# TERMINOLOGY AND METHODOLOGY IN MODELLING FOR WATER QUALITY MANAGEMENT – A DISCUSSION STARTER

Jacob Carstensen*, Peter Vanrolleghem**, Wolfgang Rauch*** and Peter Reichert†

\* *Department of Environmental Science and Engineering, Building 115, Technical University of Denmark, DK-2800 Lyngby, Denmark*
\*\* *BIOMATH, Department of Applied Mathematics, Biometrics and Process Control, University Gent, Coupure Links 653, B-9000 Gent, Belgium*
\*\*\* *Institute of Environmental Engineering, Universität Innsbruck, Technikerstrasse 13, A-6020 Innsbruck, Austria*
† *Swiss Federal Institute for Environmental Science and Technology (EAWAG), CH-8600 Dübendorf, Switzerland*

## ABSTRACT

There is a widespread need for a common terminology in modelling for water quality management. This paper points out sources of confusion in the communication between researchers due to misuse of existing terminology or use of unclear terminology. The paper attempts to clarify the context of the most widely used terms for characterising models and within the process of model building. It is essential to the ever growing society of researchers within water quality management, that communication is eased by establishing a common terminology. This should not be done by giving broader definitions of the terms, but by stressing the use of a stringent terminology. Therefore, the goal of the paper is to advocate the use of such a well defined and clear terminology. © 1997 IAWQ. Published by Elsevier Science Ltd

## KEYWORDS

Identification; modelling; model constituents; model attributes; model applications; model characterisation.

## INTRODUCTION

When the Tower of Babel was destroyed, man was forced to speak different languages. Today, fortunately, researchers all over the world are using the same language - English - for spreading their messages. This has eased communication a lot, but when it comes to understanding, researchers still appear to be using different languages. This is due to a *lack of common terminology.*

Within the field of water quality management, researchers have very different scientific backgrounds. Mathematicians, statisticians, microbiologists, environmental engineers, civil engineers, control engineers, ecologists, biologists, etc. need to collaborate. This scientific multi-cultural society has the potential for solving many of the problems threatening the environment. It should be stressed that scientific advances

often occur when different scientific fields interface and knowledge is transferred from one field of science to another. This transfer of knowledge is facilitated by establishing a common terminology for reference. Thus, the goal of this paper is to enable communication within the field of water quality management through clarifying the use of modelling terms.

Mathematicians are known to have strict use of terms, giving a definition each time a new term is introduced. As a result the communication between mathematicians is very stringent. However, many of the other scientific fields have had an unwillingness to use the general and stringent mathematical language and have simply developed their own more loose and less stringent scientific language. The aim of this paper is not to discard all other scientific terminologies but the purely mathematical, and make this the common reference of terms. This would be similar to forcing people to speak Esperanto, which is structurally and grammatically a more stringent language than English. In contrast, the emphasis of the present paper is to establish a common modelling terminology which is sufficiently stringent and easy to understand for the diversity of scientific fields within water quality management. The terminology presented in the paper is not an exhaustive list of modelling terms and descriptions, but a sample of the most widely used terms for characterising models and model building.

Problem statement

There are two main sources of scientific communication problems: (1) misuse of existing terminology and; (2) use of unclear terminology. As the scientific society constantly grows, these communication problems will become more obvious due to lack of common and widespread terminology. The misuse of existing terminology is the most serious communication problem, because it brings confusion to the recipient of the message as well as confusion to the existing terminology. This is potentially dangerous for the scientific society as it brings disorder into the established terminology. The misuse of existing terminology is either due to lack of knowledge or carelessness with the use of modelling terms. The other source of confusion is the use of terms which do not have a clear definition yet or the use of terms that do have a clear definition in a confusing way. For instance, the term a lumped model is frequently used without reference to which type of model it is lumped from. These sources of communication problems are exemplified below.

*Example 1.* A frequently observed confusion in terminology is the use of model versus program (e.g. "the HYDRODIF model" in Dakhlaoui et al., 1995, "the EFOR simulation model" in Pedersen et al., 1992 or "application of MOUSE model" in Ghafouri, 1996). Since the `Activated Sludge Model No.1' was published (Henze et al., 1987), several commercial software programs implementing the ASM No.1 with some additional features have been launched. Today, many computer programs are available for water quality management. With the misuse of terminology as given above, insight into the computer programs is becoming a prerequisite. As the number of software programs for water quality management constantly grows, the terminology confusion worsens.

*Example 2.* In the literature many terms are used for characterising the basis of a model structure. On one hand words like physical[1], mechanistic or white-box are used as model attributes to describe the model's relationship to physical, chemical and biological laws, while words like empirical, phenomenological or black-box are used as model attributes to describe the model's relationship to heuristic methods. However, no model is purely mechanistic or phenomenological in the sense that mechanistic models have elements of phenomenology and phenomenological models have elements of underlying mechanisms. It is also obvious that a mechanistic model may incorporate different specification levels with respect to the number of physical, chemical and biological laws, and that the model attributes described above have a weak or unclear definition.

