2 Distribution of the vegetation input variables

This is the most delicate step in the generation of the training data base. As a matter of fact, the training data base has to reflect the actual distribution of the vegetation types over the Earth's surface. The variable distributions are derived from available information on the actual distribution of the variables. Table 5 presents the range of variation and the actual distribution used for the input variables of the vegetation and background.

### 3.3.2.1 Using results from the literature for leaf characteristics

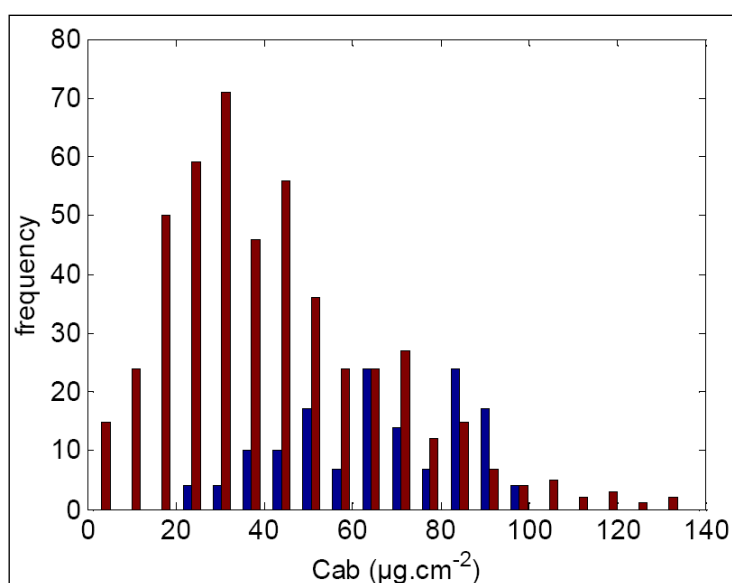
One of the main problems in generating a learning data base is to get the proper distributions for the input model variables if realistic distributions are required. Unfortunately, the distribution of most variables is not precisely known. For leaf optical properties, a brief literature review allows to characterize the distributions for chlorophyll ( $C_{ab}$ ) and dry matter contents ( $C_{dm}$ ).

Results show that  $C_{ab}$  et  $C_{dm}$  distributions are roughly Gaussian (Figure 5 and Figure 6). However, they depend on the vegetation type:

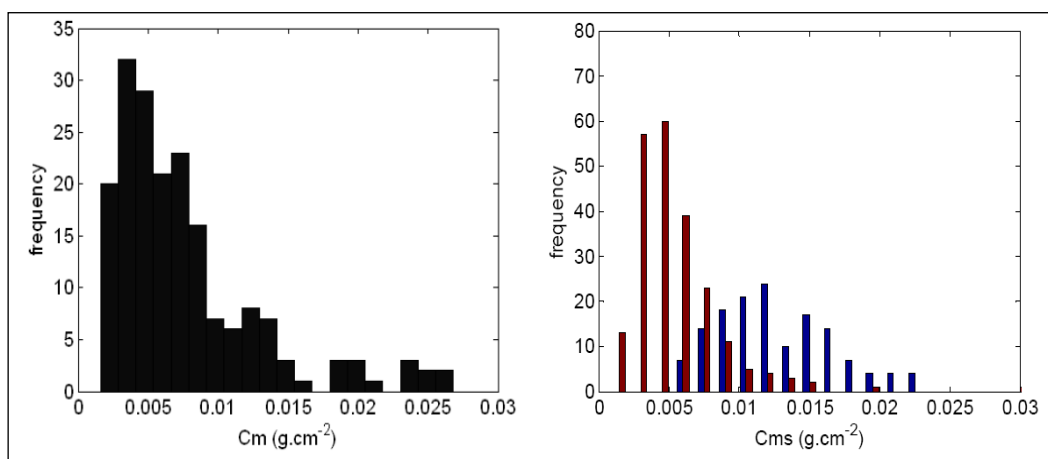
- for chlorophyll (Figure 5), the mode is close to 40  $\mu\text{g}\cdot\text{cm}^{-2}$  for non evergreen vegetation types, while it is close to 70  $\mu\text{g}\cdot\text{cm}^{-2}$  for evergreen broadleaf forests (the Hawaiian data set reported in Feret et al (2008)). Note that the chlorophyll content definition considered here is consistent with the input of the PROSPECT model, i.e. it includes chlorophyll a and b as well as carotenoids.
- for dry matter contents (Figure 6), the mode is close to 0.005  $\text{g}\cdot\text{cm}^{-2}$  for non evergreen broadleaf forests, and close to 0.012 for evergreen broadleaf forests.

The distribution derived from a compilation of literature data (Figure 6 left) with few broadleaf evergreen forest samples is very consistent with the observations of Feret et al (2008) (Figure 6 right).

Note however that the compilation of data presented here might not represent realistically the actual distribution of values around the globe. Anyway, the differences found between vegetation types should indicate that training specific algorithm for each vegetation type might improve the performances.



**Figure 5: Distribution of leaf Chlorophyll content as described by Feret et al (2008). Values of Chlorophyll include chlorophylls a and b as well as carotenoids. Distribution in blue corresponds to data mainly from tropical forests (Hawaii data set in Feret's paper). Distribution in red corresponds to other species.**

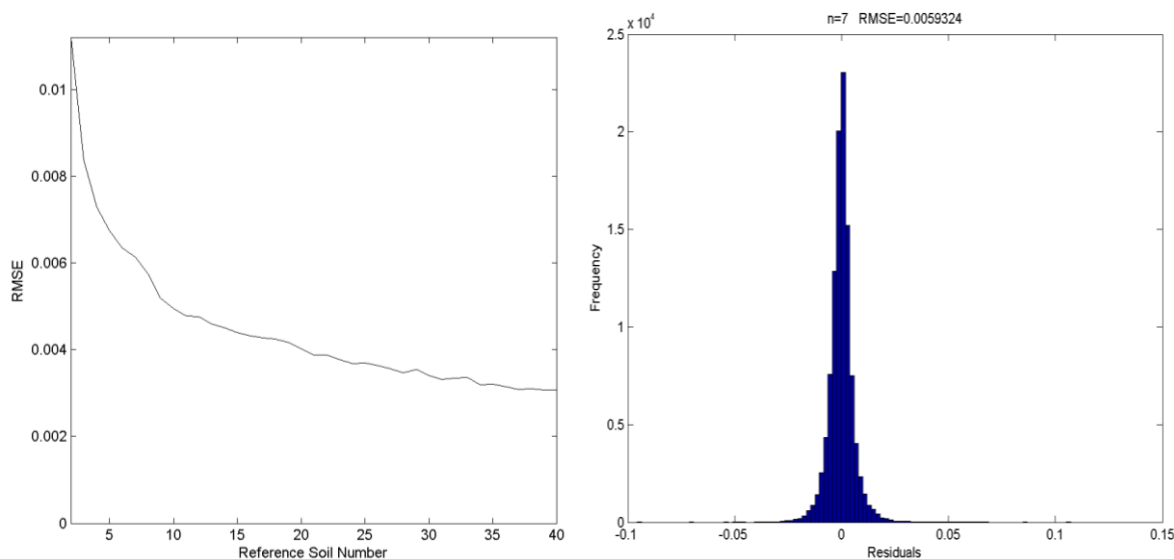


**Figure 6:** On the left, distribution of dry matter contents (Cdm) derived from literature review: (Prado and De Moraes 1997; Poorter and Evans 1998; Reich, Ellsworth et al. 1998; Poorter and De Jong 1999). On the right, distributions reported by Feret et al (2008): distribution in blue corresponds to data mainly from tropical forests (Hawaii data set in Feret’s paper). Distribution in red corresponds to other species.

The water content was tied to the dry matter content assuming that the green leaves have a relative water content close to 75%.

### 3.3.2.2 Soil characteristics

The reference soil spectra will be derived from a soil reflectance data base available at INRA Avignon representing a large variation of soil types, moisture, roughness and geometrical configurations (Jacquemoud et al. 1992; Liu et al. 2002). Considering the brightness concept will allow increasing the diversity in actual soil properties. Only 7 soil spectra were selected among the 1500 measured ones available to represent the range of spectral shapes observed within a good accuracy (Figure 7). The measurements were performed using an ASD Fieldspec Pro spectrophotometer providing a 1nm spectral sampling close to the spectral resolution in the visible and near infrared domains



**Figure 7.** RMSE values associated to the reconstruction of all the soils investigated for SENTINEL2 bands (left) and distribution of the residuals when considering 7 soils.

The corresponding soil spectra are presented in Figure 8. Note that here all the spectra were normalized so that they provide the same reflectance value when averaged over the SENTINEL2 bands.

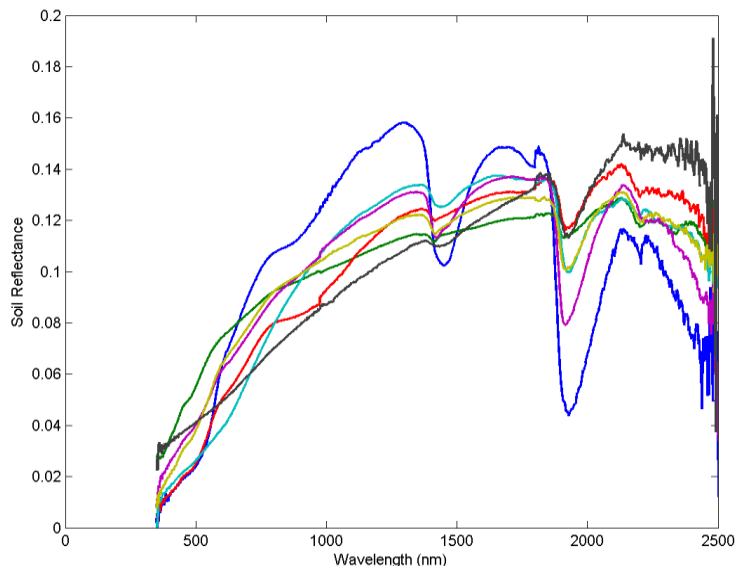


Figure 8. The 7 standard soil reflectance spectra in SENTINEL2 spectral domain as measured by Liu et al (2002).

