# IQPC 2015 TRACK: WATER DETECTION AND CLASSIFICATION ON MULTI-SOURCE REMOTE SENSING AND TERRAIN DATA

A. Olasz[a], D. Kristóf[a]\*, M. Belényesi[a], K. Bakos [a],
Z. Kovács[b], B. Balázs[b], Sz. Szabó[b]

[a] FÖMI, Institute of Geodesy, Cartography and Remote Sensing, 1149 Budapest, Hungary - kristof.daniel@fomi.hu
[b] Department of Physical Geography and Geoinformation Systems, University of Debrecen, Hungary -
szabo.szilard@science.unideb.hu

**Commission III, WG III/5**

**KEY WORDS:** IQmulus Processing Contest, Remote Sensing, Water detection and classification, Multi-source remote sensing data processing, Resource optimization, R programming

**ABSTRACT:**

Since 2013, the EU FP7 research project "IQmulus" encourages the participation of the whole scientific community as well as specific user groups in the IQmulus Processing Contest (IQPC). This year, IQPC 2015 consists of three processing tasks (tracks), from which "Water detection and classification on multi-source remote sensing and terrain data" is introduced in the present paper. This processing track addresses a particular problem in the field of big data processing and management with the objective of simulating a realistic remote sensing application scenario. The main focus is on the detection of water surfaces (natural waters, flood, inland excess water, other water-affected categories) using remotely sensed data. Multiple independent data sources are available and different tools could be used for data processing and evaluation. The main challenge is to identify the right combination of data and methods to solve the problem in the most efficient way. Although the first deadline for submitting track solutions has passed and the track has been successfully concluded, the track organizers decided to keep the possibility of result submission open to enable collecting a variety of approaches and solutions for this interesting problem.

## 1. INTRODUCTION

### 1.1 Introduction

The IQmulus EU FP7 project (http://www.iqmulus.eu) encourages the participation of the whole scientific community as well as specific user groups and research teams in the framework of IQmulus Processing Contest (IQPC) since 2013. Software performance is evaluated through the creation of benchmarks and evaluation methodologies specific for selected processing tasks (IQmulus Tracks). Beside the selection of test datasets with a ground truth, IQPC supports the usage of a common infrastructure where the executables submitted are run and results are collected and evaluated.

In 2015, IQPC is a theme of a Special Session in the ISPRS Geospatial Week in the GeoBigData Workshop and track reports will be reviewed to be included in the conference proceedings (ISPRS Archives). At the IQPC Special Session 30 minutes presentations will be held by the track organizers to report about the proposed solutions about each submission as well.

This paper introduces the track named "Water detection and classification on multi-source remote sensing and terrain data". This track has been defined by the Institute of Geodesy, Cartography and Remote Sensing, Hungary (FÖMI), a consortium partner in the IQmulus project. The challenge is to address the problem of detecting water surfaces and different categories of water-affected soils and vegetation based on a variety of remotely sensed data of different resolutions.

### 1.2 Background

Detection and monitoring of various water surfaces has been a challenge for a long time in remote sensing data processing. A large number of studies are available in the relevant literature dealing with water and wetness detection and monitoring for land management and conservation (e.g. Li et al. 2013; Rokni et al. 2014). Numerous different types of RS data are useful for some kind of water detection; however the accuracy is highly dependent on the input source, the processing methodology and in particular the combination of the two.

In real-life situations, a balance is achieved by creating a processing chain consisting of different methods and input sources and spatially aware algorithms are used to combine the tools to provide a result of sufficient information content. However a very challenging task is to optimize the resources for the tasks.

### 1.3 Objectives

The goal of this track is to detect water-related thematic classes in a specific area in Hungary. The track leader provides a wide array of remote sensing data covering the area in concern, including:

1. High-resolution aerial hyperspectral imagery
2. High-resolution visible (RGB) orthophotos
3. Terrain model (DTM) and surface model (DSM) derived from airborne LiDAR point clouds,
4. Medium-resolution (Landsat 8) satellite imagery.

---

\* Corresponding author

The participants have to provide thematic maps with a set of pre-defined categories. A set of calibration and validation samples will be provided to train and validate the various processing chains, and final evaluation will be carried out by the track leader (FÖMI) using independent samples. Competitive solutions have to fulfil the below criteria:

1. Create the best possible classification using the simplest set/combination of input sources
2. Try to reduce the number of input data for the processing
3. Develop algorithms that are fast to run
4. Find the best balance of complexity and accuracy (maximize efficiency) during the data processing

## 2. DATA DESCRIPTION

### 2.1 Study area description

The study area is located in the North-East part of Hungary, along the river Bodrog (Figure 1) which is a tributary to the river Tisza. The river Bodrog is crosses the Slovak–Hungarian border at the village of Felsőberecki (near Sátoraljaújhely) in Hungary, and continues its flow through the Hungarian county Borsod-Abaúj-Zemplén, until it joins the river Tisza in Tokaj city. The lowest point of the study area is on 77,19 m, while the highest point is on 258,73 m above the sea level. The study area is often affected by floods and inland excess water by the two rivers.