As stated previously, the goal of this paper is to clarify the use of modelling terms. For this reason a glossary of the terms used in this paper is suggested after the references. The terms in the glossary are typed in

---

[1] A **physical** model also describes a small-scale model in physical materials, which is used for scale-down experiments. This is a non-mathematical model.

boldface in the text. The glossary list only contains some of the most common modelling terms, and is by no means meant to be exhaustive (see e.g. Singh (1987) or Young (1993)). Two main topics are discussed in the following sections of the paper - terminology and methodology of modelling.

## MODELS AND THEIR APPLICATIONS

Scientific research makes the attempt to describe the mechanisms responsible for the observed behaviour of a system. The abstract representation of a real system by the ideas on its constituents and functional relationships is called a **conceptual model**, or, if they are formulated mathematically, a mathematical model - in short just a model. Mathematical formulation of these ideas leads to a mathematical model that can be used to give quantitative answers to questions about its behaviour under given external conditions. Such mathematical models are referred to simply as "models" in this paper. Because environmental systems are much too complicated to be described in detail, models used for water quality management must be drastically simplified descriptions of reality. Since the aspect of a system that is relevant depends on the question to be answered, a unique model for an environmental system does not exist, but different models must be used for different purposes, and even in a given context several adequate descriptions are possible.

### Model constituents

The discussion of the various definitions and the terminology of model constituents can never emerge from the universe but is inevitably in the context of the general structure of models. Ultimately, we think of a model as a "machine" that transforms inputs (u) to outputs (y) by defined relations (Casti, 1977; Chui and Chen, 1988), where u and y are, when discretised, sequences of either scalars or vectors. (In continuous time formulation inputs u may be also a vector of **forcing functions** from outside forces.) The features of the input-output relations determine the basic structure type of the model, which is either an input/output or a state-space description (Casti, 1977). The inputs of a model consist of disturbances and manipulated variables.
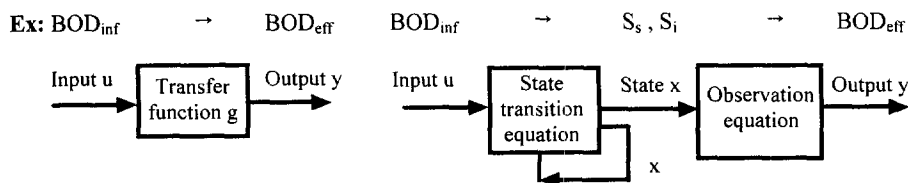


Figure 1. Pictorial description of left: input/output model and right: state-space model. An example of modelling the degradation of carboneous matter is given above, where BOD in the effluent is modelled directly as a function of BOD in the influent or through the internal state variables $S_s$ and $S_i$ denoting fractions of BOD.

An **input/output model** is, strictly speaking, only the set of **transfer functions** (g) that transform the inputs u directly to outputs y. Although this type of description of system behavior (Casti, 1977 denotes it as an external system description) is useful and sufficient in various circumstances, it omits the consideration of the mechanism by which inputs are related to outputs. In order to overcome that deficiency the state-space approach was introduced in the early 1960s (Kalman, 1960).

The most important feature of a **state-space model** is the introduction of **state variables** (vector x) which act as mediators between the inputs and the outputs. These state variables are additional model constituents and the system description is consequently addressed as internal. It is characterised by the fact that x is obtained from present and past x and u by means of the **state-transition equation** and y is generated from x by means of the **observation equation**. The **state** of the system (as described by the model) is defined as the values of the state variables at any instant of time. Note that by definition it is neither required that state variables are measurable nor that they are meaningful in terms of natural science (although they frequently

are the latter). It should also be stressed that the equations of a model can be formulated as either algebraic or differential equations.

Irrespective of the **model structure**, the mathematical equations that relate inputs to outputs contain three types of constituents: **variables, constants** and **parameters**. Inputs, outputs and states are seen as variables in the equations. The difference between constants and parameters is less evident and gave rise to some confusion in modelling terminology. In the following we denote all model constituents that never change their value throughout all possible applications of the model as constants. A parameter, on the other hand, is a model constituent whose value varies with the circumstance of its application. Hence the value needs to be determined for each particular application of the model. The value of some of those parameters may even need to be modified during a specific application, e.g. adaptive control. The current value of those parameters is related to time and location. For some applications a parameter can be replaced by a function describing its time and space dependency but in fact that should be seen as a model extension. Henze *et al.* (1987, 1995) incorrectly use the term constant for parameter in the widely used IAWQ Activated Sludge Model No. 1 and 2 as they are application dependent.

Model attributes

The model constituents describe the fundamental elements of a model, while the characteristics of a model can be addressed by a number of descriptive modelling terms. These terms are called model attributes. Some attributes have a clear and stringent definition (e.g. linear versus non-linear), while others have a weak or relative definition (e.g. phenomenological versus mechanistic). In the following the model attributes are referred to as being strong or weak depending on the stringency of their definition. It is obvious that the more strong model attributes are applied, the better a model is characterised.

*Strong model attributes.* The model attributes **linear** and **non-linear** relate to the structure of the model equations. The model may be linear with respect to the variables or to the parameters (only considered in statistics). Unless otherwise stated, the term linear relates to linearity in the variables. Thus, a model can be non-linear in the parameters but linear in the variables and vice versa. Linearity is a basic characteristic of a model that has quite some impact on the properties of solution, e.g. linear models are frequently used, because the analytical solution can be found. For non-linear models numerical solutions are predominant.