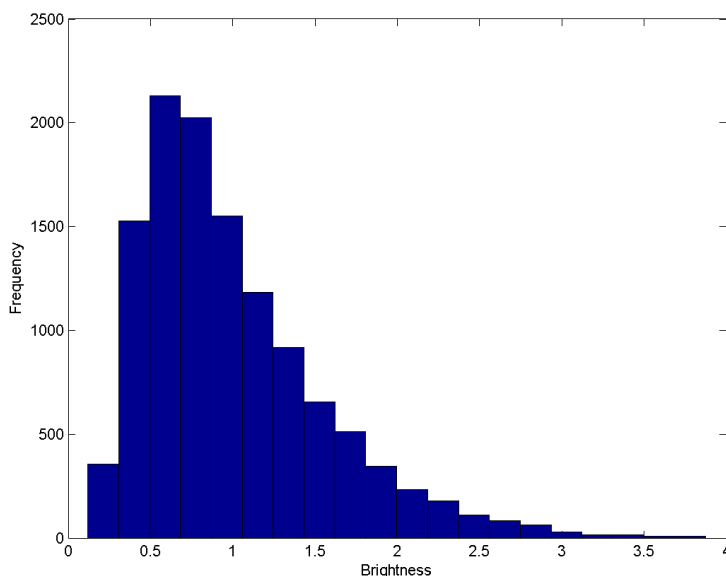


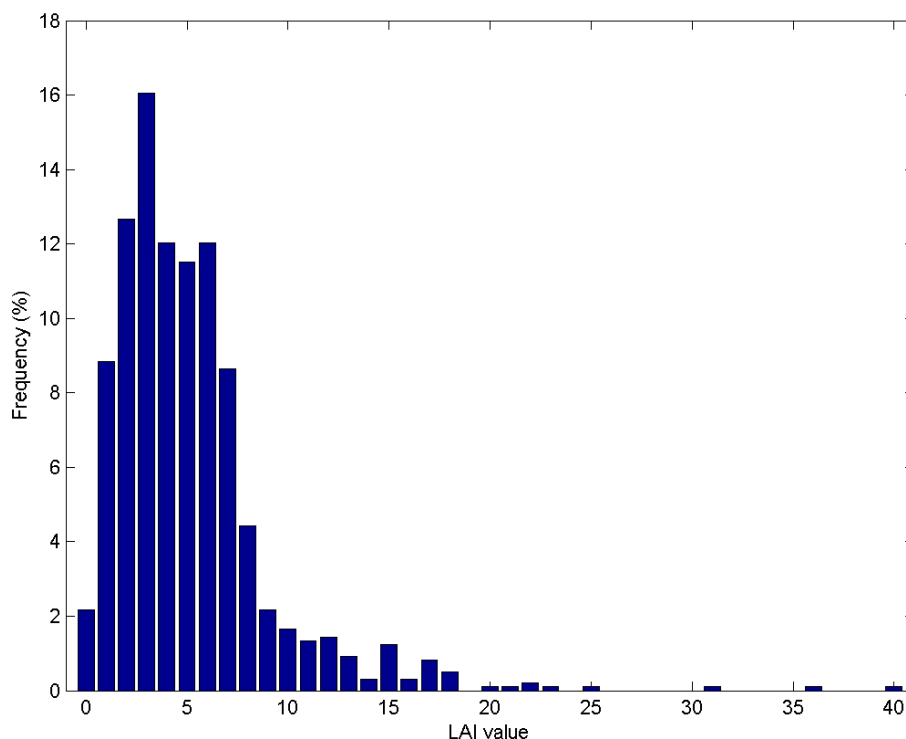
Figure 9. Distribution of brightness values as observed over Liu's data base (Liu et al. 2002) when using 7 reference spectra.

### 3.3.2.3 Canopy variables

#### 3.3.2.1 Leaf Area Index

(Scurlock et al. 2001) have compiled about 1000 LAI experimental values over 300 original-source references. These values correspond to nearly 400 unique field sites during the 1932-2000 period.

Figure 10 shows that the measured LAI is quite uniformly distributed between the 1-7 values while low LAI values (<1) are not frequently measured. On the other hand, some very high but not very frequent LAI values are represented in this data set (up to 40).



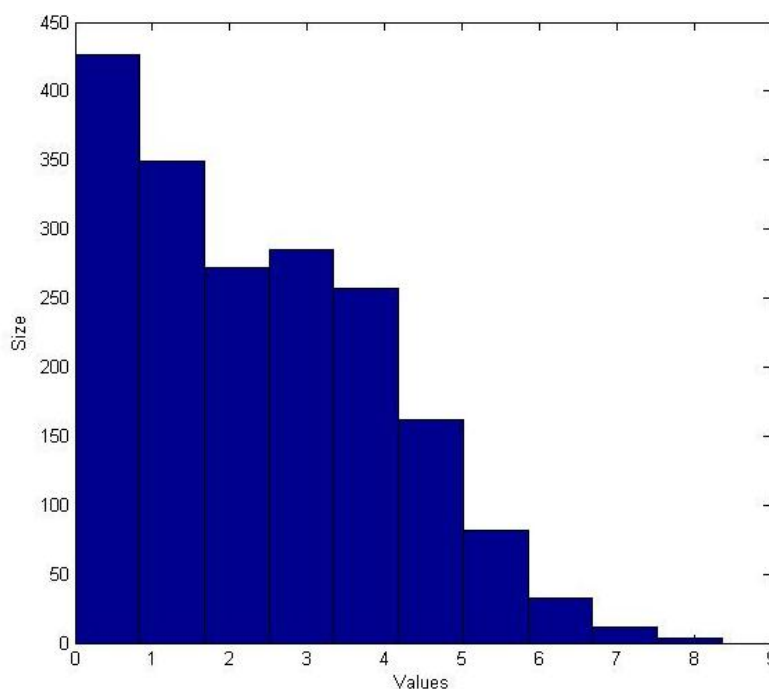
**Figure 10. Distribution of leaf area index from the 1008 ground LAI measurements provided by Scurlock et al (2002)**

However, although this data set represents a large amount of experimental data, it does not represent well the actual LAI distribution over the globe for the following reasons:

- The measurements are achieved when some vegetation is present during the experiments. Therefore, bare soils (LAI=0) are not taken into account here, whereas they represent about 30% of the emerged surface (Baret et al. 2006; Masson et al. 2003).
- The measurements used in this study are performed using different methods and LAI devices leading to different estimates of this variable. They all correspond to LAI but under different assumptions (effective or true LAI) and accuracy (Garrigues et al. 2008; Weiss et al. 2004)
- Finally the distribution derived from this study is certainly not representative of the actual distribution at the Earth surface over a year since the measurements correspond to punctual dates during the vegetation cycle depending on the objectives of the authors of the studies.

We compared the distribution obtained from Scurlock et al (2001) with the one derived from the VALERI data set (<http://www.avignon.inra.fr>). This latter corresponds to an ensemble of measurements performed between 2000 and 2008 over various vegetation types, including the main biomes, and is used for the validation of remote sensing products (Figure 11). The results presented here correspond to the local measurements of effective LAI performed at the 20m scale. The VALERI distribution is significantly different from that of Scurlock's since true LAI values are most of the time higher than effective values. Further, the VALERI data set includes a significant fraction vegetation types with low LAI values such as crops shrubs and savannahs which also represent a significant fraction of the earth surface (about 30%). Therefore, it may better

correspond to what is expected for the algorithm training data base although the spatial and temporal sampling are very limited.



**Figure 11. Effective LAI distribution over the VALERI sites**

Therefore, the only way to access the actual distribution of LAI over the Earth Surface would be to use remote sensing estimates over a representative sample of the Earth surface (Baret et al. 2006; Masson et al. 2003). Such data set is currently available only at 1km spatial resolution (MODIS, CYCLOPES, GLOBCARBON) which might very significantly depart from the required distribution at 20 m resolution, particularly regarding the geostatistical characteristics of landscapes, with typical ranges of few hundreds of metres (Garrigues et al. 2006b).

In addition, the distribution used for LAI may also depend on that of the other variables, particularly for the larger LAI values corresponding to near 'saturation' conditions. In these conditions, the neural network will adjust a model that will pass in the middle of the cloud of possible solutions. To prevent underestimations due to more values with low LAIs, it was decided to use log-normal distributions that may agree quite well with those presented in Figure 10 and Figure 11 and fixing the upper limit to LAI=15. The log-normal distribution was slightly modified to increase the frequency of low LAI values by adding 15% of cases within a uniform distribution law for the 15% lower LAI values.

### 3.3.2.2 **Average Leaf Angle (ALA)**

Very few measurements of the average leaf angle measurements have been reported within the scientific community. The VALERI project provides some information on the expected distribution of LAI at 20 m spatial resolution. Figure 12 shows that low ALA values (up to 40°) are not very frequent. A peak is observed at ALA=60° which corresponds to a spherical distribution. The ALA was derived from the inversion of a gap fraction model measured either from hemispherical photos or LAI2000 devices allowing estimating concurrently both effective LAI (Figure 11) and ALA (Weiss et al. 2004). However, due to the limited spatial and temporal sampling associated, with in additional possible biases due to the retrieval technique, these distributions might be only considered as indicative.

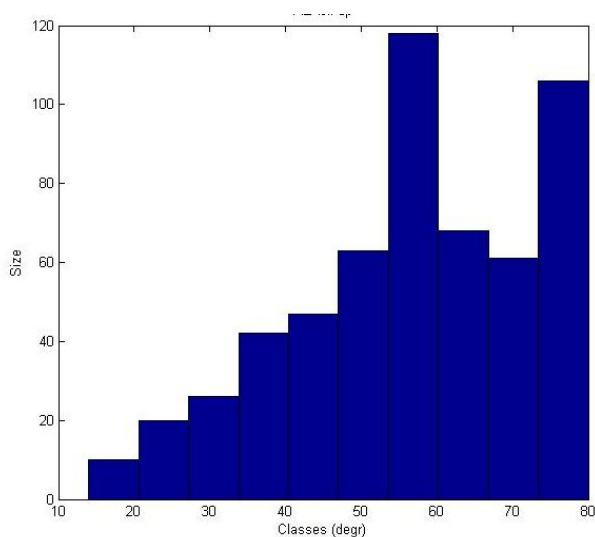


Figure 12. Distribution of the average leaf angle (ALA) over the VALERI sites

### 3.3.2.3 Hot Spot Parameter (HOT)

It is not easy to derive the hot spot parameter distribution at the Earth surface from the literature since it is a variable which is only used in radiative transfer modelling. Moreover, this variable has an impact only when remote sensing measurements are acquired near the principal plane. Figure 13 shows the distribution of typical values of the hot spot parameter for some crops provided by (Baret et al. 2009), most of the values are between 0.05 and 1. Extreme values are obtained for crops with specific architecture: low values for row planted trees (olive, peach) while the highest value is observed for young vineyards, lettuce and sugar beets.

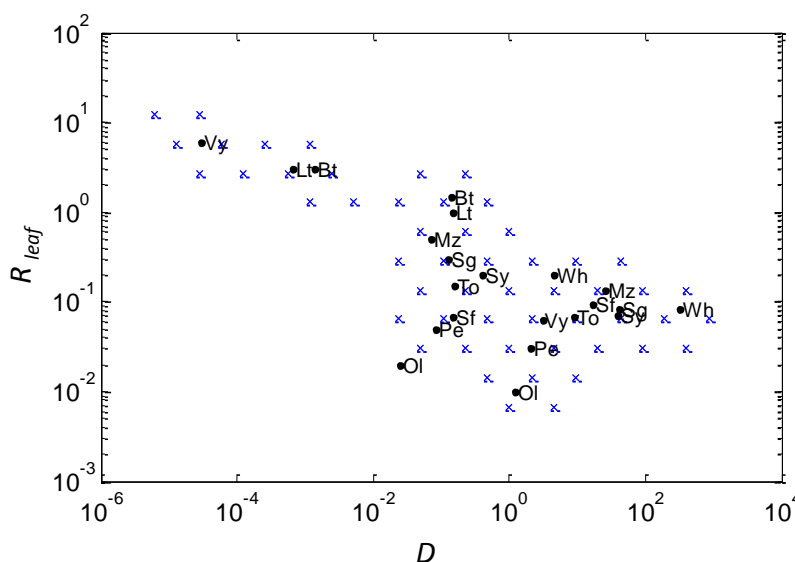


Figure 13. Relative leaf size (hot spot parameter) as a function of plant densities from (Baret et al. 2009). Typical values for 3D simulated scenes (crosses) and for actual canopies (circles): Wheat (Wh), Maize (Mz), Sorghum (Sg), Vineyards (Vy), Sunflower (Sf), Soybean (Sy), Tomatoe (To), Olive tree (Ol), Peach (Pe), Lettuce (Lt) and Beet (Bt).