Figure 1. The location of the study area in Hungary

The study area is divided into three parts (Figure 2):

1. Area A: the full study area with medium-resolution Landsat 8 coverage

2. Areas B and C: two sub-areas with full high-resolution data coverage (hyperspectral and multispectral airborne data, DTM and DSM).
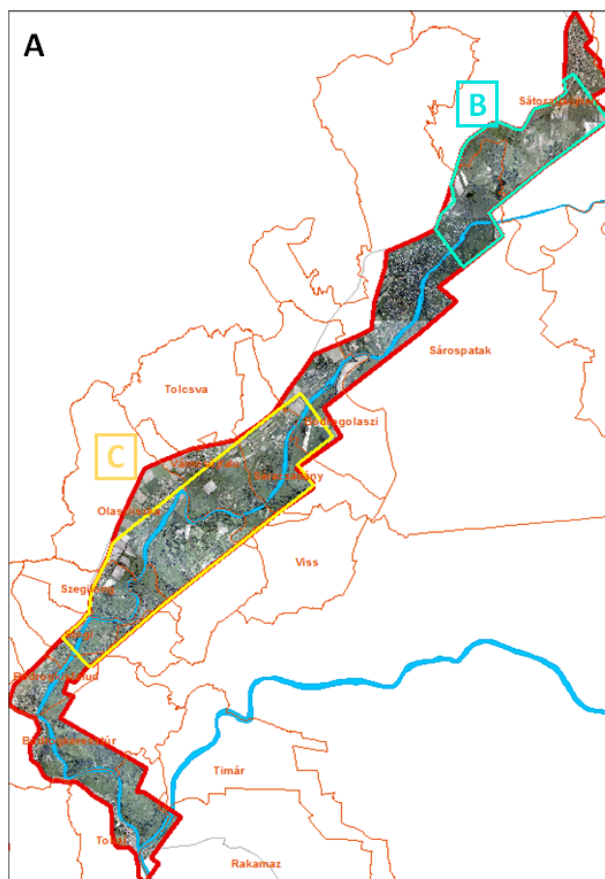


Figure 2. Study area

The participants have to provide thematic maps for each of the areas (A, B and C) with the below categories:

1. Water surfaces (code: 1)
2. Wet/waterlogged soils (code: 2)
3. Soils not directly affected by water (code: 3)
4. Vegetation standing in water (code: 4)
5. Vegetation not directly affected by water (code: 5)
6. Other (code: 0)

For each area, a set of training and verification samples being provided by the track leader in vector format, covering each of the above thematic categories. The participants can use the samples to train the classification algorithms and to verify the results. Areas with both high- and low-resolution data coverage (B and C) can be used to tune the algorithms for better performance on the larger area with only Landsat data (area A, Figure 3).
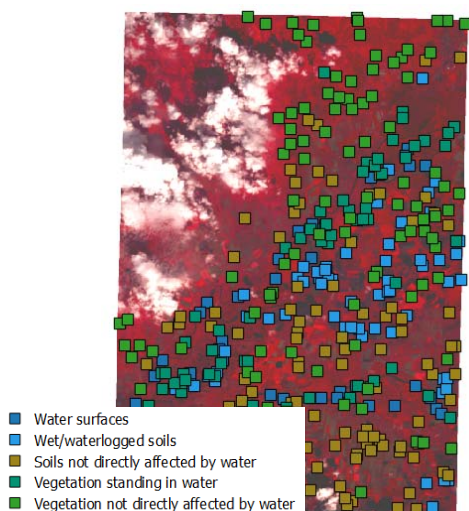
Figure 3. Distribution of categorized point dataset (reference) for the study area A.

## 2.2 Input data description

For area 'A' Landsat 8 data is provided, whereas hyperspectral images, LiDAR-derived Digital Surface (DSM), Terrain Model (DTM) and orthophotos (RGB) are additionally provided for areas 'B' and 'C'.