In water quality management models are often characterised as **dynamic** in the sense that the variables evolve over time. A model which is not **dynamic** is called **static** or **steady-state**. Thus, dynamic relates to a time dependency in the model which can be formulated as dynamic input variables and/or state variables. The output of a dynamic model is often called time series (e.g. Ljung, 1987). If the model parameters are constant in time the model is characterised as being **time-invariant**. Similarly, a model can have a space-dependency (e.g. a clarifier model). Such models are referred to as **distributed parameter** models. In water quality modelling only time- and space-derivatives are concerned, and dynamic distributed models are normally formulated as partial differential equations.

Models can be termed discrete or continuous, and in most cases these attributes relate to the model formulation of difference/differential equations with respect to time, i.e. the correct terminology should be **discrete-time** or **continuous-time** model. However, **discrete-space** or **continuous-space** are two other model attributes describing a space relationship for up to three dimensions in the model formulation. A continuous-time differential equation is either solved analytically or discretised into a discrete-time difference equation which is solved numerically. It should be stressed that most computer programs apply a discretisation to the continuous-time and continuous-space differential equations, and the discretised equations are solved as algebraic equations. Another fact is that data can only be observed at discrete time instants (Tong, 1990).

If a model contains elements of randomness, it is called **stochastic** otherwise **deterministic**. The uncertainty encapsulated in any model is due to a combination of: 1) uncertainty in input variables; 2) uncertainty in parameter values; and 3) uncertainty in model structure (O'Neill and Gardner, 1979; Beck, 1983,1987).

When all uncertainty aspects are neglected, the model is deterministic and the output is determined uniquely by input and initial conditions. The output of a stochastic model can be described as a probability density function. The terms stochastic and statistical are occasionally confused, and the use of statistical as a model attribute should be avoided, since the term statistical is referring to methods of analysing data. Likewise, the term deterministic is occasionally confused with mechanistic, physical or white-box implying that deterministic models are always based on physical, chemical and biological laws. This is not true, however, because a black-box model may also be deterministic, e.g. a neural network or a spline. Stochastic input to a model is denoted as innovations, realisations, disturbances or pertubations. For on-line processing of data the term **adaptive** is applied to models, meaning that there is a feedback from measurements to the model. This pertubation of the model is a **stochastic** input. Hence, an adaptive model is stochastic with a feedback from on-line measurements.

The strong model attributes are highly recommended for characterising models as the meaning of the terms is well defined. However, the use of additional adjectives such as purely, totally or completely (e.g. 'purely deterministic nor purely stochastic model' in Harremoës and Carstensen, 1994) has no meaning when the stringent definitions of the strong model attributes are referred to - in fact, these combinations should be avoided as they are confusing. The majority of models for water quality management so far are formulated as non-linear dynamic continuous-time deterministic models (Rose, 1987; Jørgensen, 1992).

*Weak model attributes.* These terms have less clarity in their interpretation and may in the lack of strong model attributes potentially lead to confusion. However, provided that the terms are used correctly these attributes also signify the background of the model. To a large extent many of the weak model attributes describe almost the same model property, i.e. the degree of conceptualism, basis in physical, chemical and biological laws, simplification level, etc.

As described in the introduction, words like **mechanistic, physical** and **white-box** are used to describe that the model's structure is based on physical, chemical and biological laws. The attribute **transparent** has the same interpretation as **white-box**, which means that every detail in the model has a mechanistic explanation. The extreme is **reductionist** models (not to be confused with reduced order models) that are based on the attempt to include as many details as possible. The term causal is also used with the same meaning as mechanistic, but in some scientific fields a causal model is strictly defined as a model which only depends on past observations. Thus, the use of causal as a model attribute should be avoided. **Phenomenological, black-box, empirical** and **heuristic** (by rules of thumb) are used as model attributes to describe that the model is based on empiricism rather than laws. A black-box model has not necessarily a structure compatible with the underlying physical, chemical or biological reality (Tulleken, 1992). A combination of the mechanistic and phenomenological approach is frequently called **grey-box** modelling. Holst *et al.* (1992) refer to grey-box models as reflecting a priori knowledge as well as black-box parts, while in Carstensen *et al.* (1996), grey-box models are given two virtues - the properties of parameter interpretability and parameter identifiability. Ljung (1987) refers to this approach as **semi-physical** modelling.

Another issue for characterising models is the degree of simplicity or complexity in the model. A **simple** model is characterised by few equations and parameters while a **complex** model has many equations and parameters, but it remains unclear when a given model should be termed simple or complex. As a rule of thumb, mechanistic and phenomenological models are normally formulated with a high and low degree of complexity, respectively. However, an artificial neural network is considered to be a black-box model but may at the same time have a high level of complexity (Hertz *et al.*, 1991). The terms simple or complex for characterising a model indicates that the model is derived from a basis model or compared to another model. The confusion occurs when this reference model is not given or is just assumed to be well known. The level of model complexity/simplicity can also be addressed with the attributes **aggregated/segregated**. Jeppsson and Olsson (1993) use a **reduced order** model to describe the model's derivation as being **lumped** or aggregated. A model is lumped or aggregated when model variables or equations are united in a simplified description. Even though it is obvious to some people that a model is complex, based on the state-of-the-art in modelling today, this fact is very likely to be unclear to researchers unfamiliar with the specific topic, or researchers within the field 10 years from now.