### 3.3.2.4 Co-distribution between variables

Almost no information is available on possible co-distributions between variables. However, it is likely that some of the variables are linked together. For example, a very dense forest canopy will never be associated to low chlorophyll content and planophile leaf orientation. For this reason, we



proposed to restrict the range of variation for some variables as a function of LAI value. This will be simply achieved by assuming that the range of variation linearly changes with LAI between  $V_{\min}(0)$  (respectively  $V_{\max}(0)$ ) and  $V_{\min}(LAI_{\max})$  (respectively  $V_{\max}(LAI_{\max})$ ) as illustrated by Figure 14.  $LAI_{\max}$  is the maximum LAI value considered.

This will be simply achieved by assuming that the variable dynamics linearly varies with LAI according to:

$$(V - V_{\min}(0)) / (V_{\max}(0) - V_{\min}(0)) = (V^* - V_{\min}(LAI)) / (V_{\max}(LAI) - V_{\min}(LAI))$$

$$\text{With : } V_{\min}(LAI) = V_{\min}(0) + LAI * (V_{\min}(LAI_{\max}) - V_{\min}(0))$$

$$\text{and } V_{\max}(LAI) = V_{\max}(0) + LAI * (V_{\max}(LAI_{\max}) - V_{\max}(0))$$

where  $V^*$  is the value of variable  $V$  after linking its distribution to that of LAI. The values defining the co-distributions are specified in Table 6. They were derived empirically, assuming that large LAIs corresponded to a restricted range of the other variables.

Note that on other complementary way to relate the distribution between input variables will consist in filtering the simulated outputs (reflectance, FAPAR) as it will be demonstrated in the following.

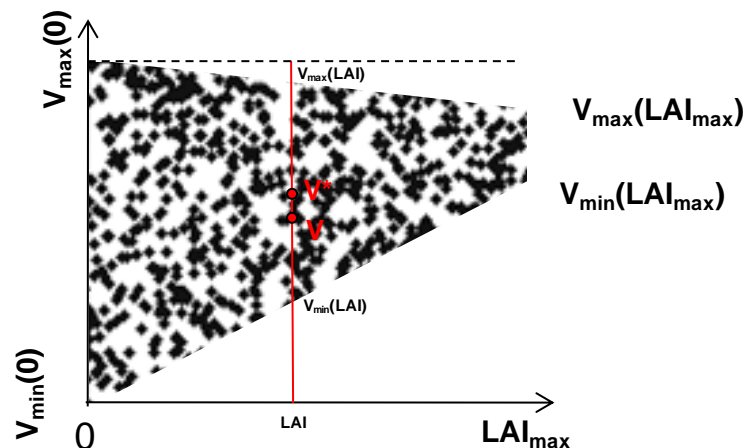


Figure 14. Scheme explaining the way the distribution of variable  $V$  is linked to that of LAI.

### 3.3.3 Actual generation of the inputs for the training data base

The training data base has to be sufficiently large to allow a robust calibration of the network, and also get a sub-set of the data base for hyper-specialization and test. The optimal size of the training data base depends on the complexity of the problem to solve. Previous studies (Combal et al. 2002) have shown that for a medium complexity problem, a training table close to 10 000 cases was satisfactory. In the current case that corresponds to a more complex algorithm, the size of the training data base should increase. We have simulated 41472 cases.



	Variable	Minimum	Maximum	Mode	Std	Nb_Class	Law
Canopy	LAI	0.0	15.0	2.0	3.0	6	Gauss
	ALA (°)	30	80	60	30	3	gauss
	Hot	0.1	0.5	0.2	0.5	1	gauss
Leaf	N	1.20	2.20	1.50	0.30	3	gauss
	Cab ( $\mu\text{g}\cdot\text{m}^{-2}$ )	20	90	45	30	4	gauss
	Cdm ( $\text{g}\cdot\text{m}^{-2}$ )	0.0030	0.0110	0.0050	0.0050	4	gauss
	Cw_Rel	0.60	0.85	0.75	0.08	4	uni
	Cbp	0.00	2.00	0.00	0.30	3	gauss
Soil	Bs	0.50	3.50	1.20	2.00	4	gauss

**Table 5. Distribution of the input variables of the radiative transfer model used to generate the training data base. Truncated Gaussian, log-normal or uniform distribution laws are used, characterized by the mode, the standard deviation (s), and minimum and maximum values. The number of classes for each variable is presented (Nb\_Class).**

	Variable	Co_Distribution	$V_{\min}(0)$	$V_{\max}(0)$	$V_{\min}(\text{LAI}_{\max})$	$V_{\max}(\text{LAI}_{\max})$
Canopy	ALA (°)	Yes	30	80	55	65
	Hot	Yes	0.1	0.5	0.1	0.5
Leaf	N	Yes	1.20	2.20	1.3	1.8
	Cab ( $\mu\text{g}\cdot\text{m}^{-2}$ )	Yes	20	90	45	90
	Cdm ( $\text{g}\cdot\text{m}^{-2}$ )	Yes	0.0030	0.0110	0.0050	0.0110
	Cw_Rel	Yes	0.60	0.85	0.70	0.80
	Cbp	Yes	0.00	2.00	0.00	0.20
Soil	Bs	Yes	0.50	3.50	0.50	1.20

**Table 6. Values used for the co-distributions of the variables with LAI (see Figure 14).**

The sampling scheme is based on a full orthogonal experimental plan (Bacour et al. 2002). This consists of identifying classes of values for each variable. Then all the combinations of classes are sampled once. Finally the actual values of each variable are randomly drawn within the range of variation defined by the corresponding class, according to the distribution law specified for the variable considered. This process allows accounting for all the interactions, while having the range of variation for each variable densely and near randomly populated. The number of classes (equally spaced) for each variable is shown in Table 5.

The following scheme was used to generate the training data base:

- Date.** Randomly draw a date. a whole year (365 days) is used. The year is decomposed into 4 periods that are successively sampled.
- Location.** The location is randomly selected within  $-56^\circ$  and  $+83^\circ$  latitudes corresponding to the mission requirements. However above  $+70^\circ$  latitude the illumination conditions are very poor most of the time, impacting the accuracy of radiance measurements. In addition, the high sun zenith angles experienced induce strong atmospheric effects due to the increase optical path.

Once date and location are sampled, the corresponding geometrical configuration is derived using SENTINEL2 orbit characteristics and swath (Table 1). Note that special patterns are observable, due to the orbit constraints.

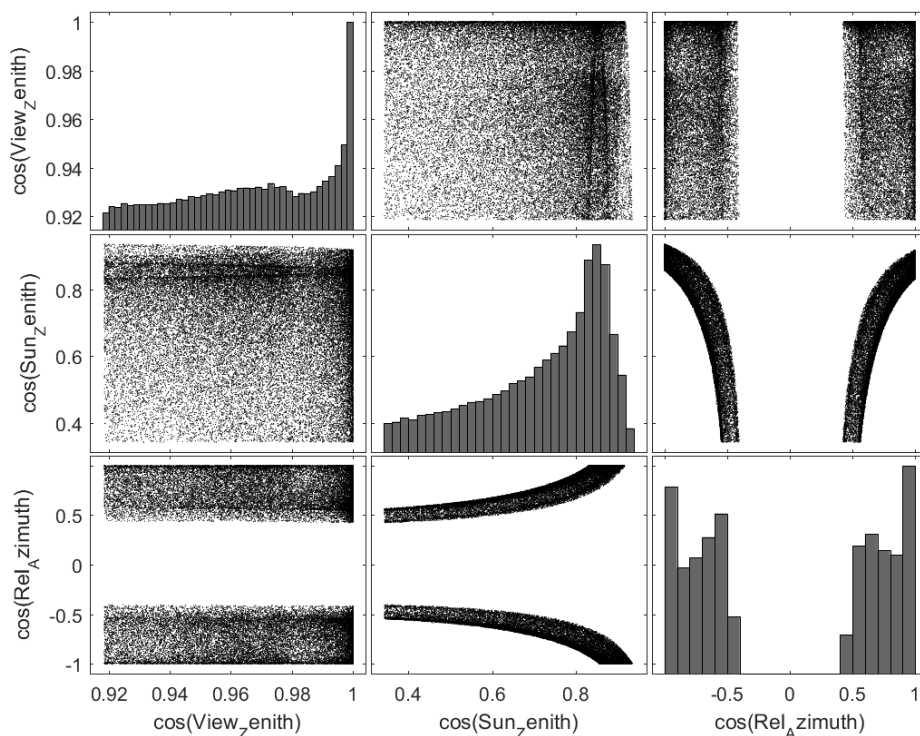


Figure 15: Distribution of values of observational characteristics used to simulate the learning data base.

### 3. Background variables.

- **Background type:** The background type is derived by randomly selecting one of the 7 soil background types (Figure 8).
- **Brightness:** The  $B_s$  coefficient is randomly drawn according to a truncated Gaussian distribution centred on  $B_s=1.0$  (Table 5). The larger frequencies for the lower  $B_s$  values are explained by the co-distribution with LAI values: bright soils are not expected under very dense vegetation.

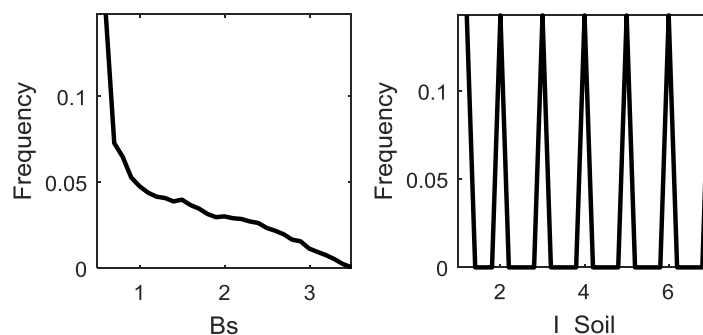


Figure 16. Distribution of the soil brightness ( $B_s$ , left) and soil type (right,  $I_{Soil}$ ) used.

### 4. Canopy variables. The following distributions were used (Table 5):

- **LAI:** The distribution of LAI values follows a gaussian distribution with a mode at LAI=2, and a relatively wide distribution ( $\sigma=3$ ). This allows sampling significantly low and high values of LAI. The higher LAI values were truncated at LAI=15. This particular selection of distribution of LAI values improved the saturation problem by introducing a significant amount of very high LAI values in the training process. The distribution law was further modified to increase the frequency of low LAI values as explained previously.
- **ALA:** The average leaf inclination angle is assumed to follow a truncated Gaussian distribution centred over the spherical widely represented one. The distribution is tied to the LAI, assuming that for large LAI values, leaf angle distribution was close to a spherical one. The peak observed between  $55^\circ$  and  $65^\circ$  is due to the co-distribution constraints for high LAI values.
- **HOT:** The hot spot parameter follows a truncated Gaussian distribution.

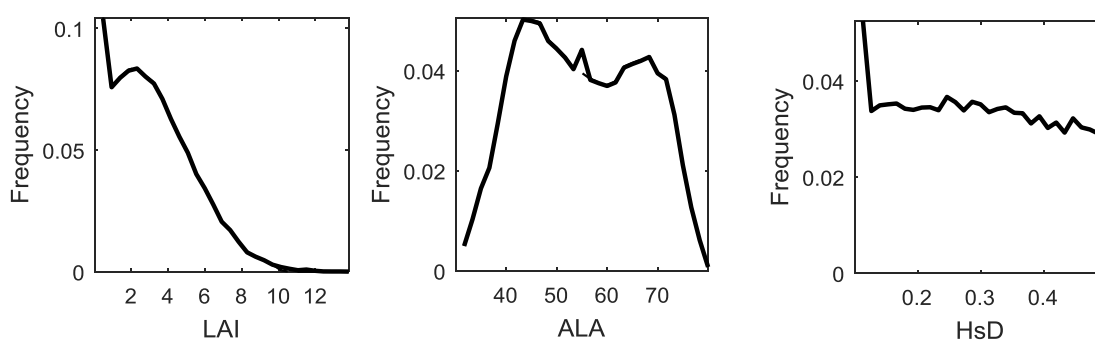


Figure 17. Distribution used for canopy variables LAI, ALA and Hot (HsD).

