

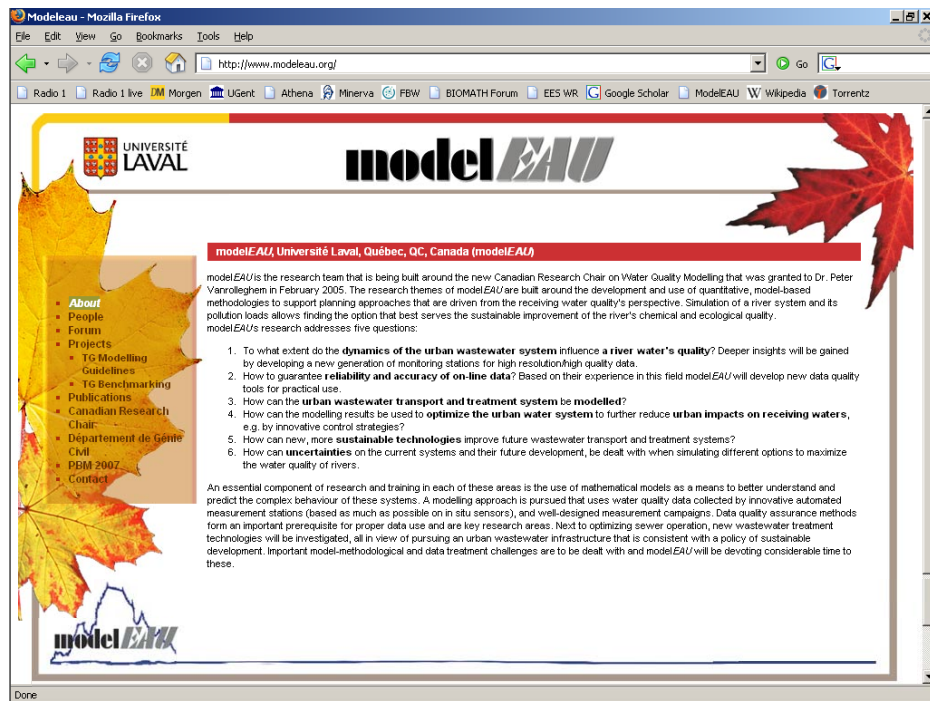
# Potential and limitations of white & black-box modelling concepts for process optimization of SBR WW treatment

**Chemical  
Eng. Seminar**

**Ecole  
Polytechnique  
Montréal**

**22 Feb 2006**

Peter A. Vanrolleghem, Kris Villez & Gurkan Sin



## Overview

- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

## Introduction

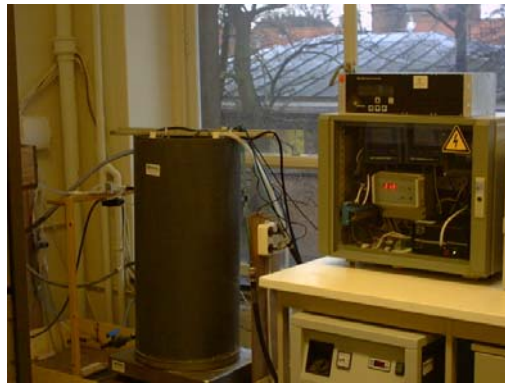
- Activated sludge systems:
  - Underpinning of the microbial community is *not fully* deciphered (yet...)
  - Involves many interacting processes ...
  - Observed behaviour is dynamic & complex
  - *Mechanistic modelling* has proven useful for better understanding & improving operation...
  - *Data driven models* are promising techniques for process monitoring (FDI)

## Overview

- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

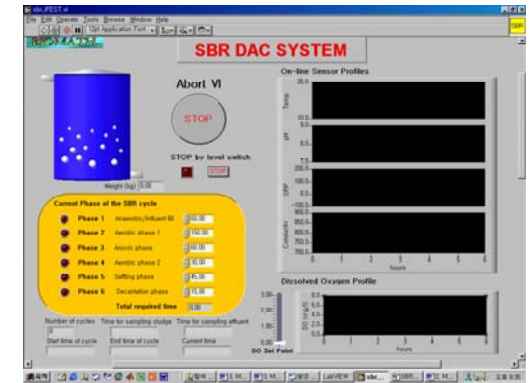
## BIOMATH's pilot SBR

- Performs N & P removal
- Synthetic WW
- $V_{\max}$  80 L
- 15 °C
- SRT 10 d
- HRT 12h



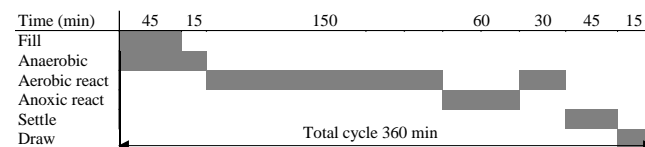
## BIOMATH's pilot SBR

- Fully automated (LabView)
- 5 years of data
- Objectives
  - Stable sludge for Sedifloc project
  - Model-based Optimization of N & P removal
  - Fault detection and diagnosis

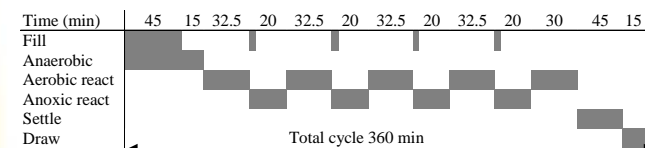


## BIOMATH's pilot SBR

First operation (DO = 2 mg/l) (run for 2 years)

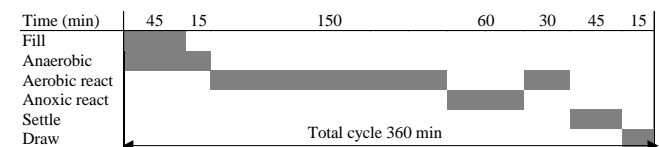


2<sup>nd</sup> operation (DO = 0.5 mg/l) (run for 3.5 months)

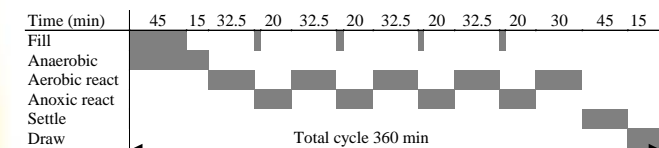


## BIOMATH's pilot SBR

First operation (DO = 2 mg/l) (run for 2 years)



3<sup>rd</sup> operation (DO = 1 mg/l) (run for 1 year)



## Overview

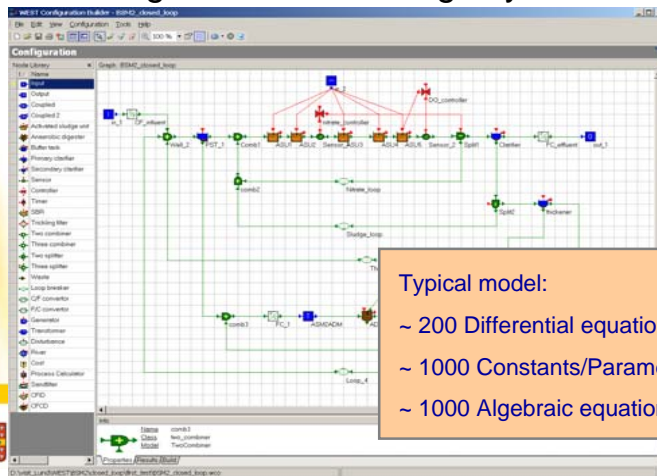
- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

## Mechanistic modelling

- Modelling of activated sludge systems
  - Model structure: internal description of the system
  - Usually constructed using
    - available prior knowledge
    - observed system behavior
  - Selection of an appropriate model structure is *very important* to successfully model the system

## Mechanistic modelling

- Modelling of activated sludge systems:



Typical model:

- ~ 200 Differential equations
- ~ 1000 Constants/Parameters
- ~ 1000 Algebraic equations

## Mechanistic modelling

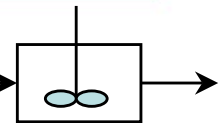
- Mechanistic models for WWTPs:

- Mass balance for compound:

$$\frac{dM}{dt} = \frac{d(VC)}{dt} = \underbrace{Q_{in} C_{in} - Q_{out} C_{out}}_{\text{transport}} + \underbrace{rV}_{\text{conversion}}$$

- with

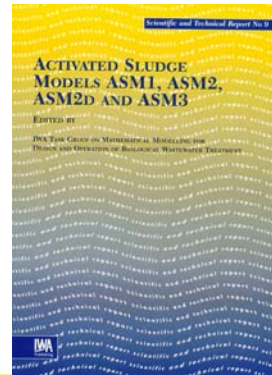
- M: Mass of compound in system (g)
- C: Concentration of compound (g/m<sup>3</sup>)
- V: Volume of system (m<sup>3</sup>)
- Q: flow rate (m<sup>3</sup>/h)
- r: volumetric conversion rate (g/m<sup>3</sup>.h)



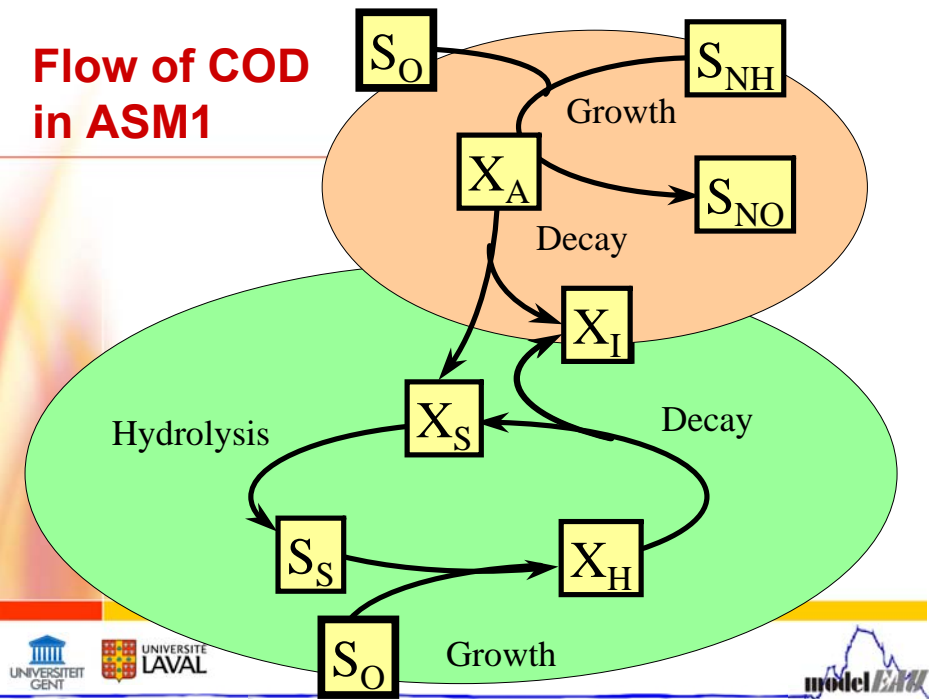
## Mechanistic modelling

- Mechanistic models for WWTPs:

- Henze, M., Gujer, W., Mino T., & van Loosdrecht, M. (2000)
- Activated Sludge Models
  - ASM1
  - ASM2 & ASM2D
  - ASM3
- Scientific and Technical Report No. 9
- IWA Publishing, London.

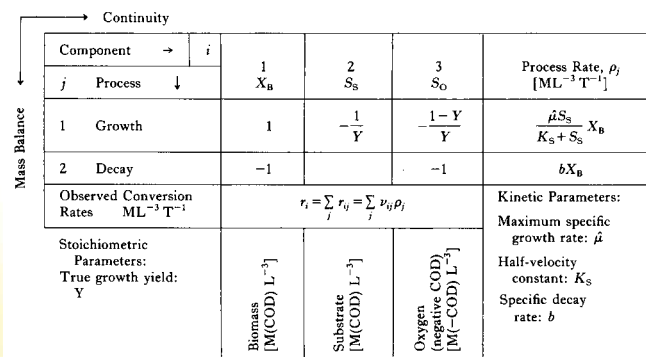


## Flow of COD in ASM1



## Mechanistic modelling

- Mechanistic model structure for WWT:

[illegible]

# The horror matrix

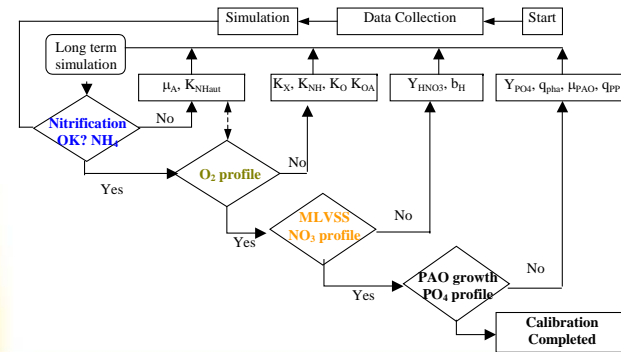


## Overview

- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

## Model-based optimization

- BIOMATH calibration protocol was applied



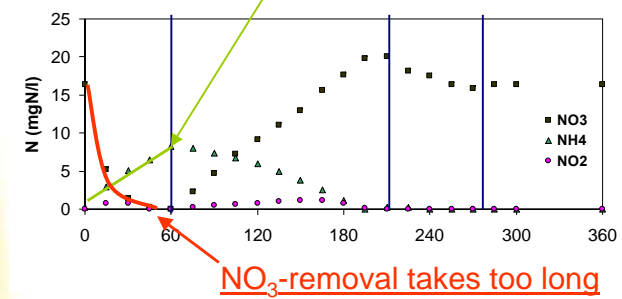
## Model-based optimization (1)

- First iteration (basic operation)
  - Complete COD-removal
  - Complete nitrification
  - Incomplete denitrification (70 % N-removal)
  - 50 % P-removal (limited due to presence of NO<sub>3</sub>)

## Model-based optimization (1)

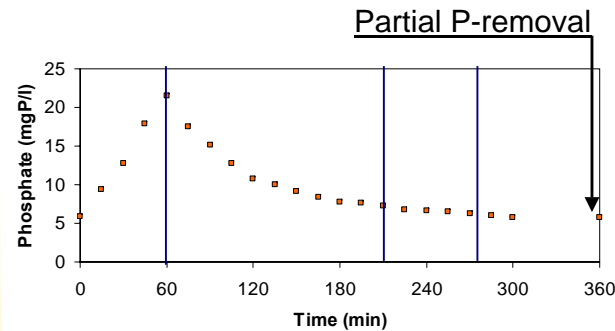
- First iteration (basic operation)

### Limitation in hydrolysis of organic nitrogen



## Model-based optimization (1)

- First iteration (basic operation)



## Model-based optimization (1)

- First iteration (basic operation)

Prior knowledge (e.g. ASM2d)  
+  
System's observation  
(Long-term performance + intensive in-cycle data)

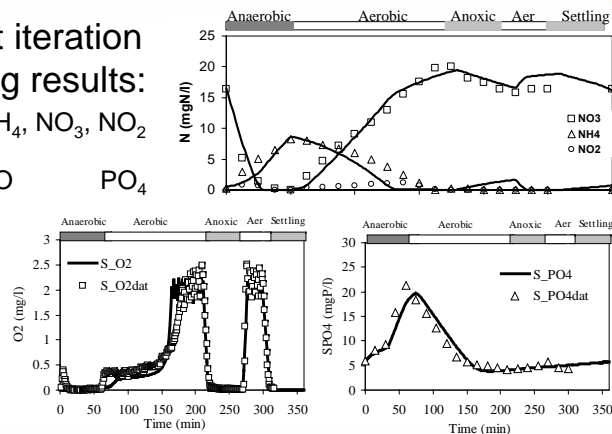


Model structure change needed (ASM2dN)  
(ASM2d is extended with ASM1 hydrolysis model)

## Model-based optimization (1)

- First iteration fitting results:

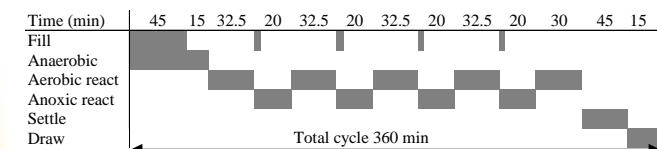
- $\text{NH}_4$ ,  $\text{NO}_3$ ,  $\text{NO}_2$
- $\text{DO}$ ,  $\text{PO}_4$



## Model-based optimization (2)

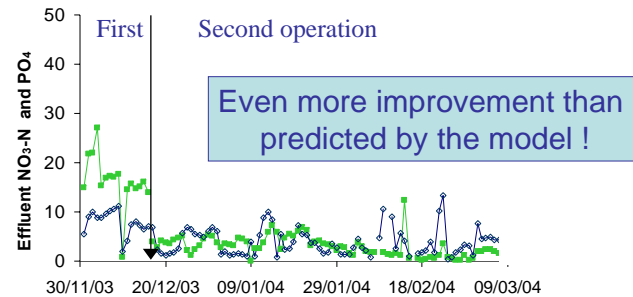
- With the calibrated ASM2dN model, process operation was optimized:
  - Low DO; Step-feed; short aerobic & anoxic phases...

2<sup>nd</sup> operation ( $\text{DO} = 0.5 \text{ mg/l}$ ) (run for 3.5 months)

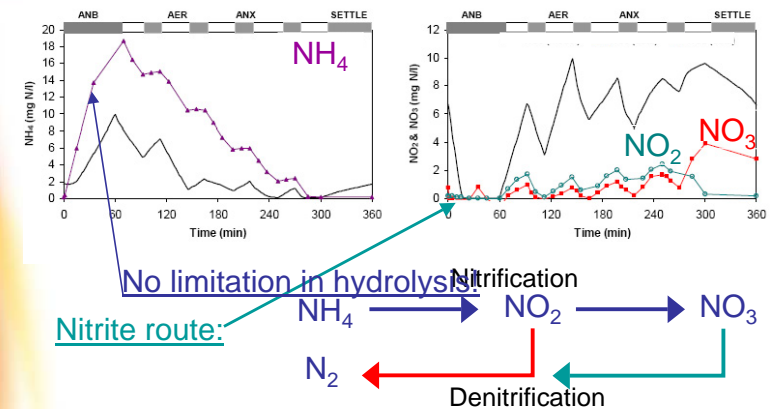


## Model-based optimization (2)

- N-removal improved by 86%
- P-removal improved by 65%



## Model-based optimization (2)



## Model-based optimization (2)

- During 2<sup>nd</sup> operation ASM2dN fails :
  - No longer a limitation of hydrolysis of organic N
  - 2-step nitrification occurs
  - Nitrite route takes place
  - Nitrate removal is better (observed !)
  - No more inhibition of P-removal (observed !)

## Model-based optimization (2)

- Necessary model modifications → ASM2d2N

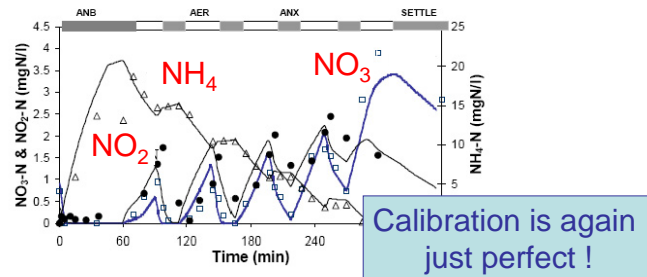
Modified processes of ASM2d (stoichiometry)

Processes	$S_{\text{NH}}$	$S_{\text{NO}_3}$	$S_{\text{NO}_2}$	$S_{\text{N}_2}$	$X_{\text{H}}$	$X_{\text{NH}}$	$X_{\text{NO}}$
$\text{NH}_4$ oxidation	$-i_{\text{NBM}} - 1/Y_{\text{NH}}$		$1/Y_{\text{NH}}$			1	
$\text{NO}_2$ oxidation	$-i_{\text{NBM}}$	$1/Y_{\text{NO}}$	$-1/Y_{\text{NO}}$				1
$\text{NO}_3$ reduction	$-i_{\text{NBM}}$	$-(1-Y_{\text{HNO}_3})/(1.14 Y_{\text{HNO}_3})$	$(1-Y_{\text{HNO}_3})/(1.14 Y_{\text{HNO}_3})$		1		
$\text{NO}_2$ reduction	$-i_{\text{NBM}}$		$-(1-Y_{\text{HNO}_2})/(1.72 Y_{\text{HNO}_2})$	$(1-Y_{\text{HNO}_2})/(1.72 Y_{\text{HNO}_2})$	1		
Lysis of $X_{\text{NH}}$							-1
Lysis of $X_{\text{NO}}$						-1	



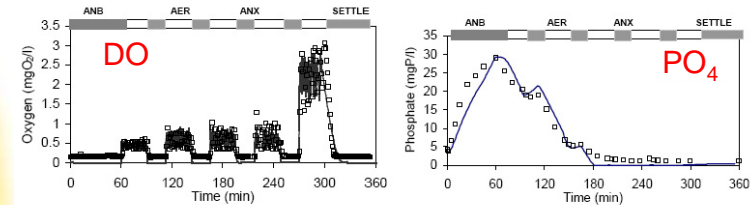
## Model-based optimization (2)

- ASM2d2N fitting results to in-cycle measurements:



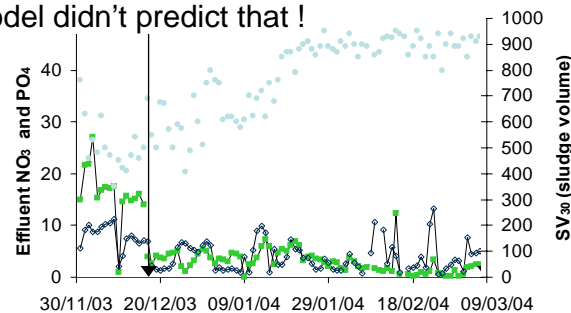
## Model-based optimization (2)

- ASM2d2N fitting results to in-cycle measurements:



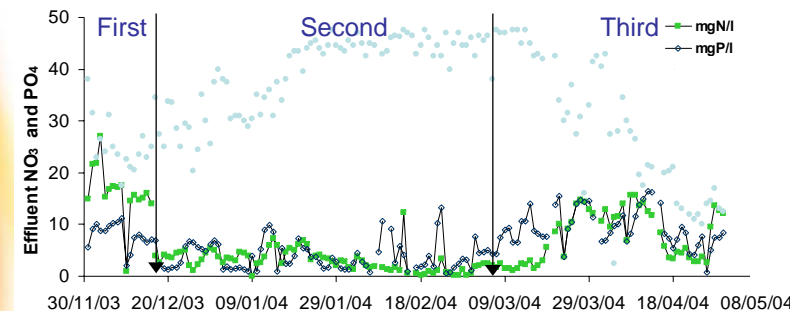
## Model-based optimization (2)

- Nutrient removal is very good, but ...
  - Severe sludge bulking occurs !
  - Model didn't predict that !



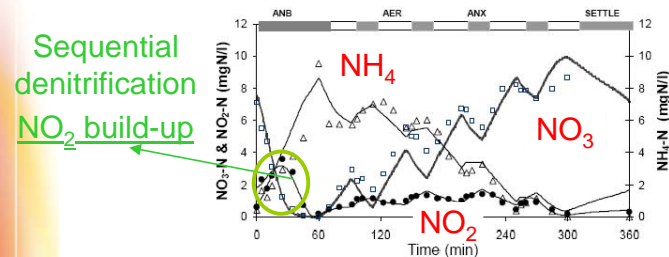
## Model-based optimization (3)

- 3<sup>rd</sup> scenario: DO ↑ (from 0.5 to 1.0 mg/l) & Ca<sup>2+</sup> ↑
- bulking solved, but N,P-removal worse again



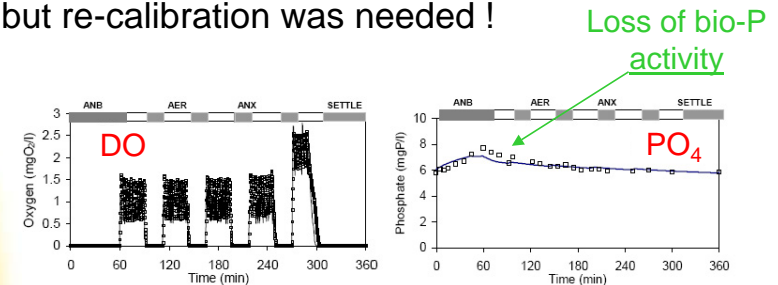
## Model-based optimization (3)

- ASM2d2N performs well
- but re-calibration was needed !



## Model-based optimization (3)

- ASM2d2N performs well
- but re-calibration was needed !



## Model-based optimization (3)

- Validation of the ASM2d2N showed:
  - The model structure remained valid, but ...
  - Parameters of the model had to be changed

## Conclusions

- The 3 models predicted the system dynamics to some extent :
  - Parallel to system changes (operation):
  - Model structure had to be changed twice
    - Hydrolysis
    - $\text{NO}_2$  route (nitrification & denitrification)
  - Parameters had to be changed every time
- Poor predictive power of mechanistic models,
- Not to mention prediction of bulking...

## Conclusions

- The underlying reasons remain unclear, but could be :
  - Unaccounted input disturbances
  - Imperfect model structure
  - ...
  - Or perhaps the system is too complex to *mechanistically* model!
- Biology was proven by DGGE analysis to change significantly after operation changes

## Conclusions

- Models that validly describe system behaviour under a wide range of conditions are not available yet
- But models appear valid within certain (narrow?) boundaries, e.g. under certain operation conditions...
- and models help to understand the system and point to optimization approaches

## Overview

- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

## Black box modelling: Intro



*Data drowning...*

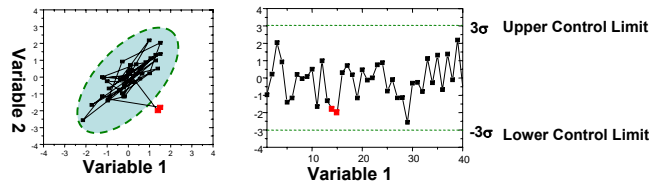
## Black box modelling: Intro

- Many data-driven approaches !
- Here we only consider multivariate statistical analysis:
  - Principal Component Analysis (PCA)
    - process monitoring (fault detection/diagnosis)
  - Partial Least Squares (PLS)
    - prediction in view of control
- Applied to the BIOMATH pilot SBR

## Process monitoring (PCA)

- Monitoring the state of the process: Statistical Process Control (SPC)
- Traditional SPC = Univariate SPC
  - One variable at a time, not efficient
  - Problem of correlation between variables

## Process monitoring (PCA)



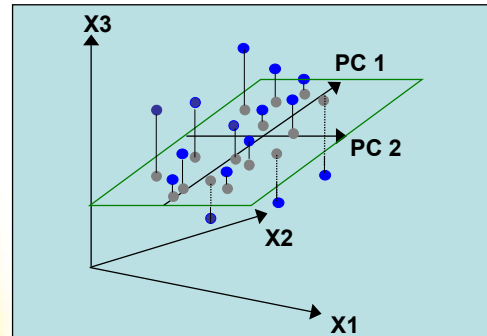
- The deviation is not detected unless the variables are combined
- Most variables are correlated
- The key to early fault detection is the correlation structure, not the original variables

## Process monitoring (PCA)

- Monitoring the state of the process: Statistical Process Control (SPC)
- Traditional SPC = Univariate SPC
  - One variable at a time, not efficient
  - Problem of correlation between variables
- Multivariate SPC
  - Account for interactions among variables
  - Detect upsets and find assignable causes

## Process monitoring (PCA)

### ▪ Geometrical interpretation

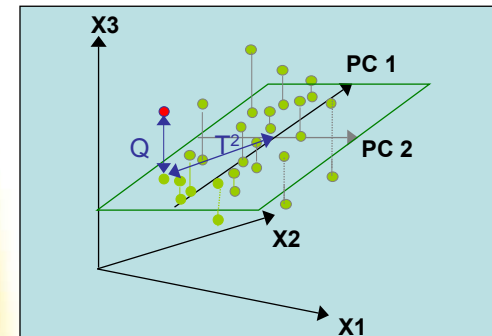


Loads (P): linear combinations of original process variables

Scores (T): projected coordinates of samples

## Process monitoring (PCA)

### ▪ Lack of Model Fit Statistics



Q: distance between the model plane and a sample

$T^2$ : distance within model plane from a sample to the origin

## Process monitoring (PCA)

$$\begin{array}{c} \text{Variables (M)} \\ \text{Samples (N)} \end{array} \begin{array}{|c|} \hline X \\ \hline \end{array} = \begin{array}{c} \text{Scores (NxA)} \\ T \end{array} + \begin{array}{c} \text{Loads (AxM)} \\ P^T \end{array} + \begin{array}{c} \text{Residuals (NxM)} \\ E \end{array}$$

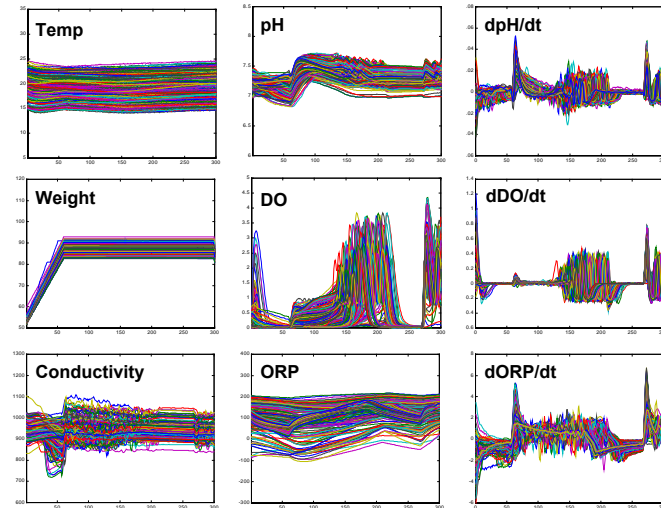
A: the number of principal components ( $A \ll M$ )

## Process monitoring (PCA)

- Applying PCA models to the BIOMATH SBR
- Objective: Develop real-time monitoring
  - Detect the major sources of process disturbances
  - Useful to keep the sludge as stable as possible
- On the basis of simple on-line data, e.g.
  - pH, temperature, weight
  - conductivity
  - dissolved oxygen (DO)
  - oxidation reduction potential (ORP)

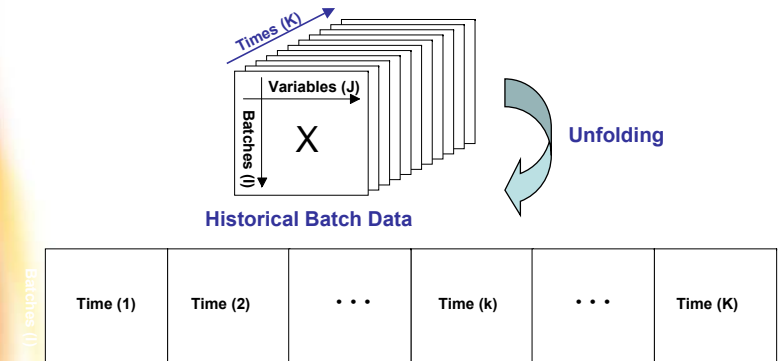


## Process monitoring (PCA)



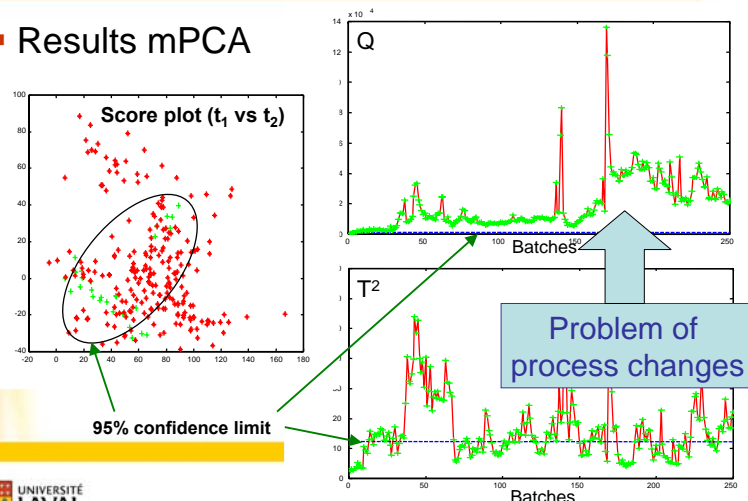
## Process monitoring (PCA)

- Multi-way PCA model



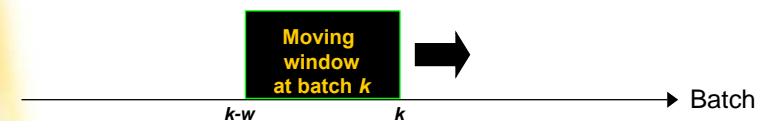
## Process monitoring (PCA)

- Results mPCA



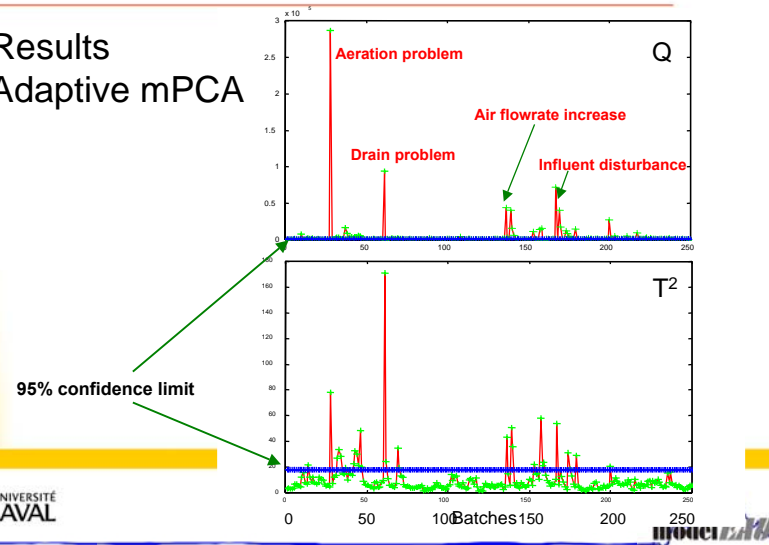
## Process monitoring (PCA)

- To overcome the problem of changing process conditions:
- Adaptive multi-way PCA



## Process monitoring (PCA)

- Results  
Adaptive mPCA



## Conclusions

- Adaptive Multi-way PCA provides more information than adaptive PCA
- Critical process disturbances are well captured in Q & T<sup>2</sup> plots
- Adaptive Multi-way PCA is a powerful tool for monitoring SBR processes

## Overview

- Introduction
- BIOMATH's pilot SBR
- White box modelling
  - Modelling approach
  - Calibration/Validation in an optimization loop
- Black box modelling
  - Background on PCA/PLS multivariate analysis
  - Performance evaluation
- Conclusions

## PLS modelling

- Why PLS ?
  - We want  $Y=f(X)$
  - When dealing with collinear inputs (X), multivariate linear regression (MLR)

$$y = \mathbf{B} \cdot \mathbf{x} = b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_n \cdot x_n$$

will lead to the unbiased regression vector but the estimated regression vectors will have a high variance (very accurate, low precision)

## PLS modelling

- PLS is one way to overcome this problem
  - It trades a bias with a decreased variance of the solution by reducing the dimension of the input space while minimizing the prediction error.

$$\mathbf{x} = \mathbf{t} \cdot \mathbf{p}^t = t_1 \cdot \mathbf{p}_1^t + t_2 \cdot \mathbf{p}_2^t + \dots + t_n \cdot \mathbf{p}_n^t$$

$$\mathbf{y} = \mathbf{u} \cdot \mathbf{q}^t = u_1 \cdot \mathbf{q}_1^t + u_2 \cdot \mathbf{q}_2^t + \dots + u_n \cdot \mathbf{q}_n^t$$

where:  $u_i = b_i \cdot t_i + h_i$ ;  $b_i$ : inner relation coefficient  
 $h_i$ : inner model error

$t_i$  ( $u_i$ ) present the transformed (lower dimensional) input (output) data

## PLS modelling

- Advantages:
  - Dimension reduction of input space
  - More robust estimates of regression vector(s)
- Disadvantages
  - Limited to linear regression conditions

## Neural Net PLS modelling

- PLS is a linear method by definition and thus fails when the relationship between inputs (X) and outputs (Y) is non-linear in nature

$$\mathbf{x} = \mathbf{t} \cdot \mathbf{p}^t = t_1 \cdot \mathbf{p}_1^t + t_2 \cdot \mathbf{p}_2^t + \dots + t_n \cdot \mathbf{p}_n^t$$

$$\mathbf{y} = \mathbf{u} \cdot \mathbf{q}^t = u_1 \cdot \mathbf{q}_1^t + t_2 \cdot \mathbf{q}_2^t + \dots + u_n \cdot \mathbf{q}_n^t$$

where:  $u_i = NN_i(t_i) + h_i$ ;  $NN_i$ : inner neural net  
 $h_i$ : inner model error

## Neural Net PLS modelling

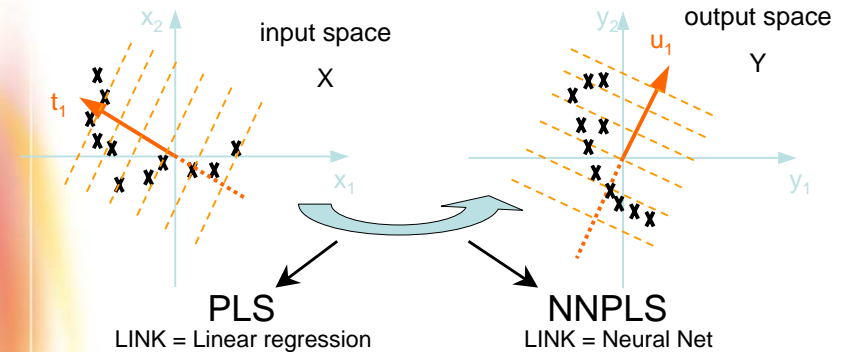
- NNPLS tackles this problem by replacing the (linear) inner relationship coefficient in PLS by a 3-layer network (1 hidden layer).

## Neural Net PLS modelling

- Advantages:
  - Not restricted to linear regression problems
- Disadvantages
  - Additional parameters (number of nodes in hidden layer)
  - => model is more complex

## Neural Net PLS modelling

- Comparison NNPLS / PLS



## Kernel PLS modelling

- KPLS also tackles the problem of nonlinearity
- not by looking for a nonlinear relation  $Y=f(t)$
- but by transforming the input space (X), prior to PLS modelling

## Kernel PLS modelling

- The transformation is chosen such that the input data become “more linear” :

$$f = \Phi(x)$$

$$f = t \cdot p^t = t_1 \cdot p_1^t + t_2 \cdot p_2^t + \dots + t_n \cdot p_n^t$$

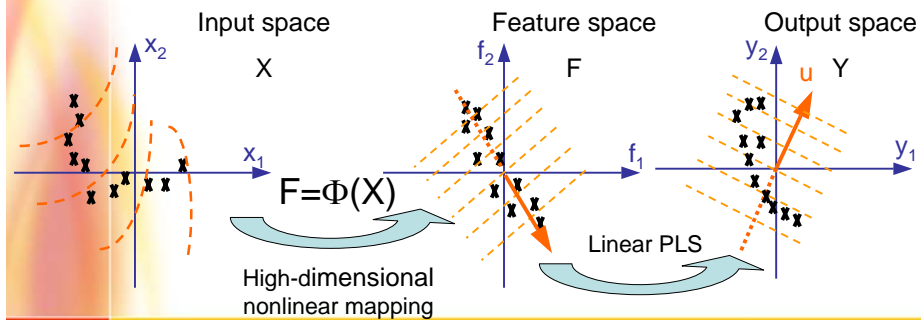
$$y = u \cdot q^t = u_1 \cdot q_1^t + t_2 \cdot q_2^t + \dots + u_n \cdot q_n^t$$

$$\text{where: } u_i = b_i \cdot t_i + h_i ;$$

$b_i$  : inner relation coefficient  
 $h_i$  : inner model error

## Kernel PLS modelling

- The transformation is chosen such that the input data become “more linear” :



## Kernel PLS modelling

- In this work the Gaussian kernel function was applied to transform the input data:

$$f_{ij} = \Phi(x_i, x_j) = k(x_i, x_j) = \exp(-||x_i - x_j||^2 / d)$$

where: d = width of the Gaussian kernel function  
( = extra tuning or meta-parameter)

## Kernel PLS modelling

- Advantages:
  - nonlinear collinearity within the X-space is dealt with
  - No non-linear optimisation required
- Disadvantages
  - Larger computational demand (x10 – x100)
  - Models are hard to interpret (as the transformed inputs are hard to interpret)

## PLS modelling: Results

- Objective: predict SBR effluent quality
  - Total nitrogen
  - NO<sub>3</sub>
  - PO<sub>4</sub>
- using on-line data (1600 batches; ΔT=1 min)
  - pH, temperature, weight
  - conductivity
  - dissolved oxygen (DO)
  - oxidation reduction potential (ORP)

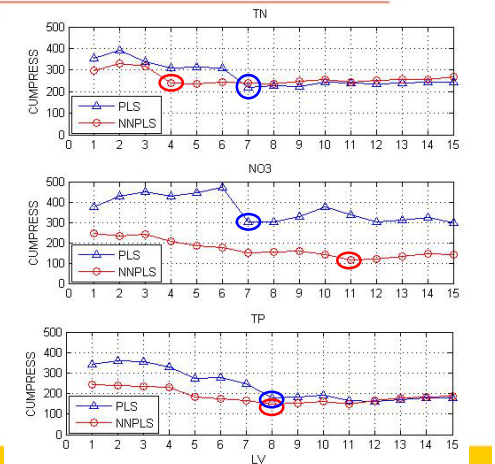


## PLS modelling: Results

- For PLS and NNPLS:
  - Degree of freedom: # of latent variables
  - Selection based on cross-validation  
CUMPRESS = Sum of SSE values

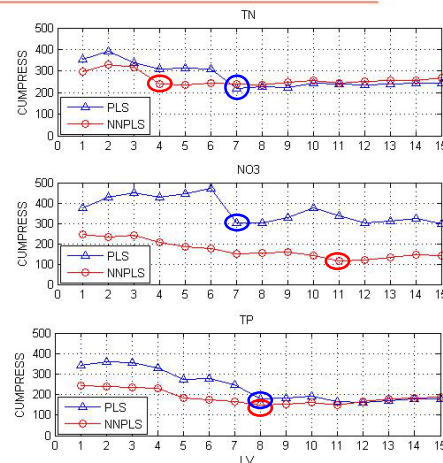
## PLS modelling: Results

- # LV's for prediction of
  - Total N
  - NO<sub>3</sub>
  - PO<sub>4</sub>



## PLS modelling: Results

- # LV's for prediction of
  - Total N
  - NO<sub>3</sub>
  - PO<sub>4</sub>
- NNPLS gives:
  - A lower dimension of the model (TN)
  - Better prediction (NO<sub>3</sub>, PO<sub>4</sub>)

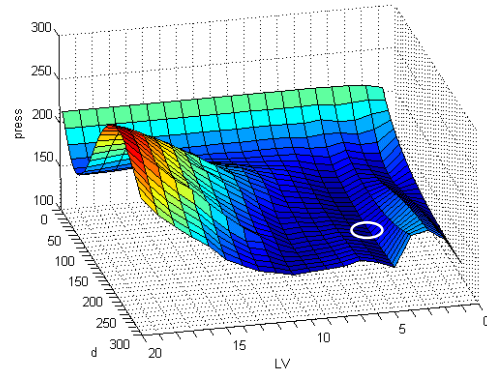


## PLS modelling: Results

- For Kernel PLS :
  - Degree of freedom:
    - # of latent variables
    - width of kernel function d
  - Selection based on cross-validation  
CUMPRESS = Sum of SSE values

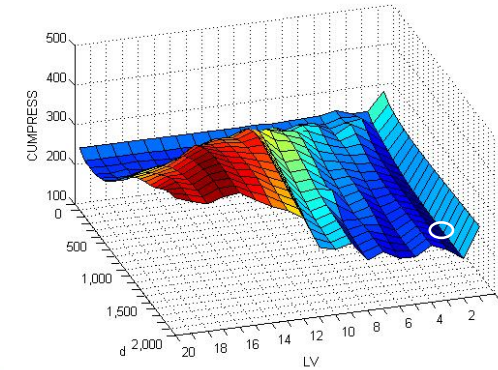
## PLS modelling: Results

- # LV's and d for prediction of Total N



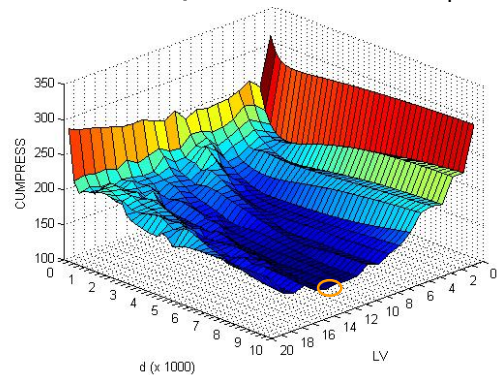
## PLS modelling: Results

- # LV's and d for prediction of  $\text{NO}_3$

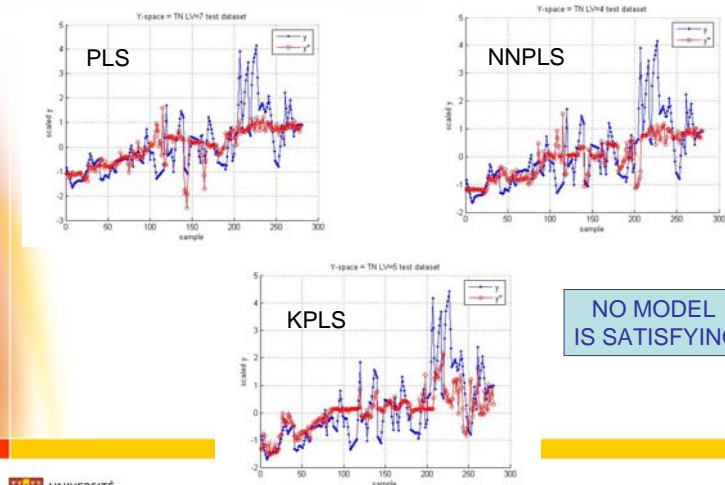


## PLS modelling: Results

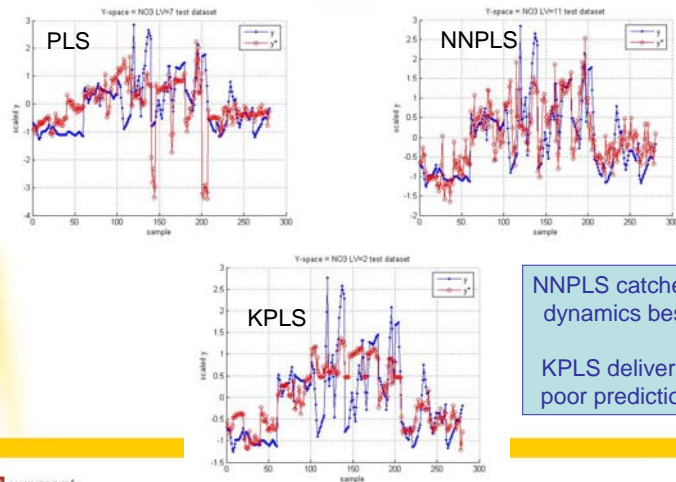
- # LV's and d for prediction of  $\text{PO}_4$



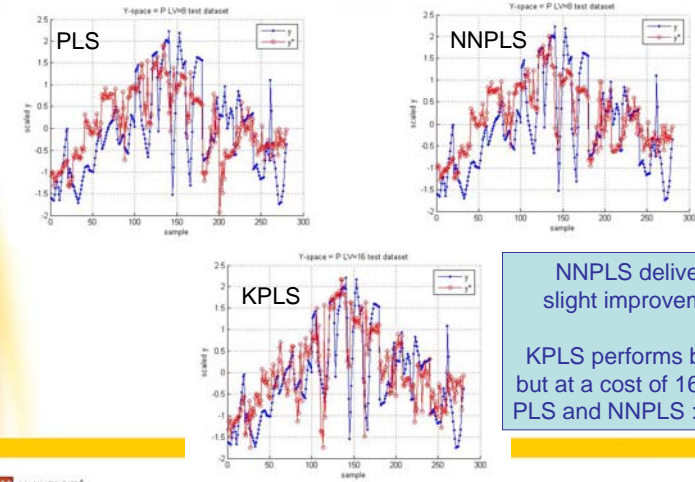
## PLS modelling: Results (Total N)



## PLS modelling: Results (NO<sub>3</sub>)



## PLS modelling: Results (PO<sub>4</sub>)



## PLS modelling: Summary

	output	PLS	NNPLS	KPLS
quality	TN	-	-	-
	NO <sub>3</sub>	-	++	-
	P	+	+	++
cumpress	TN	220	235	147
	NO <sub>3</sub>	303	116	198
	P	179	150	142
LV's	TN	7	4	5
	NO <sub>3</sub>	7	11	2
	P	8	8	16

- TN: unacceptable (dynamics)
- NO<sub>3</sub>: NNPLS is only satisfactory model
- PO<sub>4</sub>: NNPLS is selected (KPLS model needs too many LV's: 16)

## Conclusions

- Still, none of the models is really satisfying
- What did we miss?
  - The inputs data do not describe the process
  - The data are treated as independent observations
    - In fact, they represent a time series  
Autocorrelation should be accounted for
  - The data stem from a large time window (14 months)
    - Equipment, operation and biological changes may not permit a unique (overall) model





Peter Vanrolleghem  
Canada Research Chair in Water Quality Modelling

modelEAU, Département de Génie Civil  
Pavillon Pouliot, Université Laval  
Québec G1k 7P4, QC, Canada

Tel: +1 418 656 5085

Fax: +1 418 656 2928

E-mail: [peter.vanrolleghem@modelEAU.org](mailto:peter.vanrolleghem@modelEAU.org)

and of course also, please visit us at :  
[www.modelEAU.org](http://www.modelEAU.org)