# SOILSIM, a GIS-based framework for data-extensive modelling of the spatial distribution of soil hydrological characteristics in small alpine catchments

Martin MERGILI, Clemens GEITNER, Andrew MORAN, Michael FECHT and Johann STÖTTER

## Abstract

SOILSIM is a GIS-based modelling framework for the spatial interpolation of soil characteristics from individual points to a raster map. It was designed for providing a reliable approximation of the hydrological characteristics of the soils in small alpine catchments with a minimum of input required, and fully dependent on open source products (GRASS and R). A table representing the soil characteristics and a set of habitat maps are required as input. The program operates in two major steps. (1) A linear multiple regression is fitted for each soil variable (the predictors are chosen using ANOVA). The regression equations are then applied to raster maps of the predictor variables. (2) Hydrologically relevant soil characteristics (e.g. field capacity, saturated hydraulic conductivity) are calculated as response to the modelled soil variables.

SOILSIM was applied to the Stampfanger catchment (near Kitzbühel, Tyrol, Austria; 23.1 km²). The results of the study indicated that the method is quick and easy, but that it has to be applied with much caution and with qualitatively and quantitatively sufficient input data in order to provide reliable predictions for the variables under investigation.

## 1 Introduction

The characteristics of the soil of a certain habitat are determined by a complex interaction of influences like climate, substrate, relief, vegetation, soil fauna, humans, and the timespan since the beginning of soil development. The dominant conditioning variables (predictors) for the soil characteristics are not always obvious (TASSER et al. 1998). The situation becomes even more complex due to the strong interrelation between soil and vegetation. Furthermore, the soil characteristics in mountain habitats vary on different spatial scales, from the sub-meter scale (for example on rock fall deposits) to the scale of kilometres as response to vertical temperature or precipitation gradients (compare HILLER et al. 2002; HITZ et al. 2002).

Soil characteristics themselves are usually measured or estimated at individual points, using samples obtained at soil profiles and boreholes. This information is sufficient for some, but nor for all purposes. Modelling approaches for hydrological assessments of catchments, for example, require the knowledge of the soil characteristics over the entire study area which are difficult to measure. A dense network of sampling sites, in combination with sophisticated equipment, is therefore necessary for such studies. In practice, both resources time and money are limited. The scale of existing soil maps is usually too small to be used for

modelling of hydrological processes (DE GRUIJTER et al. 1997). This is particularly true for mountain regions.

An alternative way is to create a smaller set of sampling sites, including all major types of habitats and covering a wide range of combinations of habitat conditions that are supposed to be important for the soil characteristics. Such a dataset can be interpolated to the whole study area, using GIS in combination with an appropriate statistical method. A major requirement for following this approach is the full knowledge of the spatial distribution of the habitat conditions (climate, topography, substrate, landcover).

Various models were applied for this purpose. MCBRATNEY et al. (2003) provided a detailed overview of methods used in the past, discussing generalized linear models, regression and classification trees, neural networks, fuzzy systems, and geostatistics.

Many of these methods are only applicable to easy measurable soil characteristics like soil depth, soil skeleton, etc. The indirect derivation of hydrological soil characteristics has therefore been the subject of several studies (e.g. VEREECKEN 1995, UFZ-Umweltforschungszentrum Leipzig-Halle 2001, LEHMANN et al. 2005), resulting in a certain number of approaches.

Within the framework of this study FECHT et al. (2005) used geostatistics (cluster analysis) in order to delineate subareas with similar combinations of conditioning factors for the Ruggbach study area (Vorarlberg, Austria). However, it appeared to be problematic to fill these classes with content (soil variables). Hence it was aimed at developing a simple, data-extensive method (soil spatial distribution model: SOILSIM) for (1) directly modelling soil variables from predictor variables; and (2) deriving hydrologically relevant soil characteristics from simple, easily measurable variables. The model was tested for the Stampfanger catchment (Figure 1).
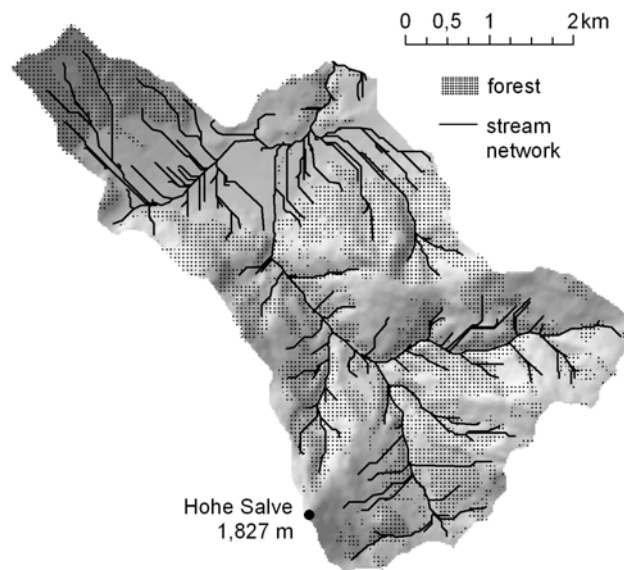


