# Uncertainty analysis and data-assimilation of remote sensing data for the calibration of cellular automata based land-use models

**J. van der Kwast**[a], **L. Poelmans**[b], **T. Van de Voorde**[c], **K. de Jong**[d], **I. Uljee**[b], **D. Karssenberg**[d], **F. Canters**[c], **G. Engelen**[b]

[a] *Department of Water Science and Engineering, UNESCO-IHE Institute for Water Education, P.O. Box 3015, 2601 DA Delft, The Netherlands (j.vanderkwast@unesco-ihe.org)*

[b] *Flemish Institute for Technological Research (VITO), Boeretang 200, 2400 Mol, Belgium (lien.poelmans@vito.be, inge.uljee@vito.be, guy.engelen@vito.be)*

[c] *Cartography and GIS Research Group, Department of Geography, Vrije Universiteit Brussel, Pleinlaan 2, BE-1050, Brussels, Belgium (tvdvoord@vub.ac.be, fcanters@vub.ac.be)*

[d] *Department of Physical Geography, Faculty of Geosciences, Utrecht University, P.O. Box 80115, 3508 TC, Utrecht, The Netherlands (k.dejong1@uu.nl, d.karssenberg@uu.nl)*

**Abstract:** A correct historic calibration of land-use models is important, because they are more and more used by decision makers. Existing calibration methods, however, do not sufficiently take into account uncertainties in input parameters. For that reason, uncertainties that propagate through simulations of future land use are mostly unknown. When uncertainties in model parameters can be estimated, Monte Carlo modelling can be used to approximate the uncertainty in simulated land-use patterns. This study shows that uncertainty information is not only indispensible for the interpretation of the output of land-use models, but that this can also be used to improve the calibration of land-use models by means of data assimilation. For this purpose the MOLAND model of Dublin has been integrated in a Python-based data-assimilation framework. Calibration of the land-use model is based on the comparison of spatial metrics derived from historic remote sensing images and land-use simulation results. Remote sensing derived probability density functions of class metrics are assimilated in the model at time steps for which they are available. The particle filter algorithm only continues successful realizations of the land-use model, thus reducing its uncertainty by removing unsuccesful particles from the ensemble. In this way, only parameters sets used in the simulation that match the patterns observed in the remote sensing imagery as quantified by the class metrics, are used in the simulation of future land use. It is expected that the automatic procedure results in an improved calibration of the land-use model. Furthermore, it provides data on the uncertainty of the results, which is important for drawing conclusions from simulations.

*Keywords:* land-use modeling; spatial metrics; stochastic modeling; remote sensing; calibration

# 1 INTRODUCTION

Cellular automata based land-use models are increasingly used as instruments for policy makers and planners to assess the impact of their policies on the sustainable development of urbanized areas. Because the use of the models has moved from an academic environment to end users who base their decisions on simulated scenarios of future land-use, the tools must be robust and reliable. In order to achieve this, the calibration of these models should be based on the best available scientific knowledge and data [van der Kwast et al., 2011].

Land-use models are normally calibrated using a historic calibration [Engelen and White, 2007]. In a historic calibration procedure, the land-use model is initiated using a historic land-use map. Next, the model hindcasts the land-use at a more recent time step for which a land-use map is available. Dedicated goodness-of-fit measures are used to evaluate the similarity between the hindcast and the actual map.

The historic calibration, however, is often hampered by the poor availability of time series of high quality land-use maps. In order to be less dependent on time series of land-use maps van der Kwast et al. [2012] proposed a calibration procedure based on time series of archived medium resolution remote sensing data. Because of the complex and indirect relation between land-use features at the earth's surface and the reflection received by the sensor it is impossible to derive land-use maps with the same thematic resolution as land-use maps produced by human interpretation. Therefore, the calibration procedure proposed by van der Kwast et al. [2012] uses spatial metrics derived from remote sensing data and simulations to quantify differences between both land-use maps. In this procedure, parameters in the land-use model are tuned in order to derive an optimal fit with spatial metrics derived from remote sensing data.

