# Systems Identification in Systems Biology

**David Šafránek**
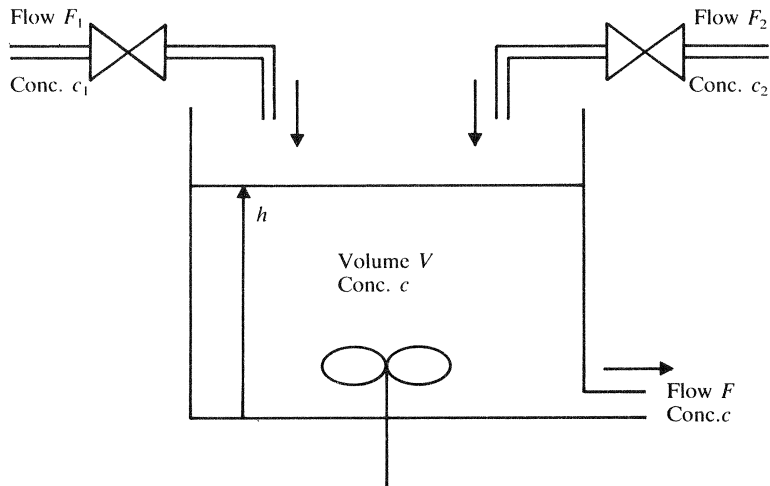
**sybila** systems biology laboratory

Masaryk University
Czech Republic

## Outline

# Outline

1 Introduction

2 The Approach: Parametric Identification
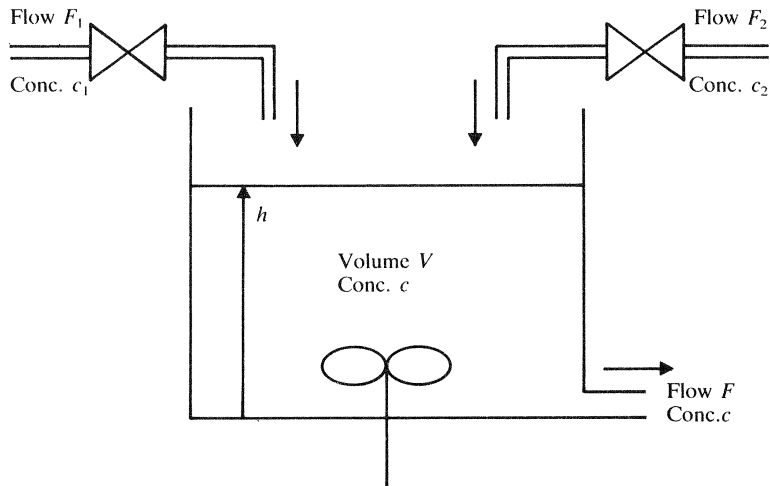
3 System Identifiability Problem

4 Overview of Approaches

$F_1, F_2$ ... input flows (controllable)

$c_1, c_2$ ... input concentrations (uncontrollable)
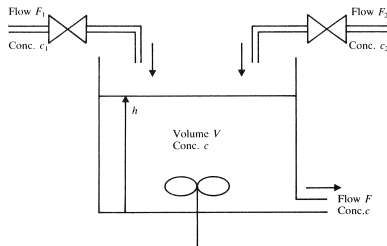
$F, c$ ... outputs of the system, can be observed (measured)

$F_1, F_2$ ... input flows (controllable)
$c_1, c_2$ ... input concentrations (uncontrollable)
GOAL: Keep the outputs $F, c$ constant.

dynamics of the volume: $\frac{dV}{dt} = F_1 + F_2 - F$

$F = a\sqrt{2gh}$ (Torricelli's law of fluid dynamics)
where $a$ ... effective area of the flow, $g \sim 10m/sec^2$

$V = Ah$ where $A$ is tank area (independent on $h$)

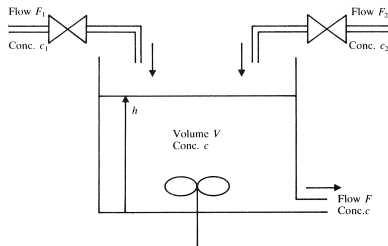dynamics of the volume: $\frac{dV}{dt} = F_1 + F_2 - F$

$F = a\sqrt{2gh}$ (Torricelli's law of fluid dynamics)
where $a$ ... effective area of the flow, $g \sim 10 m/sec^2$

$V = Ah$ where $A$ is tank area (independent on $h$)

**problems:** $a$ is difficult to obtain, does this form of the Torricelli's law really apply for the real case?

**questions:** how to control cyanobacteria to gain max ethanol

how to control *E. coli* to gain insuline, ...

joint work with P. Krejčí, Masaryk University Brno/Medical Genetics Institute, Cedars-Sinai Medical Center, L.A.