**Figure 1:** The Stampfanger catchment.

## 2 Study area

The Stampfanger catchment is located in the municipality of Söll (NW of Kitzbühel, Tyrol, Austria). It comprises an area of 23.1 km². The terrain is mountainous, but in general it is not very steep and rocky (except a small part in the N). It extends from about 640 m up to 1,827 m a.s.l. (Hohe Salve). The majority of the catchment is part of the Greywacke Zone, with a mixture of different rock types, including some volcanics. The NW edge is part of the Northern Limestone Alps. The majority of the area, however, is covered with till of different thickness, and with alluvium. The dominant landcover types are mixed forest on the lower slopes, spruce forest on the middle slopes and pastures and meadows in the upper part of the catchment. The lowest part carries meadows and human settlements. The catchment has been chosen for this study due to the availability of vegetational, geological, climatic and hydrological data (Moran et al. 2005).

## 3 Methods

SOILSIM is based on a combination of the GRASS GIS (http://grass.itc.it) and the R statistical software (http://www.r-project.org), both distributed under the open source license. The model was realized as a shell script with integrated python, R and GRASS functions. It requires soil data for various points and a number of habitat raster maps as input. The model involves the following steps (compare McBratney et al. 2003): (1) preprocessing of habitat variables, (2) choosing a set of predictors from the habitat variables, (3) fitting regression models for explaining the soil variables, (4) applying these relationships to the habitat maps, (5) relating the soil variables to hydrological characteristics and applying the relationship to raster maps.
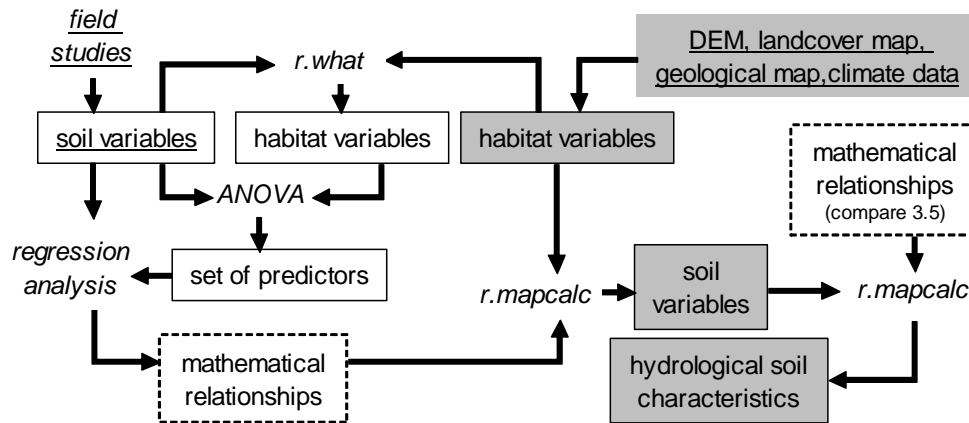


**Figure 2:** Illustration of the work flow of SOILSIM. Spatial datasets (raster maps) are shaded, input is underlined, mathematical operations are written in italic letters.

## 3.1 Soil sampling

Soil data was obtained at 52 study sites in the Stampfanger catchment. The sites were chosen subjectively in order to cover a wide range of different habitats. A soil auger and soil profiles were used to investigate the variables listed in Table 1.