### 2.2.1 Hyperspectral images

Hyperspectral images are provided as georeferenced (UTM34N/WGS84) radiance data in *.dat (ENVI) format (16 bit BSQ) containing 128 bands (Figure 4). Additional technical details on spectral and spatial resolution and accuracy are the follows:

1.  Instrumentation / camera: aerial hyperspectral instrument (AISA Eagle)
2.  Spectral range: approx. 400-1000 nm (visible and reflected (near) infrared, VIS / VNIR)
3.  Spectral resolution: 5 nm
4.  Spatial resolution: 1.5 m / pixel
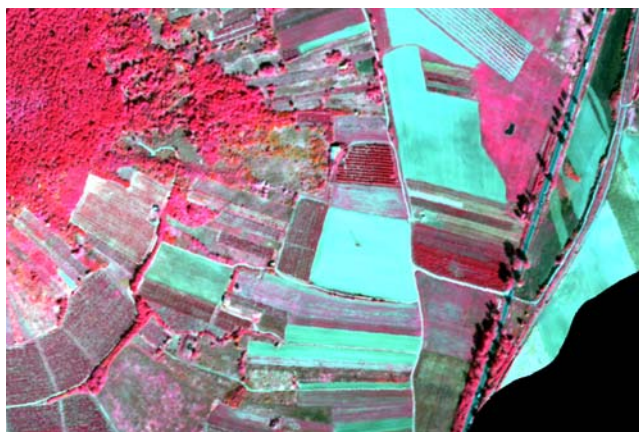5.  Spatial accuracy: 2.5 m (RMSE)



Figure 4. Hyperspectral image

### 2.2.2 Digital orthophotos

Multispectral (visible, RGB) orthophoto is provided as georeferenced (UTM34N/WGS84) and radiometrically corrected TIFF (Figure 5). The spatial resolution is 15 cm/pixel with spatial accuracy of 30 cm (RMSE)



Figure 5. Orthophoto

### 2.2.3 Digital Terrain Model (DTM)

The Digital Terrain Model is generated from aerial laser scanner (ALS) data, with gap filling by appropriate interpolation. Data is provided as georeferenced (UTM34N/WGS84), unsigned long integer TIFF. Data content is height values in cm, interpreted over Baltic Sea level. The Spatial resolution is 1 meter / pixel with spatial accuracy of 30 cm.
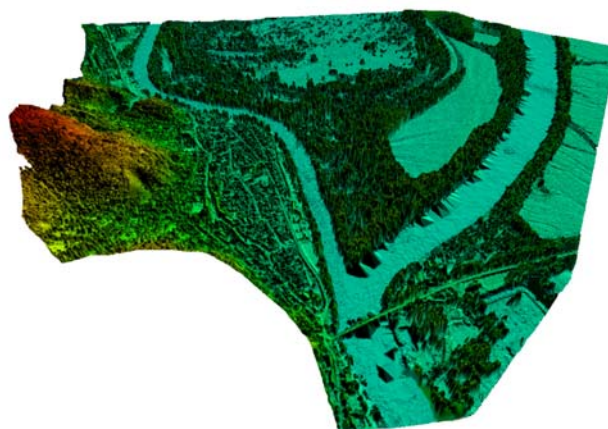


Figure 6. Digital Terrain Model for the conjunction of River Tisza and Bodrog at the city of Tokaj

### 2.2.4 Landsat 8 Data

For the whole area of interest we provide a subset of Landsat 8 imagery. The data is processed as standard Level-1 (orthorectified, terrain corrected), and provided as georeferenced (UTM34N/WGS84) and radiometrically corrected TIFF (unsigned integer 32 bit). Landsat 8 Operational Land Imager (OLI) images consist of nine spectral bands with a spatial resolution of 30 meters for Bands 1 to 7 and 9. New band 1 (ultra-blue) is useful for coastal and aerosol studies. New band 9 is useful for cirrus cloud detection.
The resolution for Band 8 (panchromatic) is 15 meters.
Thermal Infrared Sensor (TIRS) bands 10 and 11 are useful in providing more accurate surface temperatures and are collected at 100 meters, resampled to 30 meters.

## 3. EVALUATION

Participants had to submit the following material:

1. Georeferenced thematic rasters in GeoTIFF format, containing the codes of thematic categories as described above
2. Concise description of the whole methodology and processing chain (including algorithms and parameters, with references to relevant literature wherever available)

Evaluation and scoring was based upon the complexity, time- and resource efficiency of the methodology and the data requirement for processing.

## 4. TRACK RESULTS

Due to the deadline of the admission for the IQPC 2015 we have received one application to be presented at the GeoBigData Workshop. It was submitted by the team of the *Department of Physical Geography and Geoinformation Systems, University of Debrecen*. The team consist of three members: Zoltán Kovács, Boglárka Balázs and Szilárd Szabó. From this point on, in this paper we refer to them as "the Debrecen team".