## Model applications

Mankind is using models, in nearly every aspect of human life, as the principal vehicle to describe reality. The manifold applications of models can be categorised according to the modelling objective into the areas understanding and prediction. Models applied in understanding aim at the increase of knowledge of system behaviour. The objective is to develop a simple, yet universal model of the system under consideration that gives an adequate description of reality as it was observed (Reichert, 1994). The use of models for the purpose of understanding is most frequent in research and education. The prediction of either future or hypothetical system behavior is one of the basic tasks in practice. Models applied for prediction aim at providing an accurate and fast image of real systems behaviour under different conditions than those prevailing during model building. The model use can either aim at forecasting future states of the system (simulation with new inputs) or at predicting system behavior under hypothetical scenarios (simulation with new parameter values). The latter application is most frequent in design.

Based on this stringent terminology all possible applications of models can be categorised according to the basic modelling objective. In water quality management models are typically used for analysis, design, control and decision support. The goals of models are understanding for analysis and prediction for design. This relation is less apparent when considering control and decision support as model application areas, e.g. a model used as a decision support tool is aiming at identifying constraints and objectives of a problem in the initial phase (purpose of understanding) and later at choosing the most suitable solution from a set of alternative techniques (by predicting the effect of those techniques). Also in the area of control, models are used for understanding: a software sensor aims at extracting the utmost information from scarce measured data by means of a model. Hence, the implemented model aims at increasing our knowledge about the system under consideration. On the other hand, models in control are also used to predict future states of the system when applying certain control strategies (model predictive control).
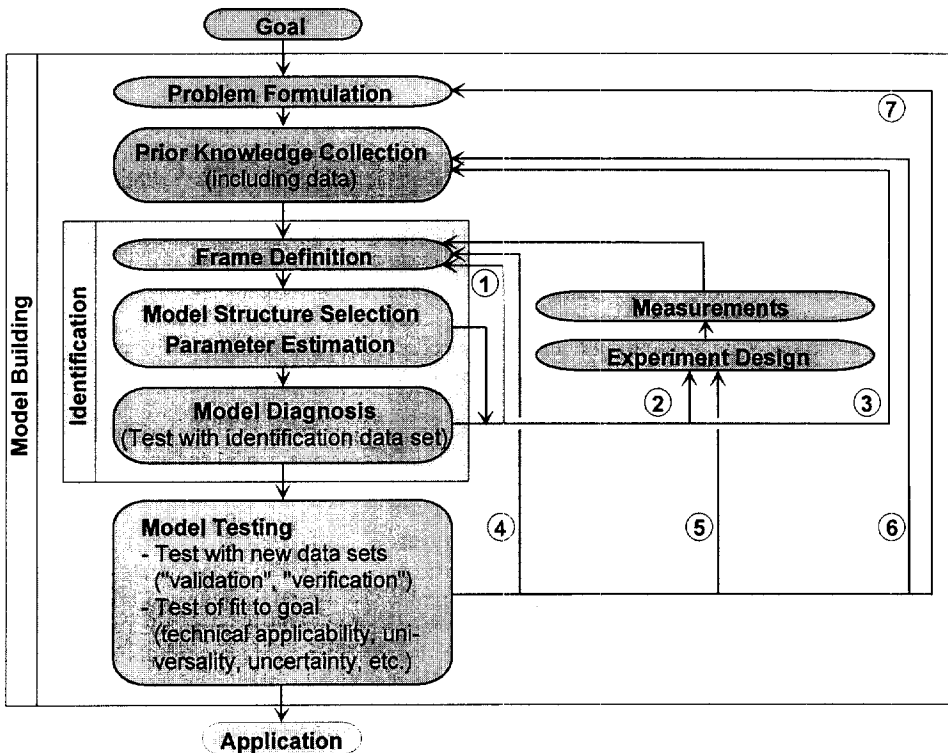


Figure 2. The model building exercise.

## MODEL BUILDING

Using the "story" of a model building exercise, a number of terms involved in this activity will be introduced within their appropriate context. At the same time a short review is given on the current state-of-the-art in modelling methodology. The diagram in Fig. 2 summarises the aspects of model building which are described in detail below. Once the steps in the figure have been fulfilled succesfully, the model can be applied for its intended purpose. These applications typically involve **simulation** that may be regarded as virtual experimentation with the virtual reality of the model.

*Problem formulation.* An often forgotten task in a model building exercise is the clear formulation of the goal of the model that is to be constructed. While in most cases this task is rather intuitively performed by the modeller when he is also the problem supplier, problem formulation or goal incorporation is much more difficult when these persons are not the same. In this case, an important effort must be made to answer such questions as accuracy of the results, degree of uncertainty in the provided answers, time scale of the solution, system boundaries, important variables, environmental conditions under which questions will be asked for which the model must give an answer, etc.