van der Kwast et al. [2011] identified that in existing calibration methods uncertainties in parameters, input data, and reference data are neglected. They proposed a methodology to incorporate the metric-based calibration method of van der Kwast et al. [2012] in a probabilistic framework in order to quantify and reduce the uncertainties in land-use simulations using data-assimilation techniques. This paper will focus on the modelling of propagation of uncertainties in a land-use model and in land use inferred from remote sensing, which forms the basis of the data-assimilation methodology. For the analysis of uncertainties in spatial metrics derived from land-use modelling a simplified version of the MOLAND model for Dublin [Engelen et al., 2007] is used. Time series of Landsat TM/ETM+ images are compared with urban masks derived from the probabilistic land-use model.

# 2 METHODS

## 2.1 The MOLAND Model of Dublin City

Dublin is the political, economic and cultural capital of Ireland. From the mid-1990s until 2008, urban growth in Dublin was characterized by a rapid sprawl. Because of the large dynamics, Dublin is one of the key cities in the MOLAND (Monitoring Land Use/Cover Dynamics) project [Barredo et al., 2003].

The MOLAND model simulates land-use dynamics at the cellular level using a constrained cellular automata [Engelen et al., 2007]. The model simulates the likely future development of land use with a temporal resolution of one year, given alternative planning and policy scenarios and socio-economic trends. Details of the MOLAND model
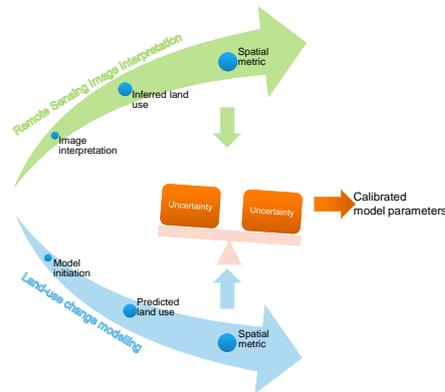
Figure 1: Concept of using spatial metrics for calibrating land-use models [van der Kwast et al., 2011]. The width of the arrows represents the propagation of uncertainty.

can be found in Barredo et al. [2003] and Engelen et al. [2007]. For this study four land-use classes are modelled for the city of Dublin: population related classes, employment related classes, non urban, other.

## 2.2 Metric-Based Calibration Framework

The historic calibration of the MOLAND model consists of four steps: (1) a first set of neighbour influence function (attraction-repulsion curves) parameters is fixed. Generally, functions are taken from a library based on previous applications of the model; (2) a second parameter determining largely the dispersal of land use and size of clusters is estimated; (3) cell-specific information represented in the suitability, accessibility and zoning maps are introduced; (4) fine tuning of the model starts (repeating of loop 1-2-3 with dedicated statistical analysis). The steps are repeated until the simulated land-use map is sufficiently similar to the actual one.

Figure 1 shows the concept of using spatial metrics derived from remote sensing in the calibration of land-use models. First, morphological/functional classes are derived from each remote sensing image that has been acquired within the calibration interval. Spatial metrics derived from the land-use map inferred from remote sensing are then compared with the same metrics derived from the simulated land-use for the same timestep. In this study the class metrics evaluated in van der Kwast et al. [2012] were calculated from urban masks, indicating the part of the landscape that is urbanized, derived from the MOLAND model and remote sensing images.

## 2.3 Probabilistic Framework for Calibration of Land-Use Models

The objective of probabilistic calibration is to calibrate model parameters given two sources of information: prior knowledge of possible ranges of model parameters and observations of the land-use system considered. Here we combine these sources of information by running the land-use model in Monte Carlo mode and solving Bayes theorem at each time step for which observations of the land-use system are available from remote sensing data, i.e. an observation time step. At the end of the model run this results in calibrated (so-called posterior) probability distributions of model parameters that are adjusted as a result of the integration of the observations in the model.