**5. Leaf optical properties.** Here also, very little knowledge is available on the actual leaf characteristics. Truncated Gaussian distributions were used for all these variables (Table 5) The artefacts observed in the distributions are due to the co-distribution constraints applied for high LAI values (Table 6).

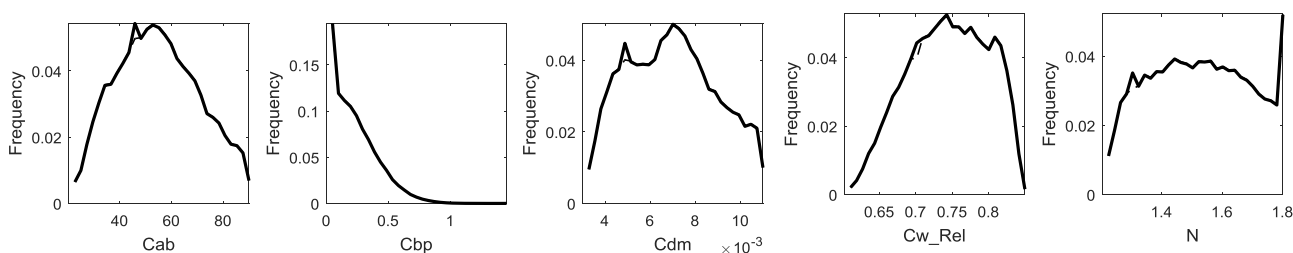


Figure 18. Distribution of the leaf characteristics used.

At the end of the simulation process, a total number of 41472 cases were simulated. Note that this number of cases is above the minimum 10 000 cases as required for a medium complexity problem and should allow good training performances for this more complex problem. This data set are then be split in two parts with a random selection process:

- **Training:** 2/3 of the simulations are affected randomly to the training of the neural network
- **Testing and Hyper-specialization:** 1/3 of the simulations are used for the hyper-specialization control and evaluation of theoretical performances.

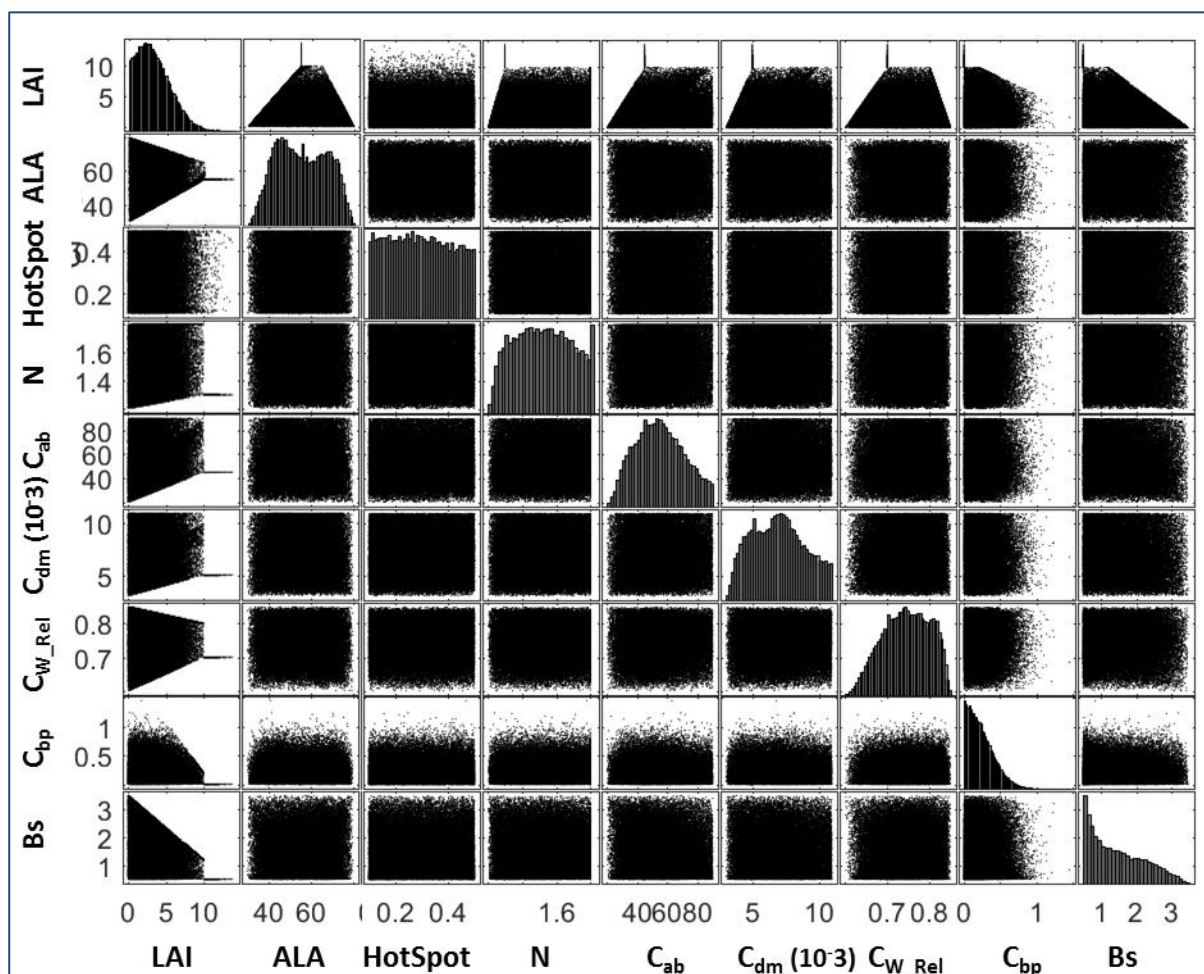


Figure 19. Distribution and co-distributions of the input soil, canopy and leaf characteristics used to populate the training data base.

## 6. Simulation of the top of canopy reflectance for the 8 SENTINEL2 bands

The previously derived table of input variables is used to simulate the corresponding SENTINEL2 top of canopy reflectance in the 8 bands using the PROSPECT+SAIL model.

A simple uncertainty model was used to better describe actual SENTINEL2 characteristics as well as the ability of the radiative transfer model used to represent actual reflectances. They are described in Table 7. The uncertainties are computed according to:

$$R^*(\lambda) = R(\lambda)(1 + (MD(\lambda) + MI)/100) + AD(\lambda) + AI$$

Where  $R(\lambda)$  is the raw simulated reflectance,  $R^*(\lambda)$  is the reflectance contaminated with noise, MD is the multiplicative wavelength dependant noise, MI is the multiplicative wavelength independent noise, AD is the additive wavelength independent noise, and AI is the additive wavelength independent noise.

The uncertainties attached to the radiative transfer model mainly derive from the representation of canopy architecture and leaf and soil background optical properties which is difficult to estimate. A posterior estimation will be issued using the reflectance mis-match criterion as computed over actual SENTINEL2 data.



Sensor	SENTINEL2			
	Add. Dep. (AD)	Add. Ind. (AI)	Mult. Dep. (MD) (%)	Mult. Ind. (MI) (%)
Band3	0.01	0.01	2	2
Band4	0.01		2	
Band5	0.01		2	
Band6	0.01		2	
Band7	0.01		2	
Band8a	0.01		2	
Band11	0.01		2	
Band12	0.01		2	

Table 7. Characteristics of the uncertainties model used.

Figure 20 shows the distribution and co-distribution of simulated reflectances. Bands 3 and 5 appear very strongly correlated. The same is observed for bands 6 and 7 with 8a, and in a lesser way bands. For each band, the distributions are roughly Gaussian.

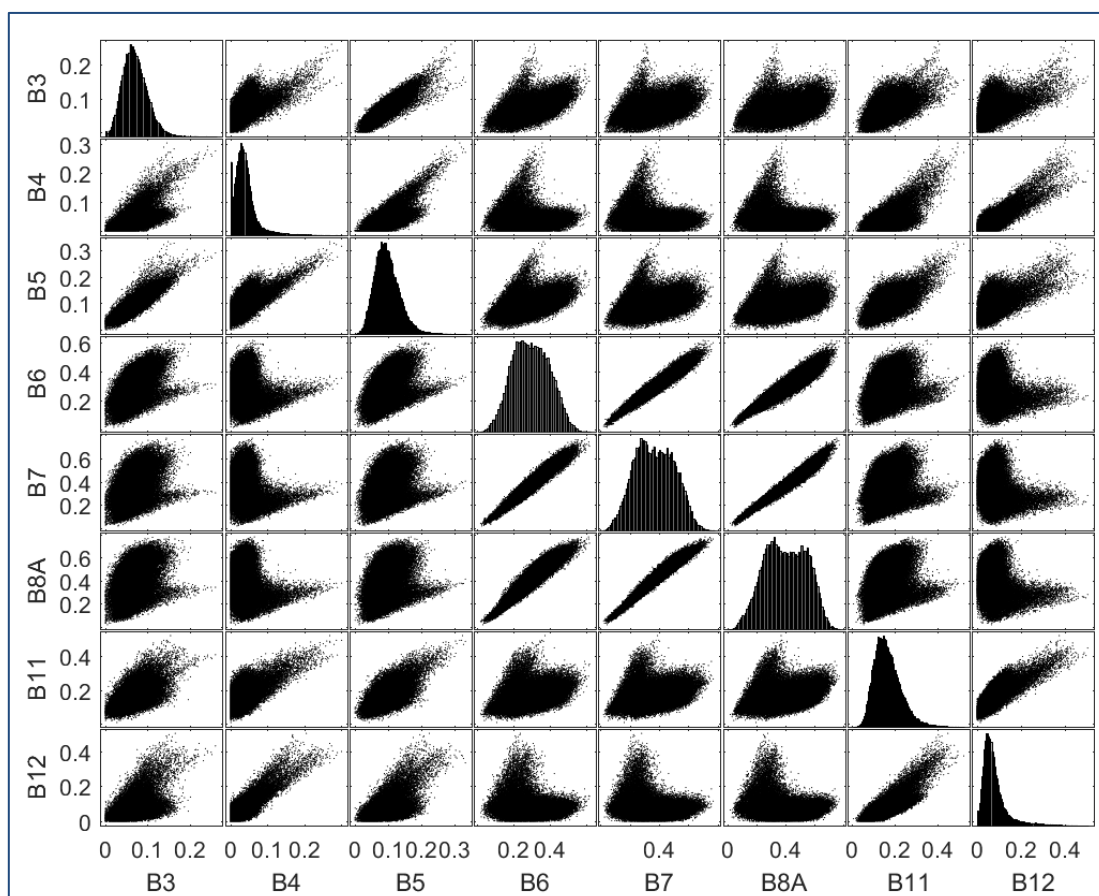


Figure 20: Distribution and co-distribution of the training database simulated reflectances, in the 8 SENTINEL2 bands.

### 3.3.4 Consequence on the distribution of the output variables

#### 3.3.4.1 Distribution of the output variables

Figure 21 shows the distribution and co-distribution of the output variables. They are highly non Gaussian. Co distributions are of interest since they demonstrate that there is no single relationship between one output variable and another. Note that, as expected, the CWC and CCC are quite linearly correlated to LAI. Further, some relationships have already been investigated in detail over a range of vegetation types such as the LAI-FAPAR relationship. It may be the basis to check the consistency of the simulated outputs with regards to the expected behaviour.