**Table 1:** Soil variables under investigation

| soil variable [unit] | method |
|---|---|
| coordinates [m UTM] | GPS device |
| total soil depth [mm] | measured with soil auger ("Pürckhauer"), values above 100 cm were recorded as 100 cm |
| depth of the A horizon [cm] | measured at profile |
| pH [-] | measured at profile using field kit; only for A horizon |
| color [-] | estimated using Munsell Soil Color Charts; only for A horizon |
| organic content [vol-%] | calculated from pH, colour and texture (RENGER et al. 1987; SCHLICHTING et al. 1995) only for A horizon, set to 0 for the other horizons |
| texture [μm] | estimated from the profiles and the soil auger using finger test and converted into average grain size using tabular data (SCHLICHTING et al. 1995) |
| bulk density [kg dm$^{-3}$] | estimated at profile (SCHLICHTING et al. 1995) separately for A horizon and remaining horizons |
| soil skeleton [vol-%] | estimated at profile separately for A horizon and remaining horizons |

## 3.2 Preparing habitat maps

Raster maps representing the potential predictor variables for the spatial distribution of the soil variables were generated using a DEM (20 m resolution), a geological map, a land-cover map and meteorological data.

- the **cumulative air temperature of the growing season** (TCUM) was calculated using the DEM and temperature data from a meteorological station in the vicinity, and a vertical temperature gradient. It was corrected for aspect using solar irradiation (r.sun) and an empirical relationship (Welpmann 2003).

- the **substrate type** (STYPE) was derived from a geological map and coded with 1 = silicate, 2 = intermediate and 3 = carbonate. It was considered numeric.

- the **substrate physical properties** (SPHYS) were also derived from the geological map and coded with 1 = bedrock, 2 = block material, 3 = finer material (alluvium, moraine). It was considered as numeric, too.

- **slope** (SLOPE) in degrees was derived from a DEM (r.slope.aspect).

- the **topographic index** (TIND) was derived from a DEM, too (r.topidx).

- the **C/N ratio of the litter** (CNRAT) was integrated assigning average values derived from the literature (Scheffer & Schachtschabel 2002) to the landcover units.

- **anthropogenic disturbance** (DIST) was roughly estimated from the landcover map:
  1 = forest (little disturbance), 2 = meadows and pastures (intermediate disturbance),
  3 = skiing slopes and cultivated places (heavy disturbance). The variable was considered as numeric.

TCUM, SLOPE, CNRAT, TIND and DIST are generated automatically within the SOIL-SIM framework. The values of all predictor variables at the study sites are extracted from the map (r.what function of GRASS) and joined to the soil variables table.

## 3.3   Analysis of variance (ANOVA)

The purpose of the analysis of variance in general is to investigate whether the variance of a continuous variable can be explained by one or more categorical independent variables (predictors). The method is based on the F-test. Each soil variable (Table 1) was tested against each predictor (as categorial variable) in order to determine the variables to be used in the subsequent regression analysis. For each soil variable the three most significant predictors were used for multiple regression as the size of the dataset did not allow for larger numbers. The critical level of significance was set to 0.10 (instead of using the standard value of 0.05), and predictors with higher values were excluded. ANOVA was performed using the R statistical software. The analysis itself is included in the SOILSIM framework. The choice of significant predictor variables and the formulation of the regression models have to be performed manually.

## 3.4 Regression analysis and map generation

Linear multiple regression is a method of multivariate statistics that is widely used in various fields of science and technology. The method was applied in order to establish mathematical relationships between the soil variables (Table 1) and a set of predictor variables, the latter selected in accordance with the results of the ANOVA. The relationships, combined with the predictor maps, were then used for the generation of raster maps representing the spatial distribution of the soil variables. This process is fully included in the SOIL-SIM framework. For the regression itself, the data is automatically handed over to the R statistical software and subsequently returned to GRASS.

## 3.5 Hydrological assessment

The hydrologically relevant characteristics of a soil - or a soil horizon - are determined by a complex interplay of texture, bulk density, organic content, and soil skeleton. Lehmann et al. (2005) provided a table relating field capacity $FC$ [mm] and saturated hydraulic conductivity $k_f$ [cm day$^{-1}$] to texture, bulk density and organic content. In order to allow the implementation of these relationships into SOILSIM a conversion into a continuous mathematical relationship was required. This was done using linear multiple regression models.

$$FC = 0.01390 \cdot t_s + 23{,}72 \cdot \rho_s - 2{,}250 \cdot k_C - 50{,}29 \qquad \text{Equation (1),}$$

and

$$k_f = 7.4849 \cdot e^{0.003880 \cdot t_s - 3.275 \cdot \rho_s + 4.466}$$ Equation (2).

$t_s$ represents the texture as average grain size [μm], $r_s$ the bulk density [kg m$^{-3}$], and $k_c$ the organic content [vol.-%]. The R² value for *FC* did not exceed 0.33, but due to lacking alternatives the approach was applied. For $k_f$, R² was 0.77. The total water storing capacity of the soil was calculated from field capacity and soil depth (soil skeleton was subtracted), separated for the A horizon and the remaining soil column.