Various approaches exist to solve Bayes theorem at each observation time step. Here we use the Particle Filter, which has the advantage that adjustment of the full state of the model at observation time steps is done by duplicating or removing realizations of the model. Thus, unlike other approaches such as the Ensemble Kalman Filter, the Particle Filter does not require adjustment of individual realizations. This can be considered a large advantage in land-use modelling, as adjustment of the state of land-use models would involve changing the pattern in land uses, which can only be done using a rather arbitrary procedure. Here, we provide a short explanation of the Particle Filter. For an extensive description, the reader is referred to Karssenberg and Bierkens [2012].

Let $\mathbf{x}_t$ be the full state of the system, described by the joint probability density of all variables (state, parameters, drivers), denoted by $p(\mathbf{x}_t)$, with the time index $t, t = 0, 1, \ldots, T$. The land-use model $\mathbf{f}$ updates this full state for each time step by evaluating $\mathbf{x}_t = \mathbf{f}(\mathbf{x}_{t-1})$. At an observation time step, the full state $\mathbf{x}_t$ is adjusted by solving Bayes theorem:

$$p\left(\mathbf{x}_t \mid \mathbf{y}_t\right) = \frac{p\left(\mathbf{y}_t \mid \mathbf{x}_t\right) p\left(\mathbf{x}_t\right)}{p\left(\mathbf{y}_t\right)} \tag{1}$$

with, $p(\mathbf{x}_t)$, the joint probability density of the full state of the model at $t$, retrieved by propagating the system model from the previous observation time step up to the current observation time step; $p(\mathbf{y}_t)$, the joint probability density of the observational data at $t$; and $p(\mathbf{y}_t \mid \mathbf{x}_t)$, the probability density of the observations at $t$ given the model state (likelihood). The Particle Filter solves this scheme with Monte Carlo simulation. A number of $N$ independent model realizations (or particles) $\mathbf{x}_t^{(n)}$ are sampled from the initial, prior, joint probability density of the model, which includes all parameters that are calibrated. At each observation time step, (1) is solved by duplicating individual particles a number of times proportional to the weight $p\left(\mathbf{y}_t \mid \mathbf{x}_t^{(n)}\right)$ of each particle, calculated as:

$$p\left(\mathbf{y}_t \mid \mathbf{x}_t^{(n)}\right) \sim \exp\left(-\frac{1}{2}\left[\mathbf{y}_t - \mathbf{H}_t\left(\mathbf{x}_t^{(n)}\right)\right]^T \mathbf{R}_t^{-1}\left[\mathbf{y}_t - \mathbf{H}_t\left(\mathbf{x}_t^{(n)}\right)\right]\right) \tag{2}$$

With $\mathbf{H}_t$, the observation operator, mostly taken as one, and $\mathbf{R}_t^{-1}$ the covariance matrix of the errors in the observations. The use of the error covariance matrix is essential, as it guarantees that observations are weighted according to their measurement error.

The Particle Filter calibration scheme is run using the PCRaster Python data-assimilation framework [Karssenberg et al., 2010], with some modifications to suspend and resume the external land-use model for each observation time step, and each Monte Carlo sample.

## 2.4 Uncertainty Propagation in Land-Use Modeling

Land-use modelling involves uncertainty caused by attribute errors, positional errors, logical inconsistencies, incompleteness and temporal errors in the model and in the reference land-use maps used for initiation and calibration [van der Kwast et al., 2011]. Positional errors of the georeferenced input maps are assumed to be smaller than the spatial resolution of the model i.e. 200 m. Logical consistency and completeness have been tested in previous applications of the model. Temporal accuracy is determined by the synchronization of temporal input data and the model time step. These uncertainties are assumed to be small, because land-use changes emerge over periods longer than the temporal resolution. These errors are therefore ignored in this study. Uncertainties

in the reference land-use maps, however, can be important, but are difficult to quantify objectively. For this reason the only uncertainties that will be considered here are uncertainties in the input parameters that are calibrated in a historic calibration. These are the neighbourhood influence functions and the stochastic $\alpha$ parameter.