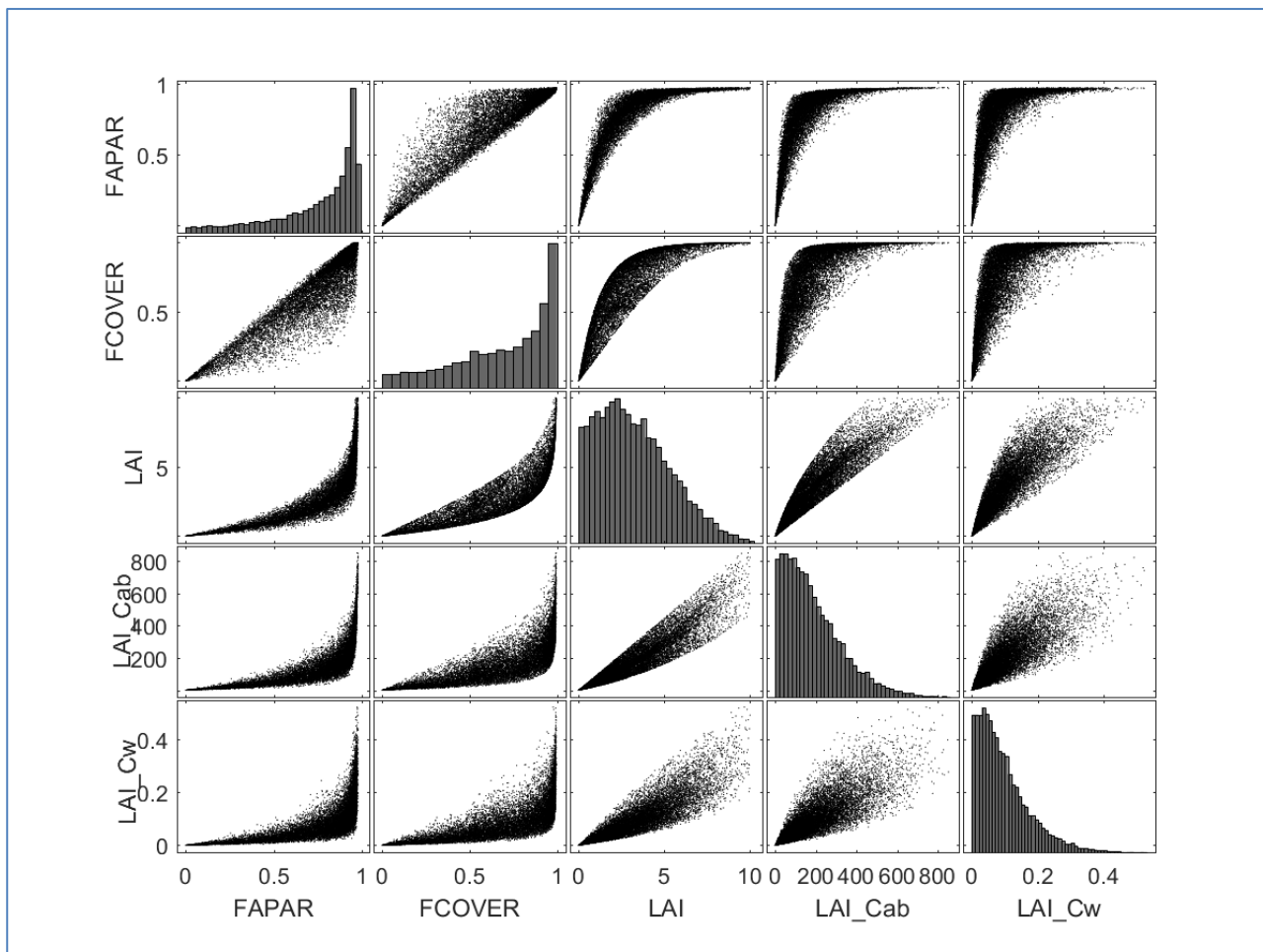
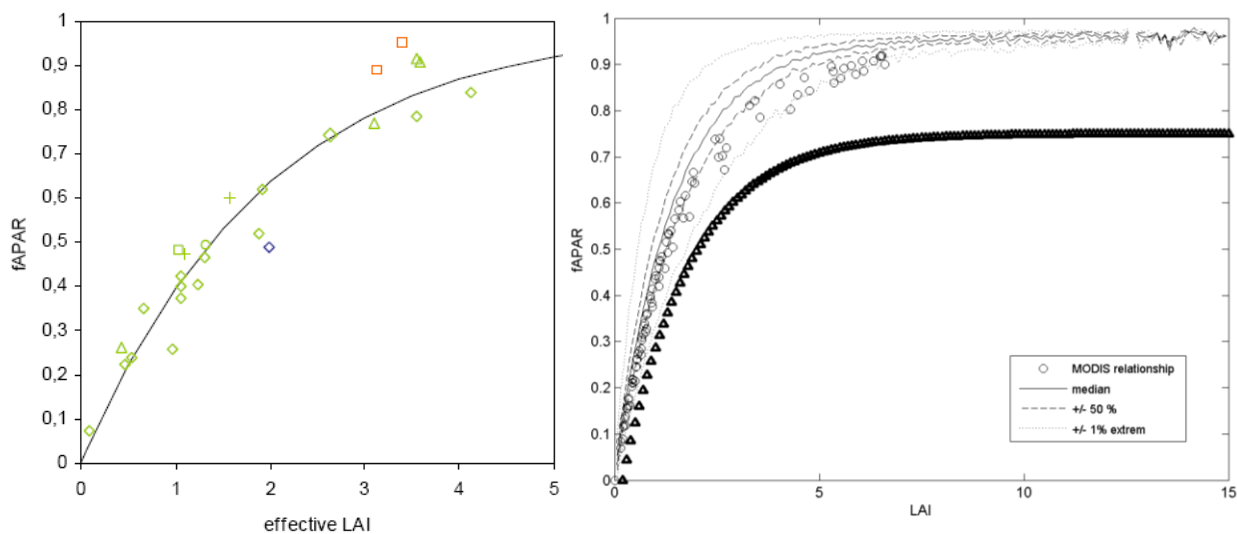


Figure 21. Distribution and co-distribution of the output variables.

#### 3.3.4.2 Relationship between LAI and FAPAR

A number of experimental observations report that the relationship between LAI and FAPAR follows an exponential law with extinction coefficients showing a limited range of variability. Figure 22 (left) shows results derived from the compilation of VALERI 3x3 km<sup>2</sup> sites covering a large range of biome types. Figure 22 right shows similar relationship as simulated for MODIS products using 3D radiative transfer model run for a range of canopy types. Both LAI-FAPAR relationships show very good agreement. The distribution of cases simulated for SENTINEL2 products shows also a good consistency.



**Figure 22:** On the left, relationship between LAI and FAPAR as established over the ensemble of VALERI sites. Symbols correspond to the sites. The solid line corresponds to the best fit model. On the right, symbols correspond to MODIS C4 relationship between LAI and FAPAR as simulated by the 3D radiative transfer model for the several vegetation types considered. The several lines correspond to the distribution within the learning data base. The line with triangles represents a threshold under which values are not expected.

### 3.4 TRAINING THE NEURAL NETWORK

Neural networks are defined mainly by the type of neurons used (the transfer function), the way they are organized and connected (the network architecture) and the learning rule. In addition, the input and output values need to be properly normalized to prevent any scaling factor or numerical problem. Back-propagation artificial neural network (Rummelhart et al. 1986) is one of the most common neural networks used to solve our radiative transfer model inversion problem.

#### 3.4.1 Normalization of the input and output values

The inputs (SENTINEL2 TOC reflectance in 8 bands and geometry) and output (the biophysical variable considered) values are first normalized according to Equation 6. Such data transformation is performed mainly to increase the performances of convergence of the training algorithm.

$$X^* = 2*(X-X_{Min})/(X_{Max}-X_{Min})-1 \quad \text{Equation 6}$$

Where  $X^*$  is the normalized input,  $X$  the original value,  $X_{min}$  and  $X_{max}$  respectively the minimum and maximum values.

#### 3.4.2 Network architecture

The connections between neurons are associated to a “synaptic” weight. Each neuron transforms the sum of the weighted signal from the previous neurons according to a given transfer function and a bias. The combination of sigmoid and linear functions is recognized as capable of fitting any type of function (Demuth and Beale 1998) .

For our more complex problem, an optimal architecture had to be determined for each biophysical variable. Several network structures have been tested. For each possible structure, three neural networks, differing by the initialization of their coefficients, have been trained. The selection of the "optimal" network architecture is then based on the RMSE between the outputs and the "true" biophysical variables as well as on the number of coefficients to be adjusted. Lower numbers are preferred because they allow faster runs of the neural networks in operational mode while precluding hyper-specialization.

Although it is possible to train a single network for all the variables considered here, previous tests have shown that the training is more easy when targeting a single variable. This decreases significantly the complexity of the neural network architecture associated to multiple output variables. Further, these tests showed also that only marginal lost consistency between the variables to be estimated was observed when using individual neural nets for each variable as compared to a single neural net for all the variables.

The neural networks investigated that way are thus composed of (Figure 14):

- one input layer made of the 11 normalized input data ( $\cos(\theta_s)$ ,  $\cos(\theta_v)$ ,  $\cos(\phi)$ , and the TOC reflectances in the 8 SENTINEL2 wavebands).
- one hidden layers with 5 neurons with tangent sigmoid transfer functions.
- one output layer with a linear transfer function.

Note that this simple network requires 65 synaptic coefficients and 6 bias coefficients to adjust.

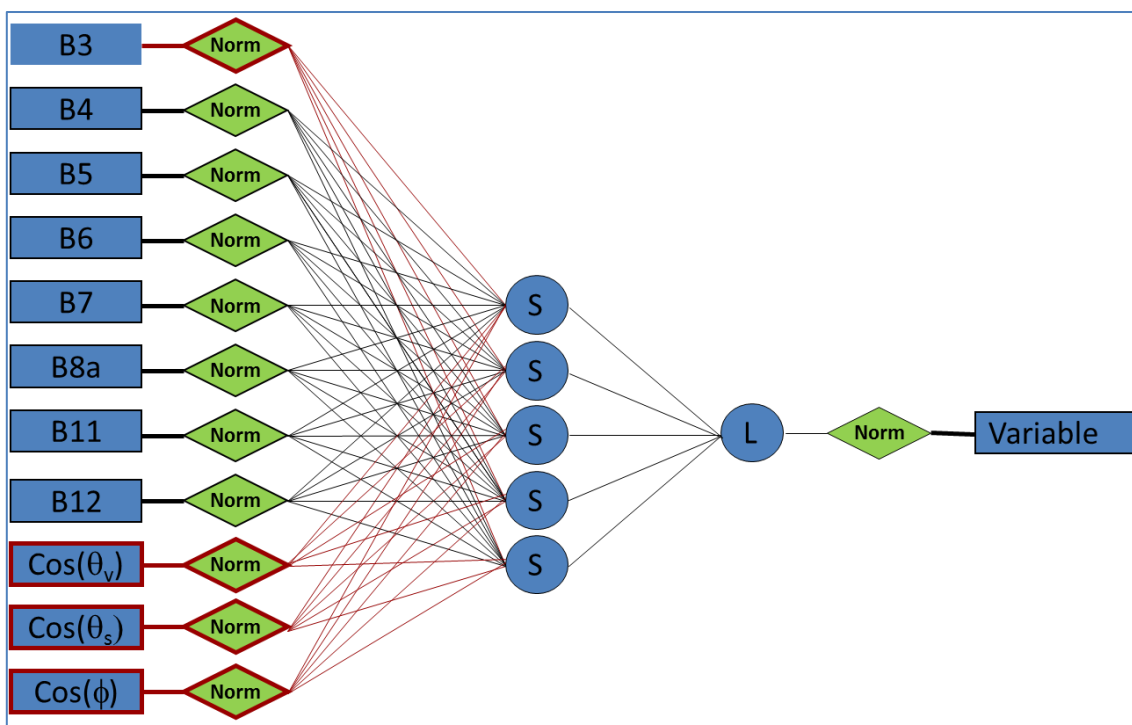


Figure 23. Neural network architecture developed for the estimation of the biophysical variables considered from the 9 SENTINEL2 bands and the 3 angles defining the geometry of observation. The network is made of 1 hidden layer of 5 neurons and 1 linear output neuron. The 'Norm' symbols correspond to the normalization process as described by Equation 10. Symbols 'S' and 'L' correspond respectively to the sigmoid (tansig) and linear transfer functions of the neurons.