*Prior knowledge collection.* The next task is to collect the available, relevant, a priori knowledge from literature and experts or from a model building environment that supports re-use of model-encapsulated knowledge. At this (early) stage of the exercise, some experiments may be conducted or some data collected during previous experiments may be retrieved and stored in the experimental database.

*Frame definition.* As soon as these two tasks have been performed once, a first iteration of the model building can start. The **frame definition** phase aims to delineate the conditions under which a model will be used (e.g. temperature), to choose the class of models that seems fit for the task (time series, state-space, distributed parameter, stochastic...), to specify the variables that seem important to find a solution to the formulated problem (inputs, outputs, states), the range of time constants that need to be covered by the model, etc.

*Model structure selection.* With this frame of work defined, candidate models may be constructed combining the collected a priori knowledge and the creativity of the model builder. The goal of **classical model structure selection (model structure characterisation)** is to select a unique model structure according to the principles of a quality of fit and parsimony (Spriet, 1985; Harvey, 1989). However, it is also possible to select a set of models that are attributed different weights reflecting their probability of appropriateness (Draper, 1995; Reichert and Omlin, 1997). Most model structure selection criteria require the parameter values to be estimated, e.g information criteria such as AIC (Akaike, 1974) and BIC (Schwartz, 1978). However, structural selection criteria that only require basic data analysis also exist for particular applications (Vanrolleghem and Van Daele, 1994).

*Parameter estimation.* **Parameter estimation** is based on the maximisation or minimisation of a goodness-of-fit criterion such as Least Squares, Weighted Least Squares, Maximum Likelihood, etc. and involve providing values for the parameters in the model and, in some cases, also values for the initial (and boundary) conditions of the state variables (in case a state-space representation is adopted). Although several powerful estimation (non-linear **regression**) algorithms are available (both for **recursive** and **batch estimation** depending on the number of data points used for a parameter update), their success is highly dependent on the experimental dataset available. The identifiability analysis performed prior to the parameter estimation itself can provide answers to the key question whether, given a set of measured variables, unique parameter values can be obtained. Two types of answers can be given depending on the applied method. When **structural** (also termed **theoretical** or **a priori**) **identifiability** (Godfrey and DiStefano, 1987; Norton, 1986) is evaluated the answer is either yes or no, respectively, meaning that the parameters can be given unique values or not at all (Dochain *et al.*, 1995). However, it is not ensured that the data always contain sufficient information to provide reliable and unique estimates (e.g. the model $Y=aX+b$ is structurally identifiable if measurements for Y are available, but it is practically unidentifiable if the dataset contains only values of Y for one particular X-value). Methods for the **practical** (also termed as

**numerical** or **posteriori**) **identifiability** study are available and allow one to evaluate the information content of the dataset intended for parameter estimation (Vanrolleghem *et al.*, 1995). The basis of these methods is also underlying methods that can provide a solution to a practical identifiability problem, i.e. **optimal experimental design.** This design procedure uses the model for which unique parameters are to be found to calculate experimental conditions such that sufficient information is contained in the data. One should note that a structural identifiability problem encountered cannot be solved without altering the candidate model or frame definition (e.g. include other variables in the system description). **Model reduction** can lead to models that become less "data-hungry" and hence their identifiability properties may improve (Jeppsson and Olsson, 1993).

*Model diagnosis.* Once the parameters are estimated it remains to investigate whether the identified model violates the assumptions made in the frame definition. For instance, statistical tests of systematic deviations between model results and measurements (residuals) and their distribution are frequently used (Box and Jenkins, 1976; Söderström and Stoica, 1989).

*Model testing.* Fitness of a model can be evaluated by comparing its performance with data obtained under different conditions than the ones prevailing at the time of the data collection performed for model identification. This process of putting the model in jeopardy (Boyle and Berthouex, 1974) or, in other words, straining the model to its limits, may reveal model inadequacies that may be sufficient to conclude that the model is no longer "fit" for the purpose it was intended for. Hence, the whole model building process may have to start all over. Sometimes this may even lead to a reformulation of the problem as the modelling exercise has provided considerable insight into the system under study and its behaviour. This process of putting the model into jeopardy by confronting it with new data is most often called **validation**, but serious arguments are put forward against this term. As a model only describes part of the reality (the one defined in the frame) in a simplified manner, it is obvious that a model can never describe reality completely. Therefore, there will always exist experimental conditions for which the model is not valid. Hence, validation of a model is utopian! A completely other approach is to term this process of jeopardising the model a model **falsification** step (Caswell, 1976; Reckhow and Chapra, 1983), which if answered negatively, provides more confidence in the selected model. However, the term falsification appears too negative and one has therefore looked for other terminology that is less pronounced (quantitative) as validation but still gives a qualitative insight into the level of confidence one has in a selected identified model. The terms put forth for this are **corroboration** and **confirmation** (Popper, 1980). Finally, the term **verification** is frequently interchanged with validation, but the use of the term validation is advocated here.

## SUMMARY

Communication in modelling for water quality management is often confused due to lack of common terminology and lack of clarity in the terminology used. This paper stresses the need for a common reference of terminology. A stringent terminology with the focus on model constituents, model attributes and model building is suggested. Definitions of the terms used in this paper are given in the glossary. In Henze *et al.* (1982) a reference for notation in biological wastewater treatment was given and within 5 years became standard for communication. The communication problem to be dealt with in the future is that of common terminology.