The permanent wilting point *PWP* [mm] depends primarily on the texture. FREY & LÖSCH (1998) provided values for clayey, silty and sandy soils in form of a diagram. By converting the texture into an average grain size a simple linear regression could be applied to express *PWP*  (R² = 0.99):

$$PWP = -17.72 \cdot t_s + 21.61$$ Equation (3).

The pore volume $V_p$ [%] depends primarily on the bulk density of the soil, and only to a lesser extent on the texture. SCHLICHTING et al. (1995) provide tabular data relating bulk density and pore volume, that was used to fit a simple linear regression (R² = 1.00):

$$V_P = -28.90 \cdot \rho_s + 86.55$$ Equation (4).

*PWP* and $V_p$ were lumped for the entire soil column, without distinguishing different horizons. The whole process is fully integrated into the SOILSIM framework, using the r.mapcalc function of GRASS as the major tool.


# 4    Results

## 4.1  ANOVA

The results of the ANOVA are summarized in Table 2. The soil variables appear to be well explained by the chosen predictors as the variance of each soil variable is successfully explained by at least one, but in most cases two to four predictors. Cumulative temperature, slope angle and substrate are the most powerful predictors for the distribution of soil characteristics.

## 4.2 Regression analysis

Table 3 represents the results of the regression analysis. The results for all soil variables, except the soil skeleton of the A-horizon, are significant at the 0.10 level, but total soil depth and texture are not significant at the 0.05 level of p. The results for all of the soil variables have in common that they scatter considerably around the regression line, resulting in low values for R².

**Table 2:** p-values of the ANOVA of the soil (vertical) and predictor variables (horizontal). Relationships that were used for the regression analysis are written in bold.

|  | TCUM | SLOPE | TIND | CNRAT | STYPE | SPHYS | DIST |
|---|---|---|---|---|---|---|---|
| **total depth** | **0,063** | 0,098 | 0,177 | 0,163 | **0,023** | **0,023** | 0,268 |
| **depth A** | **0,013** | 0,249 | 0,395 | 0,464 | 0,599 | 0,257 | 0,661 |
| **org. content** | **0,001** | 0,935 | 0,007 | 0,341 | **0,006** | **0,000** | 0,400 |
| **texture** | 0,204 | **0,027** | **0,037** | 0,354 | 0,454 | 0,677 | 0,646 |
| **density A** | 0,495 | **0,078** | 0,085 | **0,001** | 0,581 | 0,879 | **0,006** |
| **density B,C** | **0,014** | **0,001** | **0,006** | 0,373 | 0,387 | 0,127 | 0,482 |
| **skeleton A** | 0,477 | 0,504 | 0,967 | 0,599 | **0,007** | 0,187 | 0,376 |
| **skeleton B,C** | 0,793 | **0,006** | 0,958 | 0,710 | **0,023** | **0,022** | 0,844 |

**Table 3:** Predictors used in the regression analysis for the soil variables, values for multiple $R^2$ and levels of significance p.

|  | predictor 1 | predictor 2 | predictor 3 | $R^2$ | p |
|---|---|---|---|---|---|
| **total depth** | STYPE | SPHYS | TCUM | 0,133 | 0,075 |
| **depth A** | TCUM | --- | --- | 0,090 | 0,031 |
| **org. content** | SPHYS | TCUM | STYPE | 0,421 | 0,000 |
| **texture** | SLOPE | TIND | --- | 0,114 | 0,052 |
| **density A** | CNRAT | DIST | SLOPE | 0,270 | 0,002 |
| **density B,C** | SLOPE | TIND | TCUM | 0,280 | 0,001 |
| **skeleton A** | STYPE | --- | --- | 0,013 | 0,417 |
| **skeleton B,C** | SLOPE | SPHYS | STYPE | 0,199 | 0,013 |

As the results provided by the model include a large number of maps, it is only possible to discuss some of them.

The maps representing total soil depth and density of the A horizon are shown in Figure 3. Both of them correspond to the expected patterns. The measured variables were compared with the modelled variables at the coordinates of the study sites. The prediction was satisfactory for the majority of the study sites, but for some sites measurements and model results diverged considerably due to a number of restrictions of the model used (compare discussion).

