SCALABLE GENERALIZ	
STRATA FOR	
Jānis JĀTNIEKS <sup>1</sup> , Kon	ra
Faculty of Geography and Earth Sciences, Facu	   

#### Workflow overview

MySQL, Python (MySQLdb)	(SciPy, Numpy, BioPython)	Python (Biopython, Numpy)	R (deldir, maptools, sp, gpclib, rgdal, rgeos)	pyMOSYS, Python, OGR, shape	ly pyMOSYS, Python, Scipy		
Data preprocessing	Distance Matrix Calculations	Cluster analysis	Creation of Spatial Clusters	Integration into model geometry	Experiment Calculations	15-	
Load borehole logs from data	Normalized compression	List of cluster membership	Voronoi tesselation	Get model mesh element cluster membership from	Full model run for each of 162 candidate geometries		
base Normalize layer	Sum of range normalized Euclidean and NCD matrices	Sum of range	Export of flattened	Export to GIS 2D GIS clustering layer Shapefile format Cluster typical borehole	Determination of the best		
transitions depths against thickness of Quaternary cover		clustering to CSV Determination of the typical	Dissolve adjacent polygons with identical cluster memebership	column structure generation throughout the territory of the respective modelling mesh elements	Further optimization by using MOSYS autocalibration function	10-	
Normalized borehole log serialization		borehole column for each cluster (centroid by minimum average)		Assignment of conductivity values from lithology classifer		count	
			<b>C</b> 11		<b>^</b> .		1111411

2. Overview of the full experiment workflow for using Normalized Compression Distance and hierarchical clustering as an approch for scalable generalization of Quaternary sediment lithological structure in the MOSYS regional groundwater modelling Fig. 2. Complete-link system for the Baltic artesian basin territory.

# Introduction

The cover of Quaternary sediments, especially in formerly glaciated territories usually is the most complex part of the sedimentary sequences. In regional hydro-geological models it is often assumed as a single layer with uniform or calibrated properties (Valner, 2003). However, the properties and structure of Quaternary sediments control the groundwater recharge: it can either direct the groundwater flow horizontally towards discharge in topographic lows or vertically, recharging groundwater in the bedrock. Building a representative structure of Quaternary strata in a regional hydrogeological model is important as this affects all the underlying strata.

## Scope and objective

This work aims to present calibration results and detail our experience while integrating a scalable generalization hydraulic conductivity for Quaternary strata into the regional groundwater modelling system for the Baltic artesian basin – MOSYS V1.

## Methods

In this study the main unit of generalization is the spatial cluster and the main criteria for finding adequate representative borehole columns is the Normalized Copression Distance metric – a nonparametric datamining technique calculated by CompLearn ncd utility (Cilibrasi, 2007).

Spatial clusters were obtained from flattened completelink hierarchical clustering dendrogram representations, supplying cluster membership identifiers for each borehole. Using Voronoi tesselation and GIS dissolve spatial analysis functionality on this data spatial clusters were obtained in GIS format. The procedure is described in detail in "Lithological Uncertainty Expressed by Normalized Compression Distance" by Jatnieks et al. 2012. Full overview of the experiment workflow with software components and key stages is shown in **Figure 1.** 

We used two experiment matrices for hierarchical the layer structure into the model geometry, the clustering. One where spatial clusters are obtained from respective conductivity values are assigned to mesh Normalized Compression Distance (further identified as elements, created by the new placeholder layers. The NCD matrix) matrix alone and a second where matrix is cluster membership for each of mesh elements i combined with range normalized Euclidean distance matrix determined by intersecting the element centroids, shown of borehole locations in BalticTM93 projection (further here in white with spatial clusters in GIS Shapefile layer. identified as **NCD+E** experiment matrix).





Fig. 3. Algorithm for calculating Z values for including spatial clustering results into the MOSYS finite-element mesh geometry and resolving structure conflicts, arising from different lithological structures, geological columns selected from the most typical borehole log in this cluster.



The performance of resulting geometries is detailed in **Figures 8-11.** For both matrix variants Compression Distance, on the right. implementation, using hierarhical clustering of Normalized Fig. 4. After the algorithm described in Fig. 3 creates (NCD and NCD+E) 82 model geometries, based on different flattenings of cluster dendrogram (**Figure** This provided insight in their sensitivity to calibration 2) were prepared and run in MOSYS V1 modelling through further optimization of hydraulic conductivity values environment. From the inital model run results in in Quaternary strata (Figure 12). The model geometry, in Figures 8 and 10, 3 better perfoming, one which the Quaternary clustering experiment geometries were median and also the worst performing geometry integrated, had already been calibrated until convergence was run through MOSYS autocalibration routine before these experiements. Log scale for conductivity values on Y axis. (Klints et al. 2012).

# ATION OF HYDRAULIC CONDUCTIVITY IN QUATERNARY **USE IN A REGIONAL GROUNDWATER MODEL**

rāds POPOVS<sup>1</sup>, Ilze KLINTS<sup>2</sup>, Andrejs TIMUHINS<sup>3</sup>, Andis KALVĀNS<sup>1</sup>, Aija DĒLIŅA<sup>1</sup>, Tomas SAKS<sup>1</sup> Ity of Physics and Mathematics, <sup>3</sup>Laboratory for Mathematical Modelling of Environmental and Technological Processes. University of Latvia, Riga, Latvia. Contact: janis.jatnieks@lu.lv, www.puma.lu.lv



agglomerative hierarchical clustering dendrogram for Normalized Compression Distance matrix of serialized borehole log lithological

This degree of similarity and the spatial heterogeneity of the cluster polygons can be varied by different flattening of the hierarchical cluster model into variable number of clusters. Such an approach provides a scalable generalization solution. It can be scaled from broad to more detailed generalization, depending on model and structural characteristics of the modelling territory,  $_{2}$  with the help of the clustering dendrogram (Figure  $\degree$ 

Using the dissimilarity matrix of the NCD metric, a borehole, most similar to all the others from the ilithological structure point of view, can be identified from the subset of all cluster member boreholes The log structure of this borehole then is applied throughout the territory of the corresponding spati cluster. The spatial locations of the same logica clusters can be disjoint - in different parts of modelling territory and have the same borehole log column, representing the lithological structure

throughout spatial clusters with the same logical<sup>22</sup> cluster identifier.