### 3.4.3 Learning process

The learning process is mainly made of two elements: the training dataset that was described earlier, and the learning rule that is now described. The Levenberg–Marquardt optimization algorithm is used to adjust the synaptic weights and neuron bias to get the best agreement between the output simulated by the network and the corresponding value of canopy biophysical variable simulated in the training data base. The initial values of the weights and biases were set to a random value between -1.0 and +1.0. To prevent from hyper-specialization, a sub-set of the training data base is used to control whether the network starts to hyper-specialize, i.e. represents the particular features of the training data set and therefore losing its capacity to describe the general features to be extracted. When this happens, the optimization process is stopped.

Three networks were trained in parallel to retrieve the canopy biophysical variables, each corresponding to independent random drawing of the initial values of the synaptic weights and bias. Finally, the network kept for implementation uses is the one that provides the best performances over the test data set.

### 3.4.4 Theoretical Performances

Theoretical performances allow checking whether no major problems occur during the training process and provide a first glance on the capacity of SENTINEL2 to access the considered variables. The theoretical performances of the networks were evaluated over the test data set which is a fraction of the simulated training data base (1/3) not used in the calibration of the neural networks. It consists mainly in computing simple statistics such as RMSE values and exploring the dependency of residuals on the variable itself.

Results show that the training of the networks was quite efficient, with relatively small RMSE values: 0.89 for LAI, 0.05 for FAPAR, 0.04 for FVC, 56 $\mu\text{g}/\text{cm}^2$  for CCC and 0.03g/cm<sup>2</sup> for CWC. FAPAR (Figure 25) and FVC (**Erreur ! Source du renvoi introuvable.**) show the best performances as expected, with the larger RMSE values observed for the medium values of the products. The algorithms are unbiased as expected (Figure 24 to Figure 26). No early saturation effect is observed as a function of the product value: LAI seems to be well estimated up to values of LAI=6. However, uncertainties increase with LAI values (Figure 24). The same applies for CCC and CWC since they are directly correlated with LAI (Figure 27 and Figure 28).

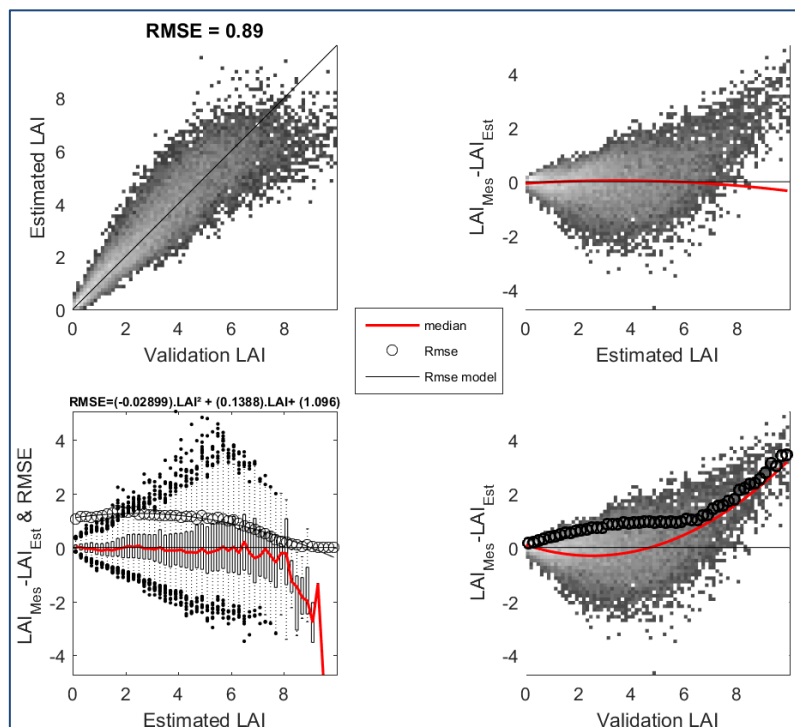


Figure 24. Theoretical performances of the neural networks for LAI on the test database

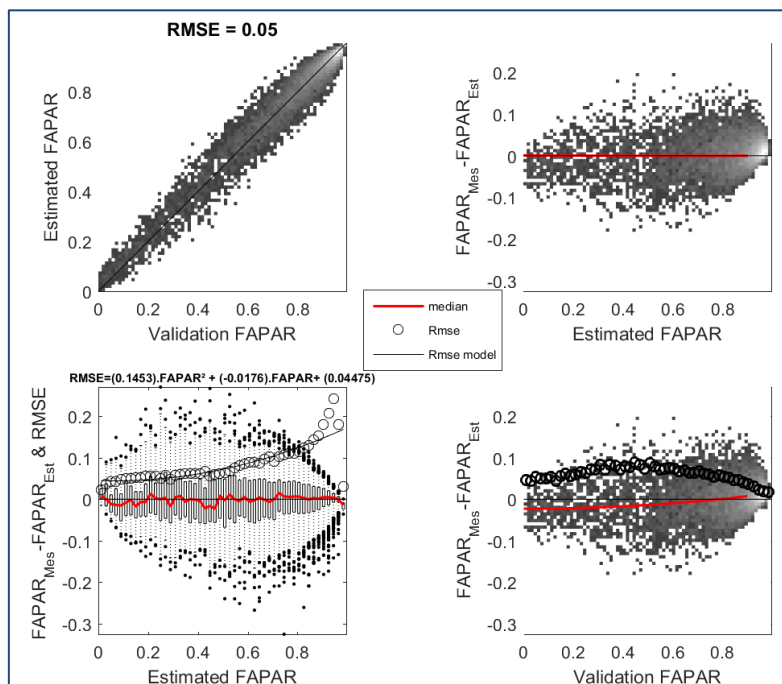


Figure 25. Theoretical performances of the neural networks for FAPAR on the test database

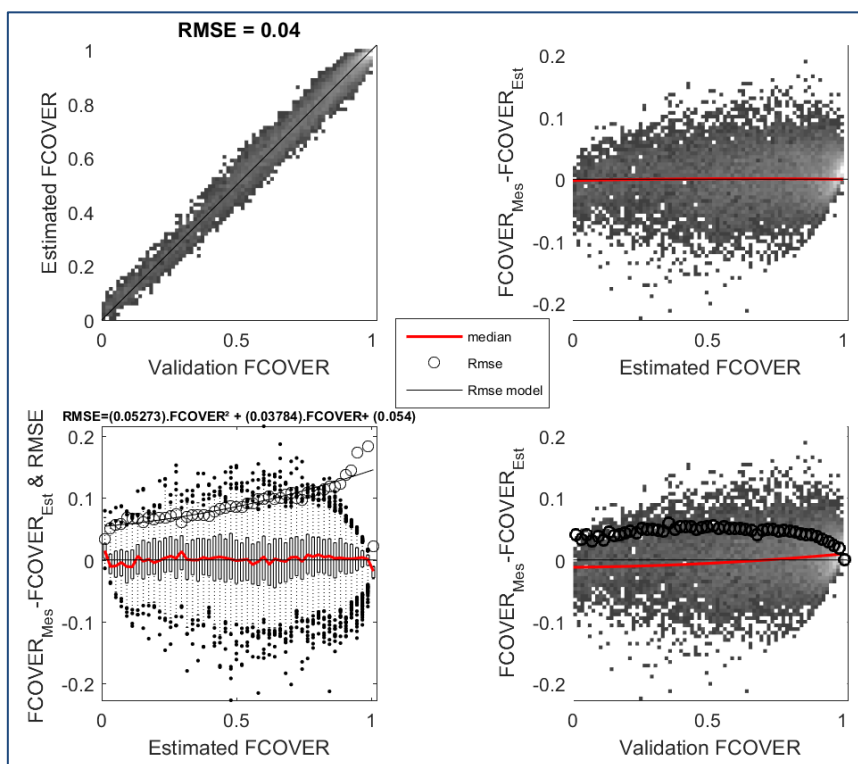


Figure 26. Theoretical performances of the neural networks for FCOVER on the test database

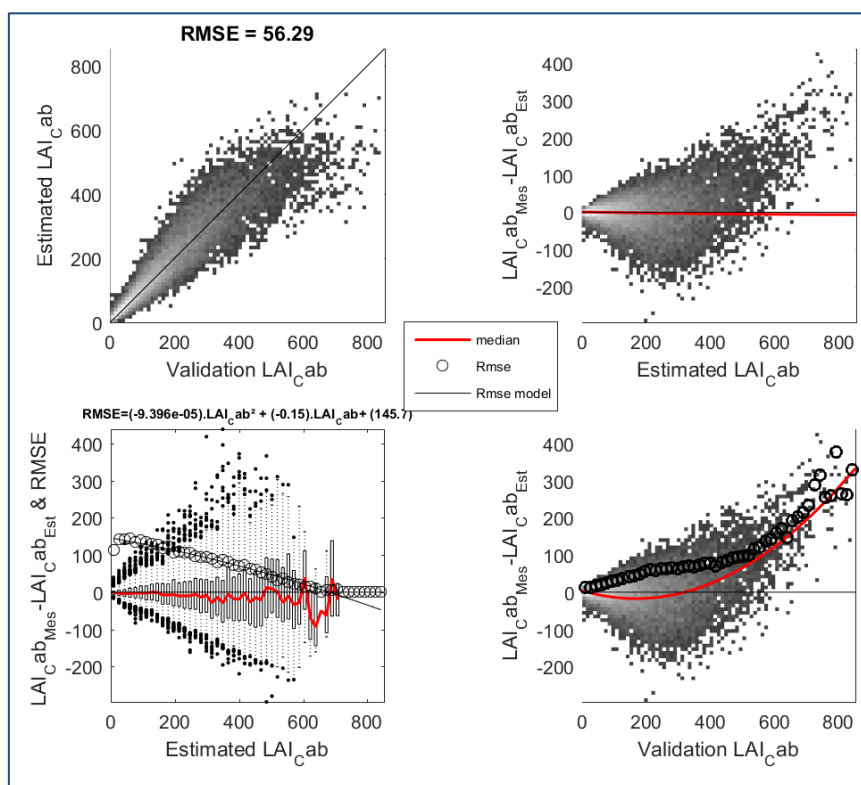


Figure 27. Theoretical performances of the neural networks for CCC on the test database