Figure 4 shows the map of the total water storing capacity of the soil, and the saturated hydraulic conductivity. The meadows and heavily influenced places in the valley showed comparatively low storing capacities and hydraulic conductivities, compared with forest, meadows and pastures at higher elevations. A surprising result was the fact that the highest capacities were related to block deposits. This pattern may be related to the less dense soils at these sites, but it is also related to the insufficient representation of skeleton content and depth of the A horizon by the model (compare $R^2$ and p in Table 3).
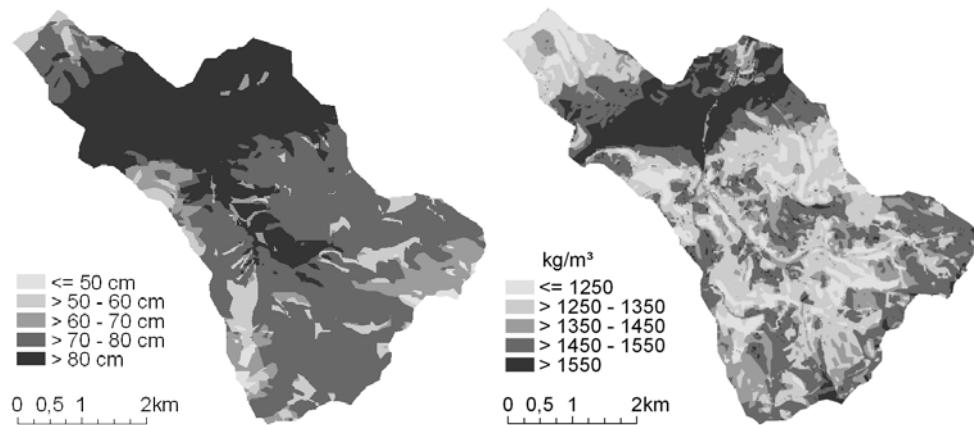
**Figure 3:** Predicted total soil depth (left) and density of the A horizon (right) in the study area.
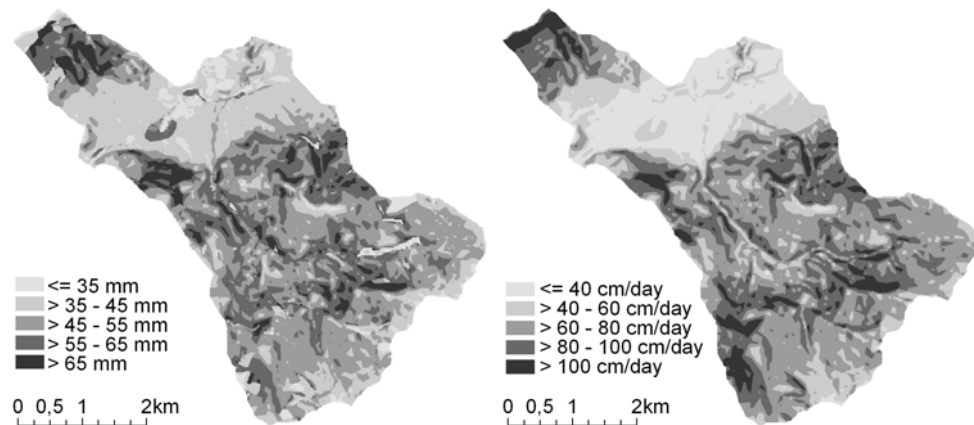


**Figure 4:** Modelled maps representing the water storing capacity of the soil and the saturated soil hydraulic conductivity in the study area.

## 5 Discussion

The application of multiple linear regression for predictive soil modelling in the Stampfanger catchment, as demonstrated using SOILSIM, illustrates the potential of this technique to be integrated into GIS environments in general and into a combination of the open source products GRASS and R in particular. It illustrates, however, also the problems connected to this technique. Although most of the relationships established in this study may

be considered as significant, they show a large amount of scatter, leading to rather unsatisfactory representations of the real conditions. This phenomenon may be related to various reasons:

- Many relationships in nature are not linear, but follow an exponential, polynomial or other nonlinear function. A dataset of 53 samples is small for an alpine catchment (HITZ et al. 2002) and not sufficient to allow the use of nonlinear multiple regression.

- Not all of the soil variables follow a symmetric distribution. Skewness of the variables involved may distort the regression results considerably.

- The transformation of tabular data into regression equations (as attempted for the field capacity) is problematic.

- The data obtained at the study site may be insufficient: total soil depth was only recorded until 100 cm (limited by the length of the Pürckhauer auger). Some other variables were estimated and are therefore susceptible to subjectivity, or they were measured using heuristic field methods in order to save time and money.

- The spatial scale at which soil characteristics undergo changes may be much larger (metres) than the spatial resolution used in the study (tens of metres).

