

# Hydraulic Engineering in Petabyte Age

Machine Learning for a Flow Resistance Problem

or

## Man vs. Machine

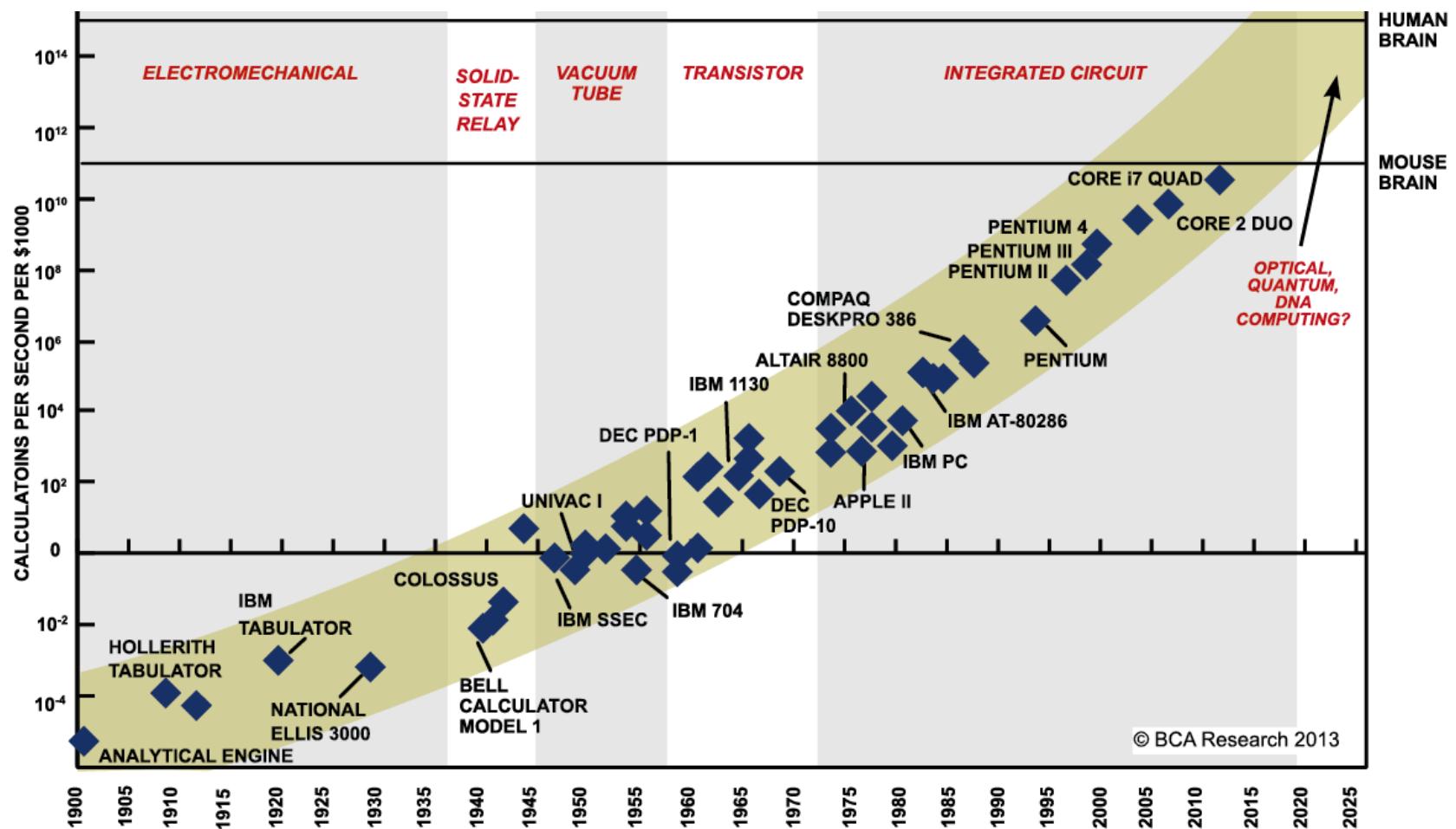
Prof. Vladan Babovic

National University of Singapore

# We Live in the Second Machine Age