As several possible participants have shown interest we are open to receive more solutions for the task published in Track 2 until the end of 2015 (but outside of the scope of the GeoBigData Workshop and IQPC 2015).

### 4.1 Solution provided by the Debrecen team

#### 4.1.1 Introduction

The aim of this contest track was to detect water surfaces using remotely sensed data. In our research we attempted to reduce the number of input data, therefore this task was approached from multispectral data source, so the medium-resolution OLI multispectral bands of Landsat 8 satellite imagery was used exclusively. Decreasing the processing time and human interaction, most of the steps were programmed in python and R programming languages.

According to the Track 2, the following categories were distinguished:

- Water surfaces (code: 1)
- Wet/waterlogged soils (code: 2)
- Soils not directly affected by water (code: 3)
- Vegetation standing in water (code: 4)
- Vegetation not directly affected by water (code: 5)
- Other (code: 0)

#### 4.1.2 Data processing

During data processing the training and two test datasets were available as ESRI point shape files, which were used to train the classification methods and verify the performance of the predictions. Intensity values of these pixels were extracted from the OLI multispectral bands (Band1-Band7, Band9) in ArcGIS 10.2. and saved as *.csv text files by using our python script.

The *.csv files were imported into an Excel workbook, where data was processed by HypDA (Hyperspectral Data Analyst) MS Excel add-in which was developed at Department of Physical Geography and Geoinformation Systems, University of Debrecen (Kovács & Szabó, 2013). It is being specially developed and designed for multispectral and hyperspectral data processing. HypDA is able to conduct hypothesis testing (e.g. Kruskal-Wallis test), separability testing (e.g. Jeffries-Matusita distance) and some classification procedures (e.g. minimum distance, Mahalanobis distance) for distinguishing the predefined classes by using self-generated spectral indices. HypDA workbooks contains specific worksheets, where the first two contains the intensity values and nominal or scale properties, the others the processing chain; one of them calculates matrices with all possibilities of bands to get the best available spectral indices, the other worksheets store the best values of matrices and show the detailed statistical background of the investigations of that certain models. In this case these methods were used to create spectral indices for each category for further investigations (classification in R). Our HypDA-based spectral indices were the following for the Landsat image:

- Code 1 <- (B5 - B3)/(B5 + B5)
- Code 2 <- (B9 - B1)/(B5 + B5)
- Code 3 <- (B7 - B2)/(B4 + B4)
- Code 4 <- (B7 - B1)/(B5 + B5)
- Code 5 <- (B5 - B3)/(B2 + B7)

R is free statistical software with command line interface (R Core Team, 2014). It is becoming more and more popular among scientific environments. It has many downloadable libraries providing widespread field of application. One of them, the Rattle was especially developed for data mining (Williams, 2011). Decision tree (DT, Mingers, 1989), Random forest (RF, Ho, 1995), Support Vector Machine (SVM, Nguyen and de la Torre, 2010) and Generalised Linear models (GLM, Bishop, 2006) can be used on training set to set classification rules for predicting classes for test sets and validation sets. For these models the target variable was the "gridcode" variable from shape files, the input variables were the spectral indices determined by HypDA. These indices were calculated for all test sets, training set and for the whole Landsat dataset. In order to verify the applied models, kappa indices and overall accuracy were calculated for all applied methods and all test sets (test1, test2 and test1+test2) from the confusion matrices of the models.

In order to create a classified geotiff map, the Landsat image was exported as ASCII file storing all intensity values of pixels. The lines of file represented individual pixels with the same structures as the training and test sets, therefore it was able to use as validation dataset for our investigations and Rattle was able to assign classes for each pixel. Beside the predicted classes, R can also determine the probabilities of classes for each pixel. In this way not only the classified pixels can be exported into text files, but the maximum probabilities of classes can be also saved and probability maps can be created for each method. Probabilities were cut by a predefined critical value (0.95), thus pixels with smaller probability were classified as "Others".

Classification procedure was summarized in the Fig. 7.

Besides masking the values having <95% classification probability, we masked the areas covered by clouds, too. Altogether 125 test data point fell in cloud covered areas; therefore, we provide two solutions for the accuracy assessment:

- (1) calculations including all test points (test1, test2 and test1+test2)
- (2) calculations omitting the cloud covered areas.