- The most important reason, however, is probably the insufficient knowledge of the predictor variables. While topographic and to some extent also climatological variables were available in sufficient detail, this was not the case for the substrate and particularly for the intensity of anthropogenic (and anthropo-zoogenic) influence. In addition, the consideration of ordinal values as numeric could bias the results as well as the conversion of the texture classes into average grain sizes.

It can be concluded that multiple linear regression shows a certain potential for predictive spatial modelling of soil variables. However, it has to be applied with caution. It is easy to produce nicely looking, colourful maps that even look realistic with this method, but generating reliable, scientifically valuable maps is a different task that requires large datasets of soil variables and carefully prepared predictor maps. Uncertainties, however, will always remain in environmental models (ODEH & MCBRATNEY 2001) as phenomena in nature are not fully predictable.

## 6 Acknowledgements

# 7    References

De Gruijter, J.J., D.J.J. Walvoort & P.F.M. van Gaans (1997): *Continuous soil maps – a fuzzy set approach to bridge the gap between aggregation levels of process and distribution models.* Geoderma 77: 169–195.

Fecht, M., C. Geitner, A. Heller & J. Stötter (2005): *Ermittlung der räumlichen Verteilung von Bodeneigenschaften – eine Kombination von GIS, Clusteranalyse und Geländearbeit.* In: Angewandte Geoinformatik 2005. Beiträge zum 17. AGIT-Symposium Salzburg, 155–164.

Frey, W. & R. Lösch (1998): *Lehrbuch der Geobotanik.*

Hiller, B., A. Müterthies, F.-K. Holtmeier & G. Broll (2002): *Investigations on Spatial Heterogeneity of Humus Forms and Natural Regeneration of Larch (Larix decidua Mill.) and Swiss Stone Pine (Pinus cembra L.) in an Alpine Timberline Ecotone (Upper Engadine, Central Alps, Switzerland).* Geographica Helvetica 57: 81–90.

Hitz, C., M. Egli & P. Fitze (2002): *Determination of sampling volumes needed for representative analysis of alpine soils.* Journal of Plant Nutrition and Soil Science 165: 326–331.

Lehmann, A., S. David & K. Stahr (2005): *Technique of Natural and Anthropogenic Soil Evaluation. Version to be tested with pilot projects within TUSEC-IP.* Hohenheim University, Institute of Soil Science. http://www.tusec-ip.org (May 2005).

McBratney, A.B., M.L. Mendonça Santos & B. Minasny (2003): *On digital soil mapping.* Geoderma 117: 3–52.

Moran, A.P., J. Lammel, C. Geitner, A. Gerik, C. Oberparleiter & G. Meißl (2005): *A conceptual approach for the development of an expert system designed to estimate runoff in small Alpine hydrological catchments.* In: Landschaftsökologie und Umweltforschung 48: 199–210 (Conference Proceedings of the International Conference on Hydrology of Mountain Environments, Berchtesgaden).

Odeh, I.O.A. & A.B. McBratney (2001): *Estimating uncertainty in soil models (Pedometrics '99).* Geoderma 103: 1.

Renger, M., G. Wessolek, B. List & R. Seyfert (1987): *Beziehungen zwischen Bodenfarbe und Humusgehalt.* Mitteilungen der Deutschen Bodenkundl. Gesellschaft 55: 821-826.

Scheffer, F. & P. Schachtschabel (2002): *Lehrbuch der Bodenkunde.* 15th edition. 494 pp. Stuttgart.

Schlichting, E., H.-P. Blume & K. Stahr (1995): *Bodenkundliches Praktikum.* 2nd edition. 295 pp. Berlin.

Tasser, E., B. Ostendorf, U. Tappeiner, C. Pöttinger, W. Bitterlich & A. Cernusca (1998): *Analysis of factors determining soil depth on an Alpine hillslope.* HeadWater'98, 58–62.

UFZ-Umweltforschungszentrum Leipzig-Halle (2001): *Modellierung von Umsatz und Transportprozessen im Boden. CANDY – Simulation der Kohlenstoff- und Stickstoffdynamik. Dokumentation.* http://www.ufz.de/data/dokumentation53.pdf.

Vereecken, H. (1995): *Estimating the unsaturated hydraulic conductivity from theoretical models using simple soil properties.* Geoderma 65, 81–92.

Welpmann, M. (2003): *Bodentemperaturmessungen und –simulationen im Lötschental (Schweizer Alpen).* Dissertation an der Rheinischen Friedrich-Wilhelms-Universität Bonn.