## ACKNOWLEDGEMENT

## REFERENCES

Akaike, H. (1974). A New Look at the Statistical Model Identification, *IEEE Transactions on Automatic Control,* 19(6), 716-723.

Beck, M. B. (1983). Uncertainty, System Identification and the Prediction of Water Quality. In: *Uncertainty and Forecasting of Water Quality*, Beck, M. B. and van Straten, G. (eds), Springer, Berlin, pp. 3-68.

Beck, M. B. (1987). Water Quality Modeling: A Review of the Analysis of Uncertainty. *Water Resources Research*, **23**(8), 1393-1442.

Box, G. and Jenkins, G. (1976). *Time Series Analysis, Forecasting and Control*, Holden-Day, 575 pp.

Boyle, W. C. and Berthouex, P. M. (1974). Biological Wastewater Treatment Model Building: Fits and Misfits. *Biotechnological Bioengineering*, **16**, 1139-1159.

Carstensen, J., Harremoës, P. and Strube, R. (1996). Software Sensors based on the Grey-box Modelling Approach, *Wat. Sci. Tech.*, **33**(1), 117-126.

Casti, J. L. (1977). *Dynamical systems and their application - linear theory*. Academic Press, New York.

Caswell, H. (1976). The Validation Problem. In: *Systems Analysis and Simulation in Ecology, Volume IV*, Patten, B.C. (ed), Academic Press, New York, pp. 313-325.

Chui, C. K. and Chen, G. (1988). *Linear systems and optimal control*. Springer-Verlag, Berlin.

Dakhlaoui, M., Gaujois, D., Raimbault, G. and Tabuchi, J.-P. (1995). Water Diffusion Device Design in Reservoir Structures. *Wat. Sci. Tech.*, **32**(1), 71-78.

Dochain, D., Vanrolleghem, P. A. and Van Daele, M. (1995). Structural Identifiability of Biokinetic Models of Activated Sludge Respiration. *Water Research*, **29**, 2571-2579.

Draper, D. (1995). Assessment and Propagation of Model Uncertainty (with discussion). *Journal of Royal Statistical Society B*, **57**(1), 45-97.

Ghafouri, R. A. (1996). Application of MOUSE model on Impervious Area Runoff Simulation in Australia. *Preprints of the 7th International Conference on Urban Storm Drainage*, University of Hannover, pp. 437-442.

Godfrey, K. R. and DiStefano, J. J. (1987). Identifiability of Model Parameters, Chapter 1 of *Identifiability of Parametric Models*, Walters, E. (ed), Pergamon Press, Oxford, pp. 1-20.

Harremoës, P. and Carstensen, J. (1994). Deterministic versus Stochastic Interpretation of Continuously Monitored Sewer Systems. *European Water Pollution Control*, **4**(5), 42-48.

Harvey, A. (1989). *Forecasting, Structural Time Series Models and the Kalman Filter*, Cambridge University Press, 554 pp.

Henze, M., Sutton, P. M., Gujer, W., Koller, J., Grau, P., Elmaleh, S. and Grady, C. P. L. (1982). The Use and Abuse of Notation in Biological Wastewater Treatment. *Water Research*, **16**, 755-757.

Henze, M., Grady, C. P. L., Jr., Gujer, W., Marais, G. v. R. and Matsuo, T. (1986). *Activated Sludge Model No. 1*, IAWQ Scientific and Technical Report No. 1, IAWQ, London.

Henze, M., Gujer, W., Mino, T., Matsuo, T., Wentzel, M. C. and Marais, G. v. R. (1995). *Activated Sludge Model No. 2*, IAWQ Scientific and Technical Report No. 3, IAWQ, London.

Hertz, H., Krogh, A. and Palmer, R. (1991). *Introduction to the Theory of Neural Computation*. Addison-Wesley.

Holst, J., Holst. U., Madsen, H. and Melgaard, H. (1992). Validation of Grey Box Models, *IFAC Symposium on Adaptive Control and Signal Processing*, pp. 407-414.

Jeppsson, U. and Olsson, G. (1993). Reduced Order Models for On-line Parameter Identification of the Activated Sludge Process. *Wat. Sci. Tech.*, **28**(11-12), 173-183.

Jørgensen, S. E. (1992). *Integration of Ecosystem Theories: A Pattern*, Kluwer Academic Publishers, Dordrecht.

Kalman, R. (1960). A New Approach to Linear Filtering and Prediction Problems, *Transactions of ASME, Series D, Journal of Basic Engineering*, **82**, 35-45.

Ljung, L. (1987). *System Identification - Theory for the User*, Prentice-Hall, Englewood Cliffs, New Jersey, 519 pp.

Norton, J. P. (1986). *An Introduction to Identification*, Academic Press, London.

O'Neill, R. V. and Gardner, R. H. (1979). Sources of Uncertainty in Ecological Models. In: *Methodology in Systems Modelling and Simulation*, Innis, G. S. and O'Neill, R. V. (eds), North-Holland, Amsterdam, pp. 447-463.