First, the uncertain input parameters of the land-use change models need to be quantified. The neighbourhood influence functions are approximated by the following equation:

$$W_{l,k,d} = s + \left( -s \times \left( 1 - exp^{-\left(\frac{d}{r}\right)} \right) \right) \tag{3}$$

In which $W_{l,k,d}$ is the weighting parameter that represents the attraction (positive value) or repulsion (negative value) between land-use class $l$ and land-use class $k$ at distance $d$. $s$ is the sill and represents the value of $W_{l,k,d}$ at distance 0 and range $r$ represents the distance-decay effect of the interactions.

A drawback of using this approximation is that only exponentially shaped neighbourhood influence functions can be taken into account. The major advantage of using (3) is, however, that the neighbourhood influence functions can be estimated by only 2 parameters (sill and range), whereas the original influence rules were represented by at least 5 different parameters. The calibration procedure can thus be limited to finding optimal sills and ranges for the interactions between the population related and the employment related land-use classes.

First of all, the search radius for all sill and range values were narrowed down based on visual evaluations and expert knowledge. Hereto, the model was run using various combinations of sill and range for the simulated land-use classes. In the same way the mean and standard deviation of the stochastic $\alpha$ parameter has been estimated. Next, the deterministic land-use models are run as stochastic models using Monte Carlo techniques in order to propagate different scenarios of uncertainties in the neighbourhood influence functions and $\alpha$ parameter through the model. This will provide land-use simulations accompanied by probability maps. Finally, the probability density function (pdf) of the spatial metrics can be derived from the land-use model. For the comparison with remote sensing data the population and employment related classes have been aggregated to derive urban masks for each time step.

### 2.5 Uncertainty Propagation in Urban Land-Use Interpretation from Remote Sensing

Van de Voorde et al. [2010] inferred urban land-use from Landsat images with high accuracy using the following processing chain. First, a binary urban mask is needed, which separates the urban and the non-urban areas. The urban mask is produced by an unsupervised classification using a Kohonen self-organising map neural network, which was improved by a knowledge-based post-classification approach [Van de Voorde et al., 2007]. A change trajectory analysis is applied to remove irrational changes between urban and non-urban areas [van der Kwast et al., 2012]. Next, pixels outside of the urban mask are considered to have zero impervious surface cover, while pixels inside the urban mask are considered as mixtures of impervious surfaces and vegetation. Therefore, the next step in the processing chain consists of classifying pixels within the urban mask using a sub-pixel classification technique. Because the impervious surface fraction of a pixel is the complement of the vegetation fraction, a multiple linear regression model is used to estimate the vegetation fraction in each pixel. A high resolution land-cover classification of a Quickbird image acquired on August $4^{th}$ 2003 has been used to train

and validate the multiple regression model [Van de Voorde et al., 2010]. A temporal filtering technique based on iterative linear regression between NDVI values [Van de Voorde et al., 2009] is used to remove pixels in the sample that were changed between the acquisition dates of the Landsat image and the Quickbird image. In the next step in the chain, land-use is inferred from the impervious surface map using a supervised classifier based on a neural network (MLP architecture). Training and validation data for the neural network classifier is obtained by deriving the predominant land-use class of each building block using the MOLAND land-use map. A stratified sampling approach is used to randomly select 200 blocks for each land-use class. Half of the sample is used for training, the other half for validation. Classification signatures derived from morphological characteristics of the training sample consist of the 4 parameters of a transformed logistic function fitted to the cumulative frequency distribution of impervious surface fractions within each building block, and the average impervious surface fraction within each building block, resulting in 5 parameters. The MLP classification results in three classes: residential land use, industry and services, urban green. Finally, the residential class is further subdivided in two classes, based on an intersection with the impervious percentage map. Building blocks with an impervious percentage between 10 and 50 % are assigned to the low density residential class. Blocks with an impervious percentage between 50 and 80 % are assigned to the medium density residential class. Blocks consisting of less than 10 % impervious percentage are assigned to the non-urban class and the remaining blocks covered for more than 80 % with impervious surfaces are assigned to the dense urban class. This results in a map with 6 land-use classes.