In our understanding, solution (2) is more reliable, as pixels covered by clouds are biased and have distorted intensity values.
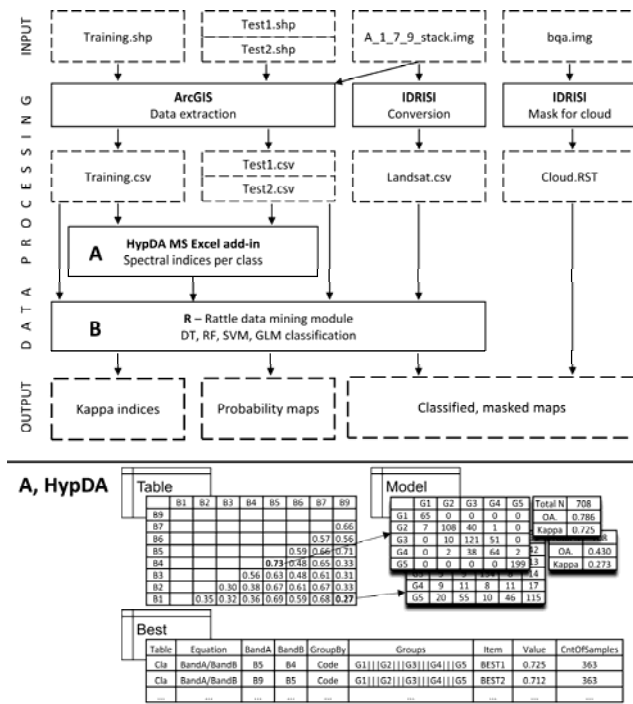


Figure 7. Classification procedure by the Debrecen team

### 4.1.3 Accuracy assessment

We compared the performance of the classification techniques based on Kappa values and found the GLM (General Linear Model) having the best performance. Confusion matrix and the accuracy statistics are shown on Figure 8.



Figure 8. Confusion matrix and accuracy statistics of the solution provided by the Debrecen team

## 5. CONCLUSIONS

The Debrecen team has provided an elegant solution to the challenge raised in the context of water detection and classification in the frame of IQmulus Processing Contest 2015. Their solution is based on Landsat-8 imagery and does not use any other input data, yet provides very accurate results for the desired categories.

Although the first deadline for submitting track solutions has passed and the track has been successfully concluded, the track organizers decided to keep the possibility of result submission open to enable collecting a variety of approaches and solutions for this interesting problem. Potential additional solutions would also be presented in the frame of the track summary presentation at the GeoBigData workshop and published at the project homepage.

## ACKNOWLEDGEMENTS

## REFERENCES

IQmulus Processing Contest 2015. Water detection and classification on multi-source remote sensing and terrain data Proposal [Online]. Available at http://iqmulus.eu/iqpc/iqmulus-processing-contest-2015 [June 25, 2015].

D.8.8.2 Report on IQmulus Processing Contest – year 2 [Online]. Available at http://iqmulus.eu/dynamic/document/D8.8.2_Report_on_the_IQmulus_Processing_Contest_Year_2.pdf [June 25, 2015].

Bishop, C.M., 2006. Pattern recognition and Machine Learning. Springer, Cambridge, pp. 729

Ho, T. K., 1995. Random decision forests, in Proceedings of the 3rd International Conference on Document Analysis and Recognition, Montreal, Canada, IEEE Computer Society, pp. 278-282.

Kovács, Z., Szabó, Sz., 2013. Interactive spectral evaluation in Excel environment – add-in for hyperspectral data processing. In Lóki, J. ed.: Proceedings of 4th Hungarian Conference and Exhibition of Geoinformation Science, Debrecen, Hungary, pp. 251-255. (in Hungarian)

Li, W., Du, Z., Ling, F., Zhou, D., Wang, H., Gui, Y., Sun, B., Zhang, X., 2013. A Comparison of Land Surface Water Mapping Using the Normalized Difference Water Index from TM, ETM+ and ALI. Remote Sensing 5(11) pp. 5530-5549.

Mingers, J., 1989. An empirical comparison of selection measures for decision-tree induction. Machine Learning 3(4), pp. 319-342.

Nguyen, M. H., De la Torre, F., 2010. Optimal feature selection for support vector machines, Pattern Recognition 43 (3), pp. 584-591.

R Core Team, 2014. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Available at http://www.R-project.org/

Rokni, K., Ahmad, A., Selamat, A., Hazini, S. 2014. Water Feature Extraction and Change Detection Using Multitemporal Landsat Imagery. Remote Sensing 6(5) pp. 4173-4189.

Williams, G. J., 2011. Data Mining with Rattle and R: The Art of Excavating Data for Knowledge Discovery, Use R!, Springer.