What is the right topology?

- western blots
- measurements of protein binding (presence of certain proteins)

Optical density as a proxy of chlorophyll content and cell count

Concentration of dissolved O2 and CO2 influenced by photosynthetic activity

Rate of respiration as an indicator of metabolic changes

Rate of oxygen evolution

Červený, J., Nedbal, L. (2009) Metabolic rhythms of the cyanobacterium *Cyanothece* sp. ATCC 51142 correlate with modeled dynamics of circadian clock. *J. Biol. Rhythms* 24, 295-303.

## Outline

# The Approach: System Identification

- INPUT: controlled perturbance of input stimuli
- OUTPUT: measurements of observed variables
- GOAL: find a system that reliably maps INPUT to OUTPUT

## The Approach: System Identification

- INPUT: controlled perturbance of input stimuli
    - typically interesting patterns exploring most of (expected) systems response
    - pulses, oscillations, ...
- OUTPUT: measurements of observed variables
    - time-series or steady state data
    - not all variables might be observed
    - measurements might be very imprecise $\Rightarrow$ **noisy data**
- GOAL: find a system that reliably maps INPUT to OUTPUT
    - mapping might be non-linear
    - extrinsic noise on both input, output side
    - system might be affected by intrinsic noise (internal stochasticity)

# System Identification Concepts

- **system** $\mathcal{S}$
    - mathematical description of the real-world **process**
    - can be an idealization
    - not necessarily required to be known
- **model structure** $\mathcal{M}$
    - non-parametric (table, mapping, frequency diagram, ...)
    - parametric (with a parameter vector $\theta$) $\mathcal{M}(\theta)$
- **identification method** $\mathcal{I}$
    - depends on available data, kind of the process, ...
- **experimental condition** $\mathcal{E}$
    - concrete setting of identification experiment
    - selection and generation of input signals
    - prefiltering of data

# Parametric Identification: Problem Statement

### Definition

**Parametric model** $\mathcal{M}(\theta)$ describing $n$ dynamically evolving *autonomous* variables is defined by a set of equations:

$$\dot{x}(t) = f(x(t), u(t), p)$$
$$y(t) = g(x(t), s) + \epsilon(t)$$

where

- $x(t) \in \mathbb{R}^n$ for $t \geq 0$ is a vector of **internal model states**
- $u(t) \in \mathbb{R}^u$ for $t \geq 0$ is a vector of **input stimuli**
- $y(t) \in \mathbb{R}^m$ for $t \geq 0$ is a vector of **observables**
- $\epsilon(t)$ is a normally distributed measurement noise

If $m < n$ we speak about *partially observable* models.
Parameter $\theta$ is defined as a vector $\langle p, x(0), s \rangle$.

## Parametric Identification: Problem Statement

$$\chi^2(\theta) = \sum_{k=1}^{m} \sum_{l=1}^{d} \left( y_{kl}^D - y_k(\theta, t_l) \right)^2$$

- $y_{kl}^D$ is $l$th measurement point of the observable $y_k$ taken at time $t_l$
- $y_k(\theta, t_l)$ is model-predicted $y_k$ at time $t_l$ by employing parameter estimate $\theta$
- parameter estimate $\hat{\theta}$ is obtained as a value that minimizes $\chi^2(\theta)$:

$$\hat{\theta} = argmin \left[ \chi^2(\theta) \right].$$

- objective function and reduction to optimisation problem

# Parametric Identification: Problem Statement
Interpretation in Biology

- internal states – biochemical substances in the cell
- observables – substances that can be measured in time (e.g., metabolites or fluorescence reporters)
- input stimuli – profile of nutrient support, signalling stimuli or light program
- differential equations define continuous-time deterministic (population-average) evolution of biochemical substances
- autonomity comes from biochemistry and thermodynamics
  - mass-action kinetics, enzyme kinetics, ...
  - in this setting $x(t)$ and $p$ are always positive