This processing chain has been applied to a time series of Landsat images covering the years 1988, 1997, 2001 and 2006. Uncertainties propagate through the entire processing chain each time it is applied and cumulate into the overall uncertainty present within the metrics derived from the inferred land-use map. The quantification of uncertainties in the Landsat products due to sensor calibration, radiometric, atmospheric and geometric correction were out of the scope of this study and should be included in future research. Uncertainties in the derivation of the urban masks are negligible compared to uncertainties in other parts of the processing chain. The pdf of the estimation errors of the impervious surface proportions is derived from validation data. A normal distribution has been assumed. Realisations of impervious surface maps only based on this pdf, however, do not include the spatial autocorrelation of the uncertainties. For this reason a first order autoregressive error model [Heuvelink, 1998] has been applied to generate spatially autocorrelated random error-fields. This autoregressive model defines the value of each cell $(i, j)$ in a regular grid as a weighted combination of the value of four neighbour cells and a random value $(r\epsilon[i, j])$ characterising the cell under consideration:

$$Z[i,j] = q(Z[i-1,j] + Z[i+1,j] + Z[i,j-1] + Z[i,j+1]) + r\epsilon[i,j] \qquad (4)$$

The resulting error field is then added to the impervious surface map. Parameter $q$ in this model controls the degree of autocorrelation and is empirically determined using an experimental semi-variogram. The experimental variogram is estimated from an error image of the impervious percentage map of 2001. The error map is calculated by subtracting the reference proportion of impervious percentage, which is derived from the downsampled Quickbird classification, from the predicted proportion of impervious percentage. Temporal variant pixels were excluded from the analysis.

Uncertainties in the previous step of the chain are propagated through the MLP classifier. The neural network has three output nodes, i.e. one for each land-use class to be derived. In a hard classification each building block is assigned to the node with the highest activation level. In reality, however, for each pixel all three nodes are activated to
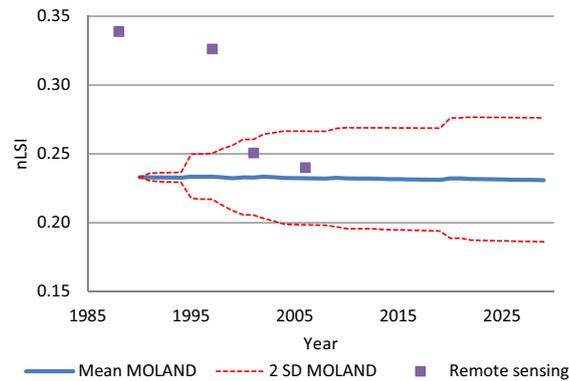
Figure 2: Temporal evolution of the nLSI metric derived from urban masks from MOLAND and remote sensing images. Uncertainty bands of 2 standard deviations are shown.

a certain level. The activation level ranges from 0 to 1 and the sum of all nodes is 1. The difference in activation level among the nodes, particularly between the highest level and the other nodes, is a measure for the degree of uncertainty in the MLP classification. This uncertainty is modelled by drawing a random number between 0 and 1 from a uniform distribution. The three activation levels can be considered as three partial intervals within the range 0–1. Based on the random number and the ranges of the partial intervals a building block is assigned to one of the three land-use classes. Obviously, a land-use class with a wider partial interval has an increased chance to be selected.

The intersection with the impervious percentage map in the final step of the processing chain in order to derive 6 land-use classes is separately applied to each realisation, because this can also result in a change in land-use classes. Finally, the spatial metrics are derived from each realisation and the pdf of the entire ensemble of spatial metrics is calculated.

The workflow of deriving pdf's from spatial metrics derived from remote sensing is currently being automated using Python scripts linking different software packages. In this paper only deterministic results of spatial metrics derived from the urban mask are compared with the probabilistic result of the land-use model.