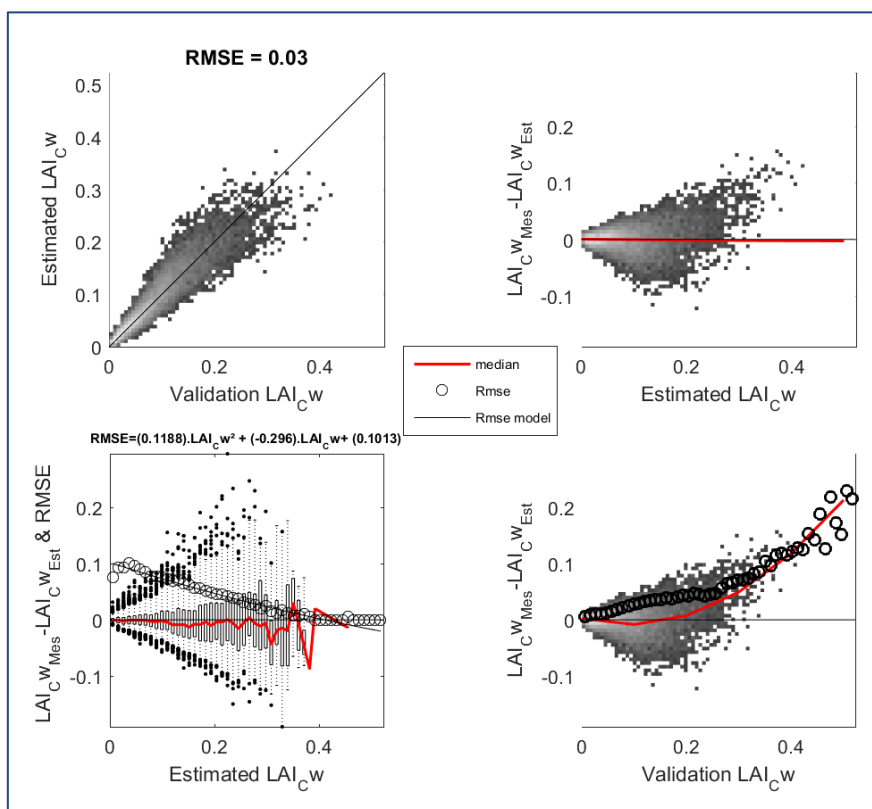


Figure 28. Theoretical performances of the neural networks for CWC on the test database

## 3.5 DEFINITION DOMAIN

### 3.5.1 Input out of range flag

When inputs are outside the convex hull defined by the simulated reflectance values of the training data base, i.e. the definition domain, then a specific 'input out of range' flag is raised. The convex hull is approximated by a hypercube with the same dimensions as those of the inputs of the neural network (i.e. 8 for SENTINEL2). Each dimension corresponding to a specific input of the neural network varies within the minimum and maximum values (Table 8). The range is split into 10 equal classes of values (see Figure 29 in the simple case of 2 inputs). However, because it was checked that geometrical configuration was more or less evenly distributed with regards to reflectance values in the considered bands, the geometrical configuration was not entered as input for the definition domain to simplify the process and make the code running faster. Results show that less than 50% of the cells the 8<sup>th</sup> dimension hypercube corresponds to expected reflectance observations over the surface.

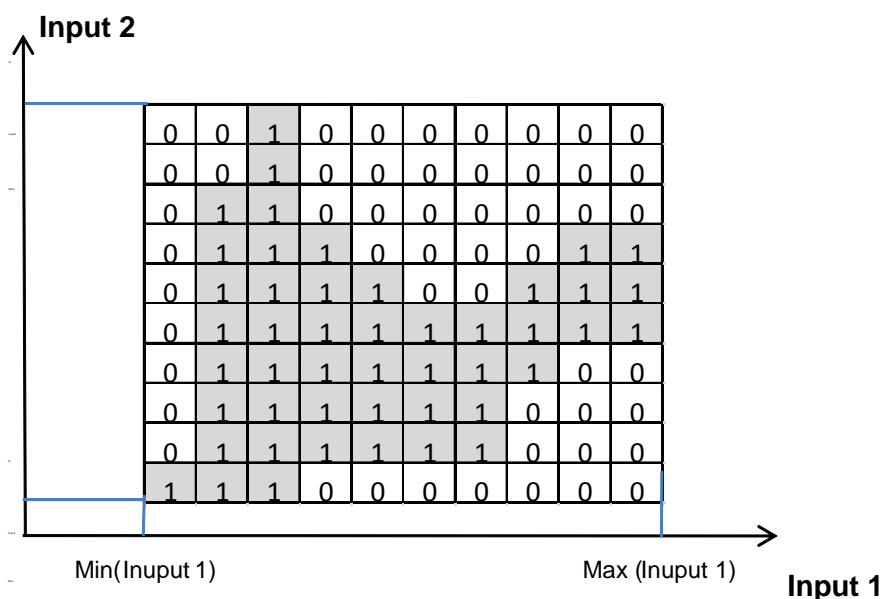


Figure 29. Schematic representation of the convex hull in the case of 2 dimensional inputs of the network. The convex hull is approximated by a regular grid. The cells with ones (the gray cells) correspond to cases of possible input range (i.e. inside the convex hull) while cells with zeros correspond to inputs out of range (outside the convex hull).

	Min	Max
Band3	0	0.26
Band4	0	0.30
Band5	0	0.32
Band6	0	0.62
Band7	0	0.75
Band8a	0	0.78
Band11	0	0.52
Band12	0	0.50

Table 8. Minimum and maximum input values (bounding box) of the definition domain for inputs from SENTINEL2 TOC reflectance.

The algorithm used to check whether the inputs are outside the convex hull is the following:

- Check if the inputs are outside the bounding box, i.e. if the value of an input is either lower than the corresponding minimum or larger than the maximum values defined in Table 8. If this is the case, the 'input out of range' flag is raised.

- If inputs are within the bounding box, then, check if values of the inputs are in a non valid cell, i.e. a cell with '0' (Figure 29). If this is the case, the 'input out of range' flag is raised.

### 3.5.2 Output out of range flag.

In the case where the ANN provides biophysical variable estimates outside their definition range, the corresponding product value will be set to the closest bound of the range, i.e. either the minimum or the maximum accepted values. However, because of the several sources of uncertainties associated to the inputs and the algorithm calibration, a tolerance is set before raising the 'output out of range' flag: Three cases are possible for each product:

1. The product value is within expected physical range of variation (Table 9): It is considered valid.
2. The value is within the tolerance limits (Table 9) but higher (lower) than the physical maximum (minimum). The value is considered valid but set to the physical maximum (minimum).
3. The value is outside the tolerance limits: it is considered invalid.

	Tolerance	$P_{\min}$	$P_{\max}$
LAI	0.2	0	8.0
FAPAR	0.1	0	0.94
FVC	0.1	0	1.0
CCC	15	0	600
CWC	0.015	0	0.55

**Table 9. Tolerance, and Minimum ( $P_{\min}$ ) and maximum ( $P_{\max}$ ) values admitted for the products.**

## 4 PRACTICAL CONSIDERATIONS FOR THE ALGORITHM IMPLEMENTATION

### 4.1 REQUIREMENTS FOR EASILY UPDATING THE ALGORITHM

To be able to easily upgrade the algorithm in the processing chain, all the coefficients used for the neural networks, normalization, quality flags and uncertainties should be set in parameter tables that can just be changed when updates will be available. This will occur possibly several times within a short period before launching the operational processing chain, when new product values will be available and when possible problems in the early versions will be quickly identified.

This information is available in excel files for each product: 'Algo\_S2\_VX\_SL2T\_VVV.xlsx', where X satnds for the version of the algorithm, and VVV stands for the vegetation variable investigated, i.e. VVV=[LAI, FAPAR, FVC].

### 4.2 ALGORITHM IMPLEMENTATION

The different steps are described hereafter:

1. **Normalization of the inputs:** for the 12 inputs X, the following normalization equation must be applied:

$$X^* = 2 \cdot (X - X_{\min}) / (X_{\max} - X_{\min}) - 1$$

where  $X^*$  is normalized input value, and  $X_{\min}$  and  $X_{\max}$  are computed over the neural network training data set. These values are provided in the Excel file under sheet 'Normalisation'

2. **Run the neural network**

The neural networks will be described by its architecture, *i.e.*, the the number of hidden layers and the output layer. Each layer is described by its number of neurons, associated weight and biases and transfer function.

For the neurons of the hidden layers, the transfer function is a tangent sigmoid function given by:

$$Y = \text{Tansig}(x) = 2 / (1 + \exp(-2x)) - 1$$

While for the output layer the transfer function is linear ( $y=x$ ).

For each neural net (one per product and per sensor), tables will be provided given the weight and biases for each neurons. Example for the NNT provided for LAI is provided hereafter. Sheet 'Weights' in the Excel file contains the weights, biases and neural network structure information.

3. **Denormalization of the output**

It simply consists in applying the inverse function used for input normalization:

$$Y = 0.5 \cdot (Y^* + 1) \cdot (Y_{\max} - Y_{\min}) + Y_{\min}$$



where  $Y^*$  is normalized output value issued from the NNT, and  $Y_{min}$  and  $Y_{max}$  are computed over the neural network training data set. These values are in the sheet 'Normalisation' in the Excel file contains the 'denormalisation' information.

#### 4. Generate quality indicator:

Generate the quality indicator QA coded over 3 bits,

QA=0 0 0 : data is OK

QA=0 0 1 : input out of range (provide product value if within tolerance)

QA=0 ,1 0 : output out of range (provide product value if within tolerance)

QA= 0 1 1 : input and output out of range (product value = fill value)

QA= 1 0 0 : bad quality of input values (depending on the TOC reflectance S2 product (provide product value if within tolerance)

## 5 CONCLUSION

This ATBD describes the algorithm used to compute *LAI*, *FAPAR* and *FVC*, from SENTINEL2 top of canopy reflectance data that will be implemented in the SENTINEL2 Toolbox. If the principles of the algorithm will be kept unchanged, it is likely that updates of the coefficients will be made when issues will be identified, data on actual models of uncertainties refined (including residual atmospheric effects). It is therefore important to keep the algorithm easy to upgrade. Refinement could be also made with actual SENTINEL2 top of canopy reflectance, allowing to filter the training data base so that only the simulated cases too different from the observed ones would be eliminated. This should result in better training and performances.

The proposed algorithm is based on specific radiative transfer models associated with strong assumptions, particularly regarding canopy architecture (turbid medium model). All the variables derived from such algorithms should be seen as effective, i.e. the variables that would correspond to the measured satellite signal reflected by a canopy verifying all the assumptions made through the radiative transfer models. Depending on the variable, this may lead to differences with ground values that may be accessed from field measurements. Further, the algorithm is 'generic', i.e. it should apply to any type of vegetation with reasonable performances. However, to better match the specificities of given canopies, either simple correction could be calibrated, or more specific algorithm could be developed.

One strong assumption embedded in any single pixel retrieval algorithm as this one, is that the pixel targeted belongs to a landscape patch presenting enough homogeneity (at the pixel scale) preventing unexpected loss or gain of radiation fluxes. Therefore, it can be applied for larger resolution than 20m. For forests with large crowns, or any pixel showing strong heterogeneity such as pixels at the intersection between two different vegetation patches, results may be uncertain. This extends also to pixels where the neighbouring ones are very different. Specific algorithms should be developed to detect such situations and possibly propose alternative retrieval methods.

## 6 REFERENCES

- Asner, G.P., Wessman, C.A., Schimel, D.S., & Archer, S. (1998). Variability in leaf and litter optical properties: implications for BRDF model inversions using AVHRR, MODIS, and MISR. *Remote Sensing of Environment*, 62:243-257
- Bacour, C., Baret, F., Béal, D., Weiss, M., & Pavageau, K. (2006). Neural network estimation of LAI, fAPAR, fCover and LAIxCab, from top of canopy MERIS reflectance data: principles and validation. *Remote Sensing of Environment*, 105, 313-325
- Bacour, C., Jacquemoud, S., Tourbier, Y., Dechambre, M., & Frangi, J.P. (2002). Design and Analysis of numerical experiments to compare four canopy reflectance models. *Remote Sensing of Environment*, 79, 72-83
- Baret, F., & Buis, S. (2007). Estimating canopy characteristics from remote sensing observations. Review of methods and associated problems. In S. Liang (Ed.), *Advances in Land Remote Sensing: System, Modeling, Inversion and Application* (pp. 171-200): Springer
- Baret, F., De Solan, B., & Weiss, M. (2009). PAI estimates from digital photos at 57.5° zenith angle over wheat crops. *Agricultural and Forest Meteorology*, soumis (Février 2009)
- Baret, F., Jacquemoud, S., & Hanocq, J.F. (1993). The soil line concept in remote sensing. *Remote Sensing Reviews*, 7, 65-82
- Baret, F., Koetz, B., & Bruguier, N. (2002). WP 1400: Characterization of maize structure and optical properties. In (p. 6). Avignon (France): INRA-CSE
- Baret, F., Leroy, M., Roujean, J.L., Knorr, W., Lambin, E., & Linderman, M. (2003). CYCLOPES User Requirement Document. In. Avignon: INRA-CSE
- Baret, F., Morisette, J., Fernandes, R., Champeaux, J.L., Myneni, R., Chen, J., Plummer, S., Weiss, M., Bacour, C., Garrigue, S., & Nickeson, J. (2006). Evaluation of the representativeness of networks of sites for the global validation and inter-comparison of land biophysical products. Proposition of the CEOS-BELMANIP. *IEEE Transactions on Geoscience and Remote Sensing*, 44, 1794-1803
- Baret, F., Weiss, M., Leroy, M., Hautecoeur, O., Santer, R., & Begue, A. (1997). Impact of surface anisotropies on the observation of optical imaging sensors. In: INRA Bioclimatologie
- Biard, F., & Baret, F. (1997). Crop residue estimates using multiband reflectance data. *Remote Sensing of Environment*, 59, 530-536
- Borel, C.C., Gerstl, S.A.W., & Powers, B.J. (1991). The radiosity method in optical remote sensing of structured 3-D surfaces. *Remote Sens. Environ.*, 36, 13-44
- Campbell, G.S. (1986). Extinction coefficients for radiation in plant canopies calculated using an ellipsoidal inclination angle distribution. *Agric. For. Meteorol.*, 36, 317-321
- Chelle, M., Andrieu, B., & Bouatouch, K. (1997). Nested radiosity for plant canopies. *The Visual Computer*, Nov1997-Mar1998, 1-24
- Chen, J.M., Menges, C.H., & Leblanc, S.G. (2005). Global mapping of foliage clumping index using multi-angular satellite data. *Remote Sensing of Environment*, 97, 447-457
- Chen, Y., & McKyes, E. (1993). Reflectance of light from the soil surface in relation to tillage practices, crop residues and the growth of corn. *Soil & Tillage Research*, 26, 99-114
- Combal, B., Baret, F., Weiss, M., Trubuil, A., Macé, D., Pragnère, A., Myneni, R., Knyazikhin, Y., & Wang, L. (2002). Retrieval of canopy biophysical variables from bi-directional reflectance data. Using prior information to solve the ill-posed inverse problem. *Remote Sensing of Environment*, 84, 1-15

- Demuth, H., & Beale, M. (1998). Neural network user's guide. In T. Mathworks (Ed.), *Matlab User's guide*
- España, M., Baret, F., Aries, F., Chelle, M., Andrieu, B., & Prévot, L. (1999). Modeling maize canopy 3D architecture. Application to reflectance simulation. *Ecological modeling*, 122, 25-43
- Fourty, T., & Baret, F. (1997). Amélioration de la précision des coefficients d'absorption spécifique de la matière sèche et des pigments photosynthétiques. In (p. 35). Avignon: INRA Bioclimatologie
- Garrigues, S., Allard, D., Baret, F., & Weiss, M. (2006a). Influence landscape spatial heterogeneity on the non-linear estimation of leaf area index from moderate spatial resolution remote sensing data. *Remote Sensing of Environment*, 105, 286-298
- Garrigues, S., Allard, D., Baret, F., & Weiss, M. (2006b). Quantifying Spatial Heterogeneity at the Landscape Scale Using Variogram Models. *Remote Sensing of Environment*, 103, 81-96
- Garrigues, S., Shabanov, N.V., Swanson, K., Morisette, J.T., Baret, F., & Myneni, R.B. (2008). Intercomparison and sensitivity analysis of Leaf Area Index retrievals from LAI-2000, AccuPAR, and digital hemispherical photography over croplands. *Agricultural and Forest Meteorology*, 148, 1193-1209
- Gascon, F., & Berger, M. (2007). Sentinel-2 Mission Requirements Document In ESA (Ed.) (p. 30): ESA
- Gastellu-Etchegorry, J.P., Demarez, V., Pinel, V., & Zagolski, F. (1996). Modeling radiative transfer in heterogeneous 3-D vegetation canopies. *Remote Sensing of Environment*, 58, 131-156
- Gausman, H.W., Gerbermann, A.H., Wiegand, C.L., Leamer, R.W., Rodriguez, R.R., & Noriega, J.R. (1975). Reflectance differences between crop residues and bare soils. *Soil science society of America proceedings.*, 39, 752-755
- Gerstl, S.A.W., & Borel, C.C. (1992). Principles of the radiosity method versus radiative transfer for canopy reflectance modeling. *IEEE Transactions on Geoscience and Remote Sensing*, 30, 271-275
- Govaerts, Y.M., & Verstraete, M.M. (1998). Raytran: a monte carlo ray tracing model to compute light scattering in three dimensional heterogeneous media. *IEEE Transactions on Geoscience and Remote Sensing*, 36, 493-505
- Green, D.S., Erickson, J.E., & Kruger, E. (2003). Foliar morphology and canopy nitrogen as predictors of light use efficiency in terrestrial vegetation. *Agricultural and Forest Meteorology*, 115, 163-171
- Houlès, V., Mary, B., Machet, J.M., Guérif, M., & Moulin, S. (2001). Do crop characteristics available from remote sensing allow to determine crop nitrogen status? In G. Grenier, & S. Blackmore (Eds.), *3rd European Conference on Precision Agriculture* (pp. 917-922). Montpellier: Agro Montpellier
- Jacquemoud, S., & Baret, F. (1990). PROSPECT: A model of leaf optical properties spectra. *Remote Sensing of Environment*, 34, 75-91
- Jacquemoud, S., Baret, F., & Hanocq, J.F. (1992). Modeling spectral and directional soil reflectance. *Remote sensing of the Environment*, 41, 123-132
- Kimes, D.S., Gastellu-Etchegorry, J.P., & Esteve, P. (2002). Recovery of forest canopy characteristics through inversion of complex 3D model. *Remote Sensing of Environment*, 79, 320-328
- Kötz, B., Baret, F., Poilvé, H., & Hill, J. (2005). Use of coupled canopy structure dynamic and radiative transfer models to estimate biophysical canopy characteristics. *Remote Sensing of Environment*, 95, 115-124
- Kötz, B., Kneubuehler, M., Huber, S., Schopfer, J., & Baret, F. (2007). Radiative transfer model inversion based on multi-temporal CHRIS/PROBA data for LAI estimation. In, *ENVISAT symposium*. Montreux: ESA

- Kuusik, A. (1991). The hot-spot effect in plant canopy reflectance. In (pp. 1-16): Tartu University
- Lauvernet, C., Baret, F., Hascoët, L., Buis, S., & Le Dimet, F.X. (2008). Multitemporal-patch ensemble inversion of coupled surface-atmosphere radiative transfer models for land surface characterization. *Remote Sensing of Environment*, 112, 851-861
- Le Maire, G. (2002). Utilisation de la télédétection hyperspectrale pour la détermination des caractéristiques biophysiques et biochimiques des couverts végétaux: de l'échelle de la feuille à l'échelle du couvert. In (p. 45). Orsay (France): Laboratoire d'Ecologie Végétale, université Paris-Sud
- Liu, W., Baret, F., Gu, X., Tong, Q., Zheng, L., & Zhang, B. (2002). Relating soil surface moisture to reflectance. *Remote Sensing of Environment*, 81, 238-246
- Masson, V., Champeaux, J.L., Chauvin, F., Meriguer, C., & Lacaze, R. (2003). A global database of land surface parameters at 1km resolution in meteorological and climate models. *Journal of Climate*, 16, 1261-1282
- Myneni, R.B., Asrar, G., & Hall, F.G. (1992). A Three-dimensional radiative transfer method for optical remote sensing of vegetated land surfaces. *Remote Sensing of Environment*, 41, 105-121
- Newnham, G.J., & Burt, T. (2001). Validation of leaf reflectance and transmittance model for three agricultural crop species. *IEEE Transactions on Geoscience and Remote Sensing*, 2976-2978
- Prince, S.D. (1991). A model of regional primary production for use with coarse resolution satellite data. *International Journal of Remote Sensing*
- Rummelhart, D.E., Hinton, G.E., & Williams, R.J. (1986). Learning internal representations by error propagation. In D. Rummelhart, & J. Mc Clelland (Eds.), *Parallel data processing* (pp. 318-362). Cambridge, MA (USA): M.I.T. press
- Sanz, C., Espana, M., Baret, F., Weiss, M., Vaillant, L., Hanocq, J.F., Sarrouy, C., Clastre, P., Bruguier, N., Chelle, M., Andrieu, B., & Zurfluh, O. (1997). Bi-directional characteristics of leaf reflectance and transmittance: measurement and influence on canopy bi-directional reflectance. In G. Guyot, & T. Phulpin (Eds.), *7th International Symposium on physical measurements and signatures in remote sensing* (pp. 583-590). Courchevel (France): Balkema
- Scurlock, J.M.O., Asner, G.P., & Gower, S.T. (2001). Worldwide Historical Estimates and Bibliography of Leaf Area Index, 1932-2000. In. Oak Ridge, Tennessee, U.S.A.: Oak Ridge National Laboratory
- Soler, C., F., S., Blaise, F., & de Reffye, P. (2001). A physiological plant growth simulation engine based on accurate radiant energy transfer. In (p. 31). Montbonnot-Saint-Martin (France): INRIA
- Verhoef, W. (1984). Light scattering by leaf layers with application to canopy reflectance modeling: the SAIL model. *Remote Sensing of Environment*, 16, 125-141
- Verhoef, W. (1985). Earth observation modeling based on layer scattering matrices. *Remote Sensing of Environment*, 17, 165-178
- Weiss, M., Baret, F., Leroy, M., Hautecoeur, O., Bacour, C., Prévot, L., & Bruguier, N. (2002). Validation of neural net techniques to estimate canopy biophysical variables from remote sensing data. *Agronomie*, 22, 547-554
- Weiss, M., Baret, F., Myneni, R., Pragnère, A., & Knyazikhin, Y. (2000). Investigation of a model inversion technique for the estimation of crop characteristics from spectral and directional reflectance data. *Agronomie*, 20, 3-22
- Weiss, M., Baret, F., Smith, G.J., Jonckheere, I., & Coppin, P. (2004). Review of methods for in situ leaf area index determination, part II: Estimation of LAI, errors and sampling. *Agricultural and Forest Meteorology*, 121, 37-53
- Zarco-Tejada, P.J., Miller, J.R., Harton, J., Hu Baoxin, N., T.L., Goel, N., Mohammed, G.H., & Sampson, P. (2001). Needle chlorophyll content estimation through model inversion using hyperspectral data from boreal conifer forest canopies. *Remote Sensing of Environment*, submitted