- mechanistic models
    - mass-action systems
        - describes rate of any elementary reaction $\sum_{i=1}^{n} X_i \to ...$:

        $$v = k \prod_{i=1}^{n} X_i^{\sigma_i}$$

        where $\sigma_i$ denotes kinetic order given by stoichiometry
        - easily obtainable model structure if reaction network is known
        - non-linearity is regular if stoichiometry $\leq 1$
        - typically leads to over-parametrised models
    - Michaelis-Menten systems
        - enzyme kinetics based on pseudo-steady-state approximation
        - reduces number of variables and parameters
        - but for general case very complicated non-linear equations
        - similar are Hill systems (generalisation of MM)

# Parametric Identification: Problem Statement
Mathematical Models in Biology

- canonical models
  - S-systems
    - for each species $X_i$ one set of influxes and one set of effluxes is specified in terms of power-law functions:

$$\dot{X}_i = \alpha_i \prod_{j=1}^{n} X_j^{\sigma_{ij}} - \beta_i \prod_{j=1}^{n} X_j^{\rho_{ij}}$$

      where $n$ is the number of all system variables, $\alpha, \beta$ are rate constants for production and degradation, $\sigma, \rho \in \mathbb{R}$ are kinetic orders
  - generalised mass-action (GMA) systems
    - for each species $X_i$ a sum of influxes/effluxes is specified (not aggregated)

$$\dot{X}_i = \sum_{k=1}^{n_i} \left( \gamma_{ik} \prod_{j=1}^{n} X_j^{f_{ikj}} \right)$$

      where $n_j$ is number of fluxes affecting $X_i$, $\gamma$ positive rate constants, and $f \in \mathbb{R}$

# Outline

## System Identifiability: Theoretical Concept

Define the (theoretical) set of exact parameter values:

$$D_T(\mathcal{S}, \mathcal{M}) = \{\theta \mid \mathcal{M}(\theta) \text{ matches the system behaviour }\}$$

Ideally this set should be a singleton. In case of higher cardinality we speak about *overparameterization*.

Assume an estimate $\hat{\theta}(N; \mathcal{S}, \mathcal{M}, \mathcal{I}, \mathcal{E})$ where $N$ is the number of measurements in observed variable $y$.

## System Identifiability: Theoretical Concept

Define the (theoretical) set of exact parameter values:

$$D_T(\mathcal{S}, \mathcal{M}) = \{\theta \mid \mathcal{M}(\theta) \text{ matches the system behaviour }\}$$

Ideally this set should be a singleton. In case of higher cardinality we speak about *overparameterization*.

Assume an estimate $\hat{\theta}(N; \mathcal{S}, \mathcal{M}, \mathcal{I}, \mathcal{E})$ where $N$ is the number of measurements in observed variable $y$.

### Definition

System $\mathcal{S}$ is (parameter) **identifiable under** $\mathcal{M}$, $\mathcal{I}$ and $\mathcal{E}$ iff $\hat{\theta}(N; \mathcal{S}, \mathcal{M}, \mathcal{I}, \mathcal{E}) \to D_T(\mathcal{S}, \mathcal{M})$ as $N \to \infty$.

## System Identifiability: Confidence Intervals

$\hat{\theta}_i$ is associated a confidence interval $[\sigma_i^-, \sigma_i^+]$ with the meaning that true value of $\theta_i$ is located in $[\sigma_i^-, \sigma_i^+]$ with probability $\alpha$

- asymptotic confidence

$$\sigma_i^{\pm} = \hat{\theta}_i \pm \sqrt{\Delta_\alpha(\chi^2) \cdot C_{ii}}$$

where

- $\Delta_\alpha(\chi^2)$ is $\alpha$-quantile for $\chi^2$
- $C = 2 \cdot H^{-1}$
- $H$ is Hessian matrix (describing curvature of $\chi^2$ around $\hat{\theta}_i$ by second-order partial derivatives)

## System Identifiability: Confidence Intervals

$\hat{\theta}_i$ is associated a confidence interval $[\sigma_i^-, \sigma_i^+]$ with the meaning that true value of $\theta_i$ is located in $[\sigma_i^-, \sigma_i^+]$ with probability $\alpha$

- asymptotic confidence

$$\sigma_i^{\pm} = \hat{\theta}_i \pm \sqrt{\Delta_\alpha(\chi^2) \cdot C_{ii}}$$

where

- $\Delta_\alpha(\chi^2)$ is $\alpha$-quantile for $\chi^2$
- $C = 2 \cdot H^{-1}$
- $H$ is Hessian matrix (describing curvature of $\chi^2$ around $\hat{\theta}_i$ by second-order partial derivatives)

- gives a good approximation of actual uncertainty of $\hat{\theta}_i$ if:
  - data have small error
  - amount of data is large wrt number of parameters
  - exact if $y(t)$ depends linearly on $\theta$

## System Identifiability: Confidence Intervals

- finite sample confidence

$$\{\theta \mid \chi^2(\theta) - \chi^2(\hat{\theta}) < \Delta_\alpha\}$$

where $\Delta_\alpha$ is $\alpha$-quantile as in the previous case

- gives an approximation of actual uncertainty of $\hat{\theta}_i$ up-to a statistically computed threshold

# System Identifiability

### Definition

Parameter $\theta_i$ is **identifiable** iff the confidence interval $[\sigma_i^-, \sigma_i^+]$ of the estimate $\hat{\theta}_i$ is finite.

# System Identifiability

### Definition

Parameter $\theta_i$ is **identifiable** iff the confidence interval $[\sigma_i^-, \sigma_i^+]$ of the estimate $\hat{\theta}_i$ is finite.

Reasons leading to non-identifiability:

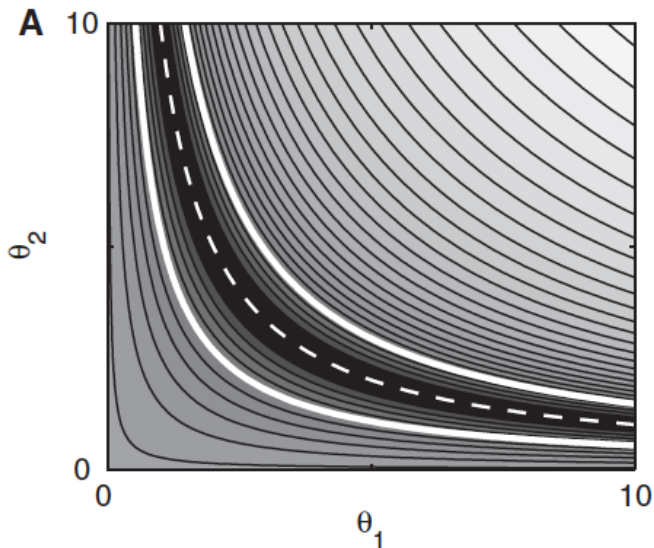- structural: model structure
- practical: precision of measured data

# Structural Identifiability

### Definition

A parameter $\theta_i$ is **structurally identifiable** if a unique minimum of $\chi^2(\theta)$ exists with respect to $\theta_i$.

- structural identifiability requires *uniqueness* of the solution
- redundant parameterisation of the model causing insufficient mapping of internal states $x$ to observables $y$
- denote $\theta_{amb} \subset \theta$ the set of ambiguous parameters
- values of $\theta_{amb}$ may be varied without any change in $y$ (and thus $\chi^2(\theta)$ keeps constant)
- in such a case there must be functional relations $h$ among the parameters in $\theta_{amb}$ that are invariant wrt $\chi^2(\theta)$, and moreover:

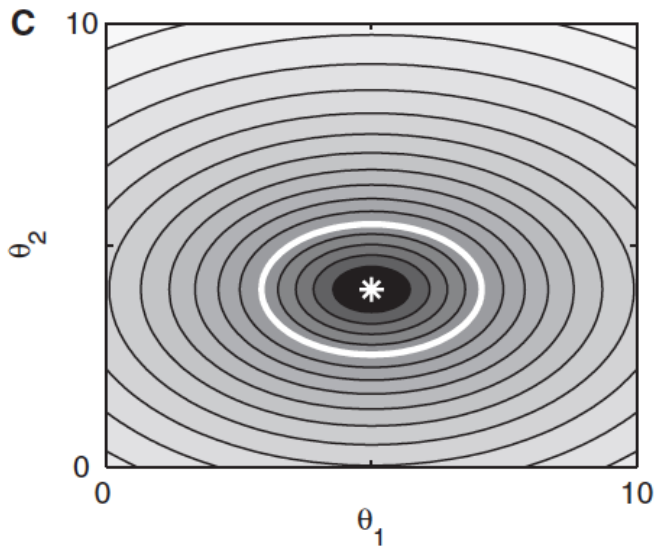$$\forall i, \theta_i \in \theta_{amb}.\, \sigma_i^- = -\infty \land \sigma_i^+ = \infty$$

functional relation between parameters: $h(\theta_{amb}) = \theta_1 \cdot \theta_2 - 10 = 0$
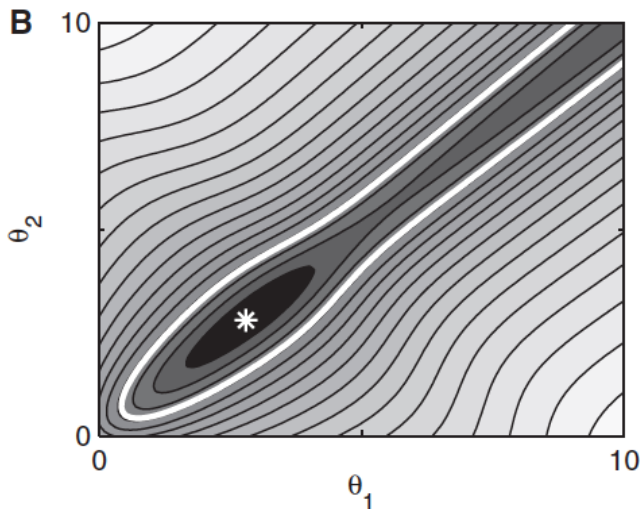
# Practical Identifiability

### Definition

A parameter estimate $\hat{\theta}_i$ is **practically non-identifiable** if the finite sample confidence interval is infinitely extended in decreasing and/or increasing direction although there exists a unique minimum of $\chi^2$.

- practical identifiability implies structural identifiablitity
- practical non-identifiability does not decide on structural identifiability
- detailed analysis can be used to improved experiment design

confidence region is infinitely extended for $\theta_1 \to \infty$ and $\theta_2 \to \infty$

## Detecting Identifiability

- differential algebraic methods to analyse the system equations can detect structural identifiability, computionally hard
- detection of $\chi^2$ flateness using simulated and experimental data
  - approximation of curvature measures by quadratic approximation of $\chi^2$ at $\hat{\theta}$
  - computation of Hessian or Fisher information matrix
  - appropriate for linear relations $h$ among parameters
  - practical non-identifiability cannot be detected

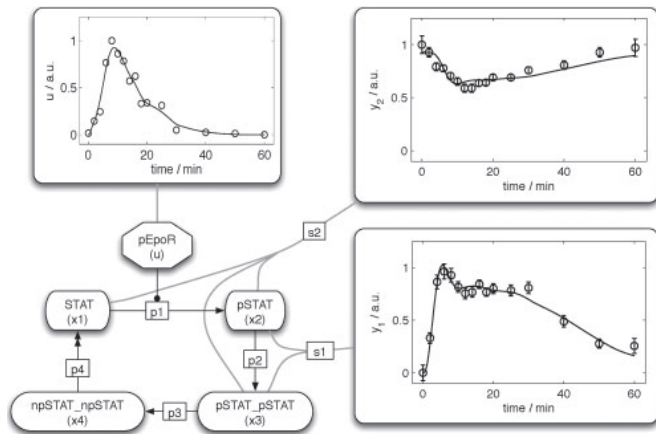# Detecting Identifiability
Profile Likelihood Method by Raue et al. 2009

- explore the parameter space for each parameter in the direction of least increase in $\chi^2$
- in particular this allows to follow the functional relations $h(\theta_{sub}) = 0$
- for practical identifiability detect crossing of the quantile threshold
- profile likelihood $\chi^2_{PL}$ is defined for each parameter $\theta_i$:

$$\chi^2_{PL}(\theta_i) = min_{\theta_{j \neq i}} \left[ \chi^2(\theta) \right].$$

- suggestion of additional targeted measurements
- need measurements that narrow the confidence interval
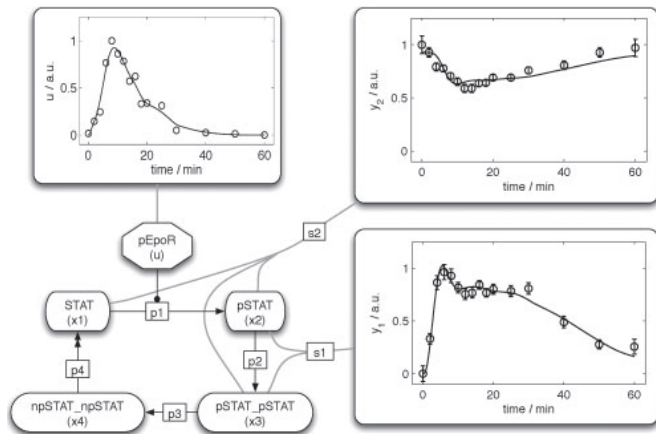- explore trajectories along PL of $\theta_i$ to improve estimation of $\theta_i$

studied system, external stimuli and measured vs. simulated data

studied system, external stimuli and measured vs. simulated data

$$\dot{x}_1 = -p_1 \cdot x_1 \cdot u + 2 \cdot p_4 \cdot x_4^\tau \qquad y_1 = s_1 \cdot (x_2 + 2 \cdot x_3)$$
$$\dot{x}_2 = +p_1 \cdot x_1 \cdot u - p_2 \cdot x_2^2 \qquad y_2 = s_2 \cdot (x_1 + x_2 + 2 \cdot x_3)$$
$$\dot{x}_3 = +\tfrac{1}{2} \cdot p_2 \cdot x_2^2 - p_3 \cdot x_3$$
$$\dot{x}_4 = +p_3 \cdot x_3 - p_4 \cdot x_4^\tau$$

profile likelihood and its quadratic approximation

relations among parameters

further PL-based analysis for experimental planning

# Outline

# Parameter Identification: Approaches Overview

- bottom-up vs. top-down modelling
    - bottom-up means detailed reconstruction from first principles
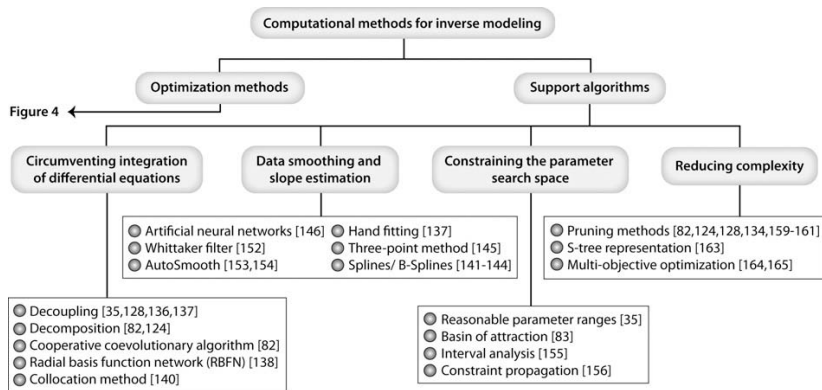    - top-down (inverse) approach starts from high-throughput data
- steady-state vs. transient modelling
    - steady-state data give simplifying assumption (time is abstracted by long-run view)
    - works well for processes with a unique stable state
    - availability of internal system variables at steady-state (e.g., metabolism)
    - transient analysis more complicated (requires detection of initial states and appropriate time-series resolution is needed to inverse modelling)

# Inverse Modelling Approach



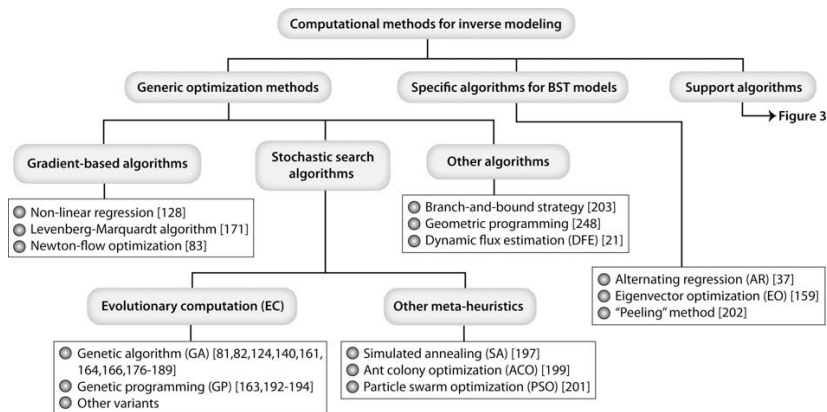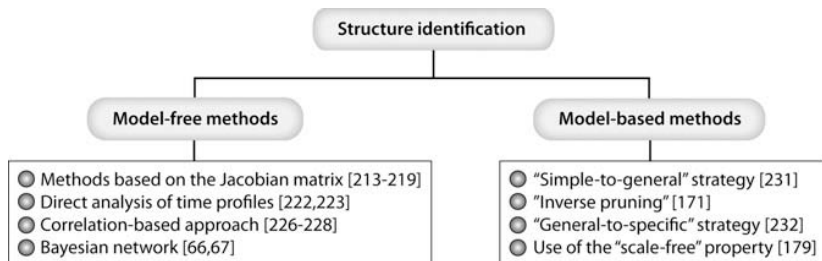| Tasks | | Challenges | Solutions |
|---|---|---|---|
| Parameter Estimation | Data | ○ Overly noisy data<br>○ Missing data points<br>○ Uncertainties about the measurements<br>○ Ill-posed data matrix<br>○ Non-informative data profile | ○ Check data consistency<br>○ Data diagnoses (*e.g.* collinearity)<br>○ Data preprocessing (*e.g.* pooling variables)<br>○ Concept map modeling |
| | Model | ○ Model selection criteria<br>  - Dynamic flexibility<br>  - Mathematical approximation<br>  - Mathematical tractability<br>  - Interpretability of results | ○ BST models: S-system, GMA<br>○ Lin-log model<br>○ Saturable and Cooperative Formalism (SC formalism)<br>○ Determination of model suitability |
| | Computation | ○ Computational capacity<br>○ Slow convergence<br>○ Lacking convergence or convergence to local minima<br>○ Computational cost for integration of differential equations | ○ Optimization methods<br>○ Supporting algorithms<br>  - Complexity reduction<br>  - Avoiding ODE integration<br>  - Data smoothing and slope estimation<br>  - Parameter search space constraints |
| Structure Identification | Math | ○ Distinctly different yet equivalent solutions<br>○ Non-equivalent solutions with similar error<br>○ Error compensation | ○ Estimation of fluxes<br>○ Data covering wide ranges of variation<br>○ Multiple datasets<br>○ Additional information about some of the parameter values |
| | | **Topology** (structure identification) | ○ Model-free, coarse methods<br>○ Model based methods |

I-Chun Chou, E.O. Voit / Mathematical Biosciences 219 (2009) 57-83

I-Chun Chou, E.O. Voit / Mathematical Biosciences 219 (2009) 57-83

I-Chun Chou, E.O. Voit / Mathematical Biosciences 219 (2009) 57-83

Structure identification

Model-free methods

- Methods based on the Jacobian matrix [213-219]
- Direct analysis of time profiles [222,223]
- Correlation-based approach [226-228]
- Bayesian network [66,67]

Model-based methods

- "Simple-to-general" strategy [231]
- "Inverse pruning" [171]
- "General-to-specific" strategy [232]
- Use of the "scale-free" property [179]

I-Chun Chou, E.O. Voit / Mathematical Biosciences 219 (2009) 57-83

# Parameter Synthesis from LTL Specifications

### Robustness

Given an LTL property $\varphi$ and a parameterized model $\mathcal{M}$ check if $\mathcal{M}(\theta) \models \varphi$ **holds for all possible parameterizations** $\theta \in \mathcal{P}$ (valuations of parameters), $\mathcal{P}$ is called the **parameter space**.

### Parameter Synthesis Problem

Given an LTL property $\varphi$ and a parameterized model $\mathcal{M}$ **find the maximal set $P \subseteq \mathcal{P}$ of parameterizations** such that $\mathcal{M}(\theta) \models \varphi$ for all $\theta \in P$.

### Problem Reduction

Robustness is reduced to Parameter Synthesis Problem by taking the set $\mathcal{P}$ of all possible parameterizations as $P$.

# References

- T. Söderström, P. Stoica. System Identification. Prentice-Hall, 1989.
- I-Chun Chou, E.O. Voit. Recent developments in parameter estimation and structure identification of biochemical and genomic systems. Mathematical Biosciences 219 (2009) 57-83
- A. Raue, C. Kreutz, T. Maiwald, J. Bachmann, M. Schilling, U. Klingmüller and J. Timmer. Structural and practical identifiability analysis of partially observed dynamical models by exploiting the profile likelihood. Bioinformatics, Vol. 25 no. 15 2009, pages 1923-1929.
- discussions with Stephan Müller, Jan van Schuppen, Alessandro Abate