## 3  RESULTS

Figure 2 shows a preliminary result of the error propagation of the normalized Landscape Shape Index (nLSI), indicating the degree of aggregation of the landscape, calculated from urban masks derived from the Monte Carlo simulation of the MOLAND model. Remote sensing derived nLSI is also shown in a deterministic way.

## 4  CONCLUSIONS AND RECOMMENDATIONS

This paper discussed the methodology for estimating the propagation of uncertainty through the MOLAND model for Dublin and reference data derived from a remote sensing processing chain. Furthermore, a probabilistic framework has been developed for the calibration of the MOLAND model using the probability distribution functions of spatial metrics derived from MOLAND and remote sensing images. Further research should confirm if the proposed particle filter data assimilation method results in an improved cal-

ibration of the MOLAND model. Different (sets of) metrics and thematic resolutions will be evaluated for this purpose.

## REFERENCES

Barredo, J. I., M. Kasanko, N. McCormick, and C. Lavalle. Modelling dynamic spatial processes: simulation of urban future scenarios through cellular automata. *Landscape and Urban Planning*, 64(3):145–160, July 2003.

Engelen, G., C. Lavalle, J. Barredo, M. van der Meulen, and R. White. The MOLAND modelling framework for urban and regional land-use dynamics. In Koomen, E., Stillwell, J., Balkema, A., and Scholten, H., editors, *Modelling land-use change progress and applications*, pages 297–320. Springer, Dordrecht, 2007.

Engelen, G. and R. White. Validating and calibrating integrated cellular automata based models of land use change. In Albeverio, S., Andrey, D., Giordano, P., and Vancheri, A., editors, *The dynamics of complex urban systems. An inderdisciplinary approach*, pages 185–212. Physica-Verlag, Heidelberg, 2007.

Heuvelink, G. *Error propagation in environmental modelling with GIS*. CRC Press, 1998.

Karssenberg, D. and M. F. P. Bierkens. Early-warning signals (potentially) reduce uncertainty in forecasted timing of critical shifts. *Ecosphere*, 3(2):art15, February 2012.

Karssenberg, D., O. Schmitz, P. Salamon, K. de Jong, and M. F. Bierkens. A software framework for construction of process-based stochastic spatio-temporal models and data assimilation. *Environmental Modelling & Software*, 25(4):489–502, April 2010.

Van de Voorde, T., W. De Genst, and F. Canters. Improving pixel-based VHR land-cover classifications of urban areas with post-classification techniques. *Photogrammetric Engineering and Remote Sensing*, 73(9):1017–1027, 2007.

Van de Voorde, T., T. De Roeck, and F. Canters. A comparison of two spectral mixture modelling approaches for impervious surface mapping in urban areas. *International Journal of Remote Sensing*, 30(18):4785–4806, 2009.

Van de Voorde, T., J. van der Kwast, I. Uljee, G. Engelen, and F. Canters. Improving the calibration of the MOLAND urban growth model with land-use information derived from a time-series of medium resolution remote sensing data. In Taniar, D., Gervasi, O., Murgante, B., Pardede, E., and Apduhan, B. O., editors, *Computational Science and Its Applications ICCSA 2010*, volume 6016 of *Lecture Notes in Computer Science*, pages 89–104. Springer Berlin Heidelberg, Berlin, Heidelberg, 2010.

van der Kwast, J., F. Canters, D. Karssenberg, G. Engelen, T. Van de Voorde, I. Uljee, and K. de Jong. Remote sensing data assimilation in modeling urban dynamics: Objectives and methodology. *Procedia Environmental Sciences*, 7:140–145, January 2011.

van der Kwast, J., T. Van de Voorde, M. Binard, G. Engelen, Y. Cornet, I. Uljee, C. Lavalle, F. Canters, L. Poelmans, H. Shahumyan, B. Williams, and S. Convery. Using spatial metrics derived from remote sensing for the calibration of land-use change models. *Landscape and Urban Planning*, 2012. In review.