Pedersen, J. and Sinkjær, O. (1992). Test of the Activated Sludge Models Capabilities as a Prognostic Tool on a Pilot Scale Wastewater Treatment Plant. *Wat. Sci. Tech.*, **25**(6), 185-194.

Popper, K. R. (1980). *The Logic of Scientific Discovery*, Hutchinson, London, 480 pp.

Reckhow, K. H. and Chapra, S. C. (1983). Confirmation of Water Quality Models, *Ecological Modelling*, **20**, 113-133.

Reichert, P. (1994). *Concepts Underlying a Computer Program for the Identification and Simulation of Aquatic Systems*, Ph.D. thesis, Swiss Federal Institute of Environmental Science and Technology (EAWAG), Dübendorf, Switzerland.

Reichert, P. and Omlin, M. (1997). On the Usefulness of Overparameterized Ecological Models *Ecological Modelling*, **95**, 289-299.

Rose, M. R. (1987). *Quantitative Ecological Theory - An Introduction to Basic Models*, Croom Helm, London.

Schwarz, G. (1978). Estimating the Dimension of a Model. *The Annals of Statistics*. **6**(2), 461-464.

Singh, M. G. (ed) (1987). *Systems & Control Encyclopedia*, 8 volumes + supplements, Pergamon Press, Oxford.

Spriet, J. A. (1985). Structure Characterization - An Overview. In: *Identification and System Parameter Estimation 1985*, Barker, H. A. and Young, P. C. (eds), Proceedings of the 7th IFAC/IFORS Symposium, Volume 1, pp. 749-756.

Söderström, T. and Stoica, P. (1989). *System Identification*, Prentice Hall, New York.

Tong, H. (1990). *Non-linear Time Series - A Dynamical System Approach*, Clarendon Press, Oxford, 564 pp.

Tulleken, H. (1992). *Grey-box Modelling and Identification Topics*, Ph.D. thesis, Delft University of Technology, Holland, 164 pp.

Vanrolleghem, P. A. and Van Daele, M. (1994). Optimal Experimental Design for Structure Characterization of Biodegradation Models: On-line Implementation in a Respirographic Biosensor. *Wat. Sci. Tech.*, **30**(4), 243-253.

Vanrolleghem, P. A., Van Daele, M. and Dochain, D. (1995). Practical Identifiability of a Biokinetic Model of Activated Sludge Respiration. *Water Research*, **29**, 2561-2570.
Young, P. C. (1993). *Concise Encyclopedia in Environmental Modelling*, Pergamon Press.

## GLOSSARY

### Model constituents

**constant**: model constituent, whose value is constant throughout all possible applications of the model.
**forcing function**: function used as model input.
**input/output model**: model that describes system behavior as being a function of only present input and past inputs and outputs.
**model**: abstraction of reality.
**model structure**: the relations between inputs, outputs and eventually states formulated as equations.
**observation equation**: equation in state-space model, that relates the state variables to the outputs (sometimes also denoted as output equation).
**parameter**: model constituent, whose value needs to be determined for each specific application of the model.
**state**: present situation of the system as described by the model.
**state variable**: model constituent in state-space models, acting as mediator between inputs and outputs and used for a descriptive representation of the system.
**state-space model**: model that includes a descriptive representation of the system by means of an additional set of state variables.
**state-transition equation**: function that relates the future state of the system to the present state and inputs.
**transfer function**: same as input/output function.
**variables**: inputs, outputs and eventually state variables in model equations.

### Model attributes

**adaptive**: model that interacts with the real system and changes the values of its inputs or state variables depending on past output values.
**aggregated**: model that contains state variables that represent functional classes of different constituents (e.g. organisms) or that simplifies the spatial configuration of a system by lumping it together (cf. segregated).
**black-box**: model that describes the observed behaviour of the corresponding subsystem without being based on the mechanisms of this subsystem (cf. white box).
**complex**: a relative attribute that values whether the model contain more state variables, parameters, forcing functions, etc. or an attribute that qualifies that (irrespective of the number of variables) there exist chaotic solutions of the model equations (cf. simple).
**continuous space**: the model resolves the spatial domain of the system continuously (cf. discrete space).
**continuous time**: the model resolves the time axis continuously; the time evolution is usually described by differential equations (cf. discrete time).
**conceptual**: a model that contains a description of the ideas/hypothesis on system behaviour without giving a mathematical formulation.
**deterministic**: the time evolution of the model solution is uniquely determined by the initial state (for state-space models) and the time evolution of inputs (cf. stochastic).
**discrete space**: the model approximates the spatial domain of the system by a number of mixed compartments (cf. continuous space).
**discrete time**: the model divides the time axis into periods of finite length and the output or the state of the model in the next period is given as an algebraic equation depending on the old inputs or states (cf. continuous time).
**distributed parameter**: model with more than one independent variable, i.e. model behaviour is governed by partial differential equation in time and space.

**dynamic**: the model describes the time evolution of a system; a solution of the model may anyway be in steady-state (cf. static, steady-state).

**empirical**: the model equations are not based on generally accepted laws but are just of a descriptive nature (cf. phenomenological, mechanistic).

**grey-box**: the model consists of submodels that are partly based on mechanistic, partly on phenomenological descriptions (cf. black-box, white box).

**heuristic**: model not based on rigorous development but on rules of thumb, feeling, qualitative reasoning.

**linear**: model equations are linear in input variables (for input-output models) or in state variables (for state-space models).

**lumped**: equal to aggregated.

**mechanistic**: model with the goal of describing the mechanisms that lead to the observed behaviour (cf. phenomenological).

**non-linear**: model equations are non-linear in input variables (for input-output models) or in input and state variables (for state-space models).

**phenomenological**: describing the observed phenomena without representing the mechanisms governing the behaviour (cf. mechanistic).

**physical**: model that is based on a description with physical, chemical or biological laws; sometimes also used for small-scale reproductions of a system made in physical materials.

**reduced order**: model of reduced complexity obtained by direct deduction (e.g. by aggregation/lumping) from a more complex basic model.

**reductionist**: hierarchical description of a system by resolving it in subsystems that again are resolved in sub-subsystems until a description level is reached at which a satisfying description is possible without empirical assumptions (in the ideal case down to a description that is based on natural laws).

**segregated**: model that separates variables in more functional classes (cf. aggregated).

**semi-physical**: equal to grey-box .

**simple**: relative attribute that describes that the model equations contain only few state variables, parameters, forcing functions, etc. and the solutions show simple behaviour (periodic or quasi-periodic; cf. complex).

**static**: the model only describes the steady-state solution of a system (cf. dynamic).

**steady-state**: the model only describes the steady-state solution of a system (cf. dynamic).

**stochastic**: the time evolution of the model contains random elements (cf. deterministic).

**time-invariant**: the way the model processes input to output does not change with time.

**transparent**: equal to white-box.

**white-box**: model that describes a system by one or several submodels (white boxes) that describe the observed behaviour of the corresponding subsystem by describing the relevant mechanisms of this subsystem (cf. black box).

Terms of model building

**calibration**: the same as parameter estimation but not necessarily by using statistical methods.

**corroboration**: the same as validation; the term introduced by Popper (1980) makes it clearer that the correctness of the model cannot be proved and that each successful test only increases the belief that the model is correct (cf. confirmation, falsification, validation, verification).

**confirmation**: the same as validation; the term makes it clearer that the correctness of the model cannot be proved and that each successful test only increases the belief that the model is correct (cf. corroboration, falsification, validation, verification).

**falsification**: demonstrating the invalidity of a model by showing that the model results deviate significantly from the measurements; confirmation, corroboration, validation and verification are failed trials of falsification (cf. confirmation, corroboration, validation, verification).

**frame definition**: Selection of which components of a system are to be described and specification of classes of models to be included in the model structure selection process and specification of the experimental conditions for use of the model (experimental frame).

**identifiability analysis (structural, practical)**: evaluation of the uniqueness of the estimates of model parameters from measured data. Structural (theoretical, a priori) identifiability analysis assesses the

uniqueness of parameter estimates from ideal data for a given experimental frame, practical (a posteriori) identifiability analysis assesses the accuracy with which parameters can be estimated with a given data set. In the latter case identifiability is not an objective property, but it depends on the required accuracy.

**model building**: the process of finding an adequate model of a system by (iteratively) processing the following model building steps: problem formulation, prior knowledge collection, system identification and model testing (see Fig. 2).

**model reduction**: simplification of an existing model in order to improve its identifiability without loosing the description of the most important phenomena.

**model (structure) selection**: selecting out of a given set of model structures the structure that makes an optimal (as simple as possible but as complicated as required for the intended purpose) description of measured data possible.

**optimal experimental design**: using a (preliminary) model of a system in order to plan an experiment that maximises the possible gain in information.

**parameter estimation (batch, recursive)**: process of finding parameter values that lead to an optimal agreement of model results with measured data by using statistical methods. Time series of data can be used as a whole (batch estimation) or data points from within a moving data window can be used. In the latter case the parameters become time-dependent and the algorithm can be implemented to modify the previous estimate by considering the omitted and the new data points (recursive estimation).

**regression (linear, non-linear)**: the same as parameter estimation, however, usually used for the special case of algebraically given linear or non-linear model equations and using the (weighted) least squares technique goodness-of-fit criterion.

**simulation (interactive, real-time, Monte Carlo)**: calculating the solution of a model for given values of the parameters, inputs and intial values (usually numerically). Interactive simulations are processed on a computer which allows the user to interact with the program (stop, change parameter values, etc.). Monte Carlo simulation is a method to propagate probability distributions of parameters, inputs and initial values to probability distributions of model results by performing a lot of simulations with parameter values sampled randomly from the probability distribution of the parameters, inputs and initial values.

**system identification**: finding a model to solve a given problem by (iteratively) processing the following identification steps: frame definition, model structure selection, parameter estimation, model diagnosis. (see Fig. 2).

**structure characterisation**: the same as model (structure) selection.

**uncertainty analysis**: estimating the uncertainty of model predictions and analysing the sources of uncertainty.

**validation**: test of a model with a data set not used for identification; note that such tests only increase the belief in the correctness of the model, it is not possible to prove that the model is correct (cf. confirmation, corroboration, falsification, verification).

**verification**: the same as validation (cf. confirmation, corroboration, falsification, validation).