162 model geometries were prepared, incorporating different candidate representations of Quaternary strata made from spatial clusters with the most typical selected borehole column, The NCD metric. using borehole log column lithological structure aggregate composed using lithology types as shown in Table conductivity values.

these integrate results into the geometry groundwater algorithm, shown in Figure calculates mesh point heights in meters above sea level for every mesh point in created newly every Quaternary layer.

After the layer surface calculations, the mesh element cluster memebership is determined by intersecting borehole element centroids (white dots in Figure 4) with the spatial clusters.



conductivity values lithology used generalization structures for studv.



Fig. 5. Spatial distribution of vertical permeability values in the MOSYS V1 modelling system mesh in the clustering experiment territory in Latvia. Screenshot fragment of MOSYS visualization platform HiFiGeo, log scale, EPSG:25884 projection.



Fig. 6. Middle fragment of vertical crossection AB (overview in Figure 5) across the Vidzeme region, showing horizontal permeability values on the right Y axis in Quaternary represented by 8 modelling system layers (only Q layers shown, for simplicity). Left Y axis - Z values in meters above sea level. Distance from crossection A point on X axis.



log **Fig. 7.** Previous implementation of Quaternary strata in the MOSYS V this modelling system consisted of 4 layers - two sets of aquifer-aquitard transitions. The arrow indicates the border between the previous Quaternary implementation (left) and the new Quaternary





The best performing model geometry before applying the NCD metric based clustering solution, had squared error sum of 814 meters. The best NCD clustering geometry provides 9.7% **improvement** to 735. This may seem modest in nominal terms. It isn't, because the experimental Quaternary strata were inserted in a previously calibrated model geometry with initial conductivity values as estimates. Model autocalibration works by multiplying layer material properties with coefficients determined by the the Scipy L-BFGS-B routine (Zhu et al. 1997). Since this works on a **Fig. 11.** Squared error distribution by model calibration layers for layer basis, the proportional differences in conductivity values each of the geometries derived from different logical clustering between elements in each layer remain constant. Ideally, the solutions using NCD+E matrix. model geometry area built using clustering and the rest of the The response of the selected new model geometries to territory would be calibrated seperately. This could yield further improvements.

calibration indicates their influence on the underlying model strata, not just improved performance in the part of the model representing the Quaternary deposits.

Geometries NCD 12c, NCD 13c were calibrated further, including the remaining bedrock layers as optimization parameters and yielding further decrease of squared error sum, used as overall fitness function for model performance against observed pressure heads.

shown in **Figures 6 and 7.** 





This work was supported by the European Social Fund project "Establishment of interdisciplinary scientist group and modelling system for ground-water research", 2009/0212/1DP/1.1.1.2.0/09/APIA/VIAA/060

Examples of crossections from NCD 12c structure are

#### Results

The experiment geometries based on NCD matrix perform better overall.

The various Quaternary representations tested have a measurable effect on model results - the worst geometries, such as NCD+E 7, 3 and others with low logical cluster count, perform more than twice as bad compared to the best NCD clustering based geometries.

NCD clustering based geometries with the better initial calculation results (in Figures 8 and 10), such as the 12 and 13 cluster variants, also tend to respond better to optimization with MOSYS auto-calibration routine (Figure 12).

Some the geometries with lower results from intial model run, show very minimal improvement with calibration (Figure 12).

For our experiment territory, consisting of formerly glaciated territory of Latvia, the **best performing generalizations have** 12 and 13 clusters.

This could have been successfully deduced from the clustering dendrogram in Figure 2.

- 1. Bennett, C. H., Gacs P., Li M., Vitanvi P., Zurek W. 1998, Information Distance, IEEE Transactions on Information Theory,
- 44(4), 1407-1423, IFFF 2. Cilibrasi, R. 2007., Statistical Inference Through Data Compression, ILLC Dissertation Series DS-2007-01, Institute for
- Logic, Language and Computation, Universiteit van Amsterdam. 3. Klints I., Virbulis J., Timuhins A., Sennikovs J. and Bethers U., Calibration of the hydrogeological model of the Baltic Artesian Basin, EGU Geophysical research abstracts EGU2012-10003 [this section poster board A255]
- 4. Takčidi. E. 1999. Datu bāzes "Urbumi" dokumentācija [Documentation of the database "Boreholes"]. Valsts ģeoloģijas dienests, Rīga, [In Latvian]. 5. Valner, L. 2003. Hydrogeological model of Estonia and its applications. Proc. Estonian Acad. Sci. Geol., 52, 3, 179-192.
- 6. Virbulis, J., Bethers, U., Saks, T., Sennikovs, J. & Timuhins, A. Hydrogeological model of the Baltic Artesian Basin Hydrogeology Journal, [Under revision

C. Zhu, R. H. Bvrd and J. Nocedal, L-BFGS-B: Algorithm 778: L-BFGS-B, FORTRAN routines for large scale bound constrained optimization (1997). ACM Transactions on Mathematical Software, 23, 4, pp. 550 - 560.