# Integrating APSIM and PROSAIL to improve prediction of crop traits in various situations from hyperspectral data using deep learning

**Chen Q**[1,2*], Bangyou Zheng[2], Tong Chen[3] and Scott Chapman[1*]

[1] School of Agriculture and Food Sciences, The University of Queensland, St Lucia, QLD 4067, Email: qiaomin.chen@uq.edu.au, scott.chapman@uq.edu.au

[2] CSIRO Agriculture and Food, Queensland Biosciences Precinct 306 Carmody Road, St Lucia, QLD 4067, Email: bangyou.zheng@csiro.au

[3] School of Information Technology and Electrical Engineering, The University of Queensland, St Lucia, QLD 4067, Email: tong.chen@uq.edu.au

## Abstract

A major challenge of high-throughput phenotyping is to build the relationship between data-derived traits and empirically measured traits that usually are biologically related. More general relationships unrestricted to specific situations can be inferred from synthetic dataset generated by radiative transfer models (RTMs). However, this approach can lead to an "ill-posed" problem, resulting in unsatisfactory inversion results for target traits retrieval. This research investigated a practical way to introduce biological constraints in 'synthetic' training data by integrating a crop growth model with an RTM to alleviate this problem. Our result shows this integration made a much more accurate estimation for target traits. Furthermore, this simulation framework allows us to determine the precision of prediction from reflectance of traits in different environments and situations. These findings can benefit the application of high-throughput phenotyping in precision agriculture and plant breeding.

## Keywords

model integration, neural network, hyperspectral data, variable retrieval

## Introduction

Imagery methods have been used to monitor and investigate vegetation since the 1960s. Recently, these methods have been deployed in more proximal sensors (planes, drones, vehicles) that allow the analysis of vegetation at higher resolutions (sub-centimetre scale) in a research field that is sometimes referred to as 'high-throughput phenotyping' (HTP) (e.g. Chapman et al. 2018).  HTP methods based on sensor and imaging technologies can rapidly measure a large number of crop traits across time and space in a cost- and labour-efficient way, which can benefit applications in precision agriculture and plant breeding. A major challenge of high-throughput phenotyping is to build the relationship between data-derived traits and empirically measured traits that usually are biologically related. For example, some vegetation indices (e.g. NDVI, EVI, etc.) computed from canopy reflectance had been developed to be used to predict LAI (e.g., Dong et al. 2019).There is an increasing interest of the application of 'model inversion methods' to radiative transfer models (RTMs) (e.g., Berger et al. 2018). These methods provide an easier way to develop more general relationships unrestricted to situations for variable retrieval by using RTMs to generate a training dataset that represents the entire range of possible situations varying in crop types and growth status as well as observation configurations (e.g., Baret and Buis 2008; Dorigo et al. 2007).

Although model inversion methods provide a reasonable way for estimating crop or vegetation variables from remote sensing data, none of them can avoid the "ill-posed" problem.  However, this problem can be alleviated by using prior knowledge to strengthen constraints on individual variables or between variables. Linking a crop growth model (CGM) to RTM provides a more straightforward solution to address this "ill-posed" problems by directly constraining the sets of RTM input parameters that contribute to canopy reflectance in two ways. The first way is to calibrate/sparameterise CGM using canopy reflectance and then using the calibrated CGM to predict target crop traits.  Such an application mode can directly retrieve those variables not included in RTMs, such as crop yield, and provide dynamic estimation across the whole growth season (e.g., Thorp et al. 2012; Zhang et al. 2016). The other way is to convert CGM output variables into input variables of RTM and then apply the model inversion method on these constrained input variables and corresponding canopy reflectance. Unfortunately, this approach has been rarely discussed or explored.

This research focused on the estimation of total leaf area index (LAI, $m^2$ $m^{-2}$), leaf chlorophyll content (Cab, µg $cm^{-2}$), dry leaf weight (Cm, g $cm^{-2}$) and leaf water content (Cw, g $cm^{-2}$) of wheat in four locations across the Australian Wheatbelt. The overall objective was to investigate inversion procedures based on a deep

learning approach (feedforward neural network, FFNN) for crop trait estimation, with a special focus on alleviating the "ill-posed" problem in model inversion through linking a CGM (APSIM) and a RTM (PROSAIL) to generate a higher quality training dataset.
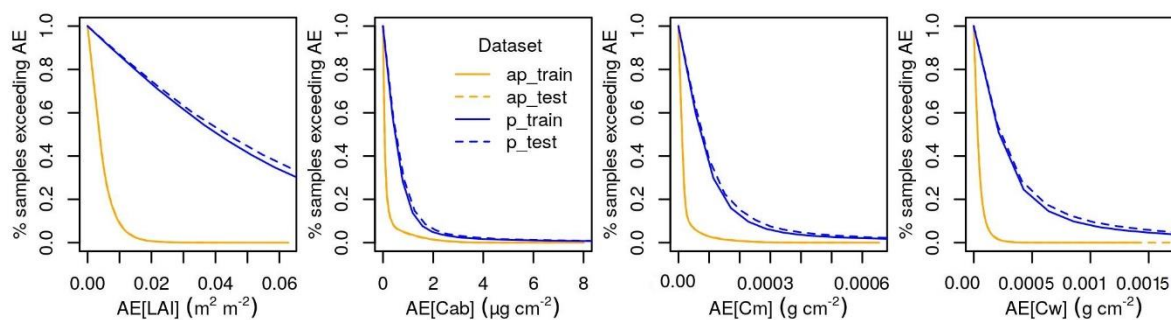
## Methods

*Couple APSIM and PROSAIL:* The Agricultural Production System Simulator (APSIM) Next Generation (https://www.apsim.info/apsim-next-generation/) is the new version of APSIM, which is simpler and faster than the classic version (i.e. 7.10., D. Holzworth et al. 2018). The application of APSIM to simulate the dynamics of many crop traits has been validated in many regions worldwide. PROSAIL combines PROSPECT and SAIL, which can simulate directional canopy reflectance (Jacquemoud et al. 2009). The current version of the PROSAIL (PROSAIL_D) can be downloaded from http://teledetection.ipgp.jussieu.fr/prosail/. The coupling of APSIM and PROSAIL is srealised by passing output variables of APSIM to PROSAIL as input variables. This permits the coupling model to estimate canopy reflectance from 400 to 2500 nm in 1 nm interval at defined observation conditions (determined by latitude, day of the year, and daytime) given that required parameters are specified. The transformation of variables is based on a series of equations and more details referred to Chen et al (2021).

*Generate synthetic dataset:* A defined set of conditions for wheat growth (scharacterised by genotype, environment and management) and observation (determined by local latitude, day of year and daytime) were set up to run APSIM and PROSAIL for simulation of crop traits and canopy reflectance, which resulted in two types of synthetic datasets. The first dataset (p_data) uses the ranges of the input parameters converted from APSIM outputs but allows PROSAIL to be run using samples from full parameter space for any combination of inputs. The second dataset (ap_data) directly uses input data converted from APSIM outputs to explore a sub-space of input parameters (i.e., limited by the APSIM biology) to run PROSAIL. More details about the generation of these two datasets can be found in our previous work (Chen et al. 2021). In total, p_data contain 100 000 unique samples (90 000 in p_train, 10 000 in p_test) while ap_data contain 2 149 226 unique samples (90 000 in ap_train, 10 000 in ap_test, 2 049 226 in ap_rest).

*Build, train and evaluate FFNN:* The FFNN model was implemented using Keras in TensorFlow 2.3.0 (https://www.tensorflow.org/). Two simulation experiments were designed to evaluate the effect of limiting the PROSAIL input parameters to the sub-space as determined by the APSIM. The FFNN model in two simulation experiments shared the same model structure (hyperparameters): 3 hidden layers and 512 units for each hidden layer and using 0.001 as learning rate, 'softplus' as activation, 'Adamax' as optimiser. More details about the determination of FFNN's optimal structure could be found in Chen et al (2021). The input layer included 1101 dimensions of 1-nm reflectance bands from 400-1100 nm and the output layer included four target variables (i.e. Cab, Cm, Cw, LAI). The first simulation experiment used p_train as a training set. It evaluated the trained model on p_test, while the second simulation experiment used ap_train as a training set and evaluated the trained model on ap_test and ap_rest.
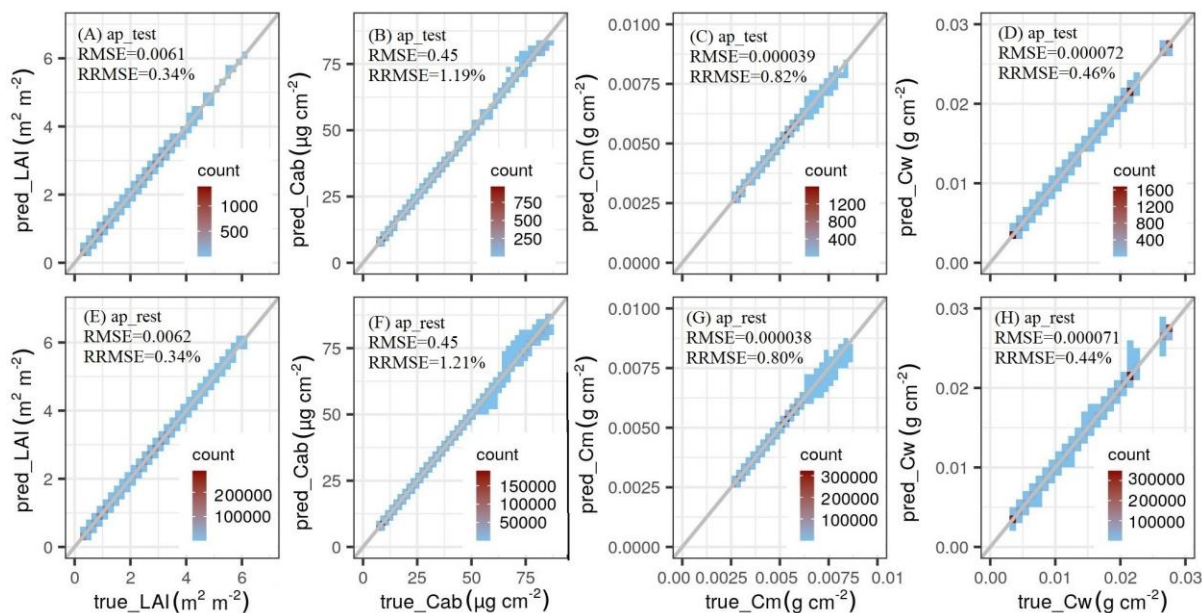
## Results

Compared with the performance of the FFNN model trained on p_data, Figure 1 indicates that the use of ap_data generated by coupling APSIM and PROSAIL significantly improved FFNN's performance for estimation of all target variables using reflectance bands in 400-1100 nm. This improved model performance presented in both average precision (smaller MAE and RMSE) and stability (narrower uncertainty range, smaller standard deviation of AE) according to statistical results. For instance, compared with test results of the trained model using p_data, the use of ap_data for total LAI estimation narrowed the uncertainty range of AE from 0~1.093 $m^2$ $m^{-2}$ to 0~0.041 $m^2$ $m^{-2}$ and also reduced the standard deviation of AE (from 0.063 to 0.004 $m^2$ $m^{-2}$), MAE (from 0.061 to 0.005 $m^2$ $m^{-2}$) and RMSE (from 0.087 to 0.006 $m^2$ $m^{-2}$ ). Compared with simulated results from other model inversion studies (Atzberger 2004; le Maire et al. 2008; Upreti et al. 2019), our results from p_data have ~10 times smaller RMSE for estimation of LAI/Cab and ~3 times smaller RMSE for Cm/Cw due to the use of better architecture and algorithm used in a neural network as well as complete information included in massive hyperspectral bands. Furthermore, our results from ap_data had even higher precision than results from p_data due to the advantages mentioned above plus the biological constraints imposed on the co-distribution of PROSAIL input variables.
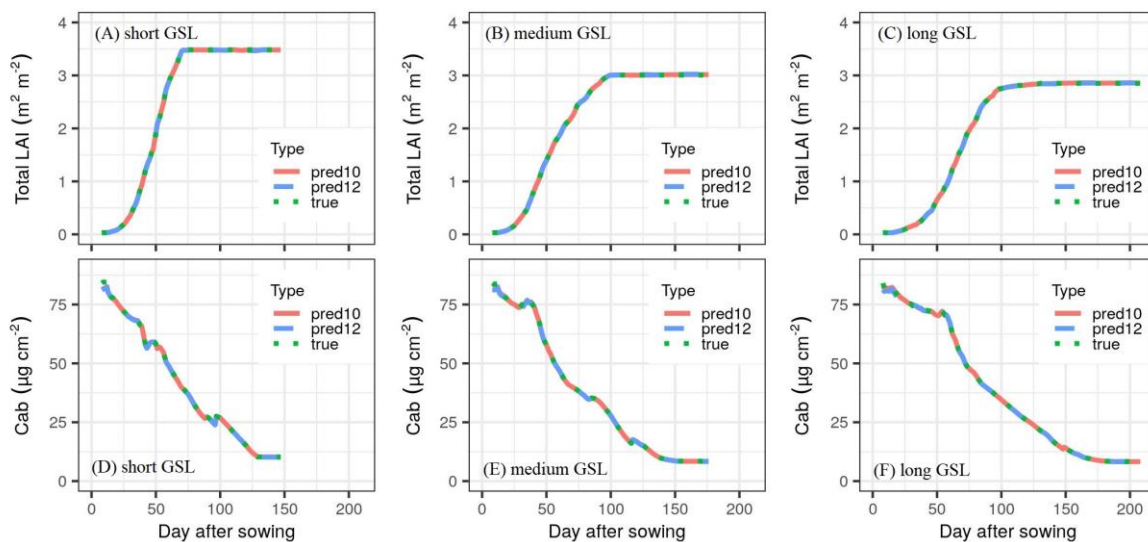
**Figure 1 Empirical accumulated density distribution for absolute error (AE) of prediction on different datasets. AE is the difference between the actual value of the target variable and its prediction.**

The fully trained FFNN model trained on ap_train were used to estimate target variables on ap_test and ap_rest. Figure 2 shows that the model can perform well on broad unseen samples from the same situation, resulting in low RMSE and RRMSE (RRMSE is the ratio of RMSE divided by the mean of true values) of all target variables. Additionally, the model continued to make good estimations at different levels of the true value, although larger true values tended to have a higher probability of larger absolute errors.



**Figure 2 True value (generated by APSIM) of target variables and their predictions on two test datasets (A~D for ap_test and E~H for ap_rest). RRMSE is the ratio of RMSE divided by the mean of true values.**

The model trained on ap_train performed similarly well in various growing situations varying in genotypes, sites, years and sowing dates (omitted for space). These slight differences in box plot of prediction error across conditions were likely resulted from the varying density distributions of true values of target variable in corresponding situations as prediction accuracy for target variables at different true value were slightly different (see Figure 2). At the seasonal scale, the seasonal prediction error of target variables (seasonal RMSE or seasonal RRMSE) was not associated with the growth pattern characterised by growing season length (GSL) (omitted for space). Figure 3 shows true and predicted values of total LAI and Cab during growing season (from emergence to maturity) with short, medium and long GSL. During the growing season, higher prediction error for total LAI usually appeared at the end of the growing season where true value were larger (Figure 3A~C), while higher prediction error for Cab, Cm and Cw tended to occur at the beginning of the growing season (Figure 3D~F for Cab and results for Cm and Cw are not shown due to limited space).

**Figure 3** True and predicted value for total LAI (A~C) and Cab (D~F) during growing season with varying growing season length (GSL). 'pred10' and 'pred12' indicate the predicted value inverted from reflectance captured at 10:00 and 12:00, respectively. 'true' indicates the simulated value generated by APSIM.

## Conclusion

This research demonstrated that the integration of APSIM and PROSAIL could generate higher quality training data, which can better characterise the real canopy srealisation and lead to more accurate estimation for target variables in theory. Although this trained FFNN model might not perform as well as presented here when it is applied to retrieve variables from real observation data due to measurement and model uncertainties, it is expected to be able to make relative good performance according to the difference of estimation precision on simulated and observed data from other model inversion studies.

## References

Atzberger C (2004). Object-based retrieval of biophysical canopy variables using artificial neural nets and radiative transfer models. Remote Sensing of Environment 93, 53–67.

Berger K, Atzberger C, Danner M, et al. (2018). Evaluation of the PROSAIL model capabilities for future hyperspectral model environments: A review study. Remote Sensing 10, 85.

Buis S and Baret F (2008). Estimating canopy characteristics from remote sensing observations: Review of methods and associated problems. S. Liang (ed.). Advances in Land Remote Sensing, 173–201.

Chapman SC, Zheng B, Andries B, et al. (2018). Visible, Near Infrared, and Thermal Spectral Radiance On-Board UAVs for High- Throughput Phenotyping of Plant Breeding Trials, in: Biophysical and Biochemical Characterization and Plant Species Studies. pp. 275–299.

Chen Q, Zheng B, Chen T, et al. (2021). Integration of APSIM and PROSAIL models to develop more precise radiometric estimation of crop traits using deep learning. bioRxiv: 2021.02.02.429471. (doi: https://doi.org/10.1101/2021.02.02.429471)

Dong T, Liu J, Shang J, et al. (2019). Assessment of red-edge vegetation indices for crop leaf area index estimation. Remote Sensing of Environment 222, 133-143.

Dorigo WA, Zurita-Milla R, de Wit AJW, et al. (2007). A review on reflective remote sensing and data assimilation techniques for enhanced agroecosystem modeling. International Journal of Applied Earth Observation and Geoinformation 9, 165–193.

Holzworth D, Huth NI, Fainges J, et al. (2018). APSIM Next Generation: Overcoming challenges in modernising a farming systems model. Environmental Modelling and Software 103, 43–51.

Jacquemoud S, Verhoef W, Baret F, et al. (2009). PROSPECT + SAIL models: A review of use for vegetation characterization. Remote Sensing of Environment 113, S56–S66.

le Maire G, François C, Soudani K, et al. (2008). Calibration and validation of hyperspectral indices for the estimation of broadleaved forest leaf chlorophyll content, leaf mass per area, leaf area index and leaf canopy biomass. Remote Sensing of Environment 112, 3846–3864.

Thorp KR, Wang G, West AL, et al. (2012). Estimating crop biophysical properties from remote sensing data by inverting linked radiative transfer and ecophysiological models. Remote Sensing of Environment 124, 224–233.

Upreti D, Huang W, Kong W, et al. (2019). A comparison of hybrid machine learning algorithms for the retrieval of wheat biophysical variables from sentinel-2. Remote Sensing 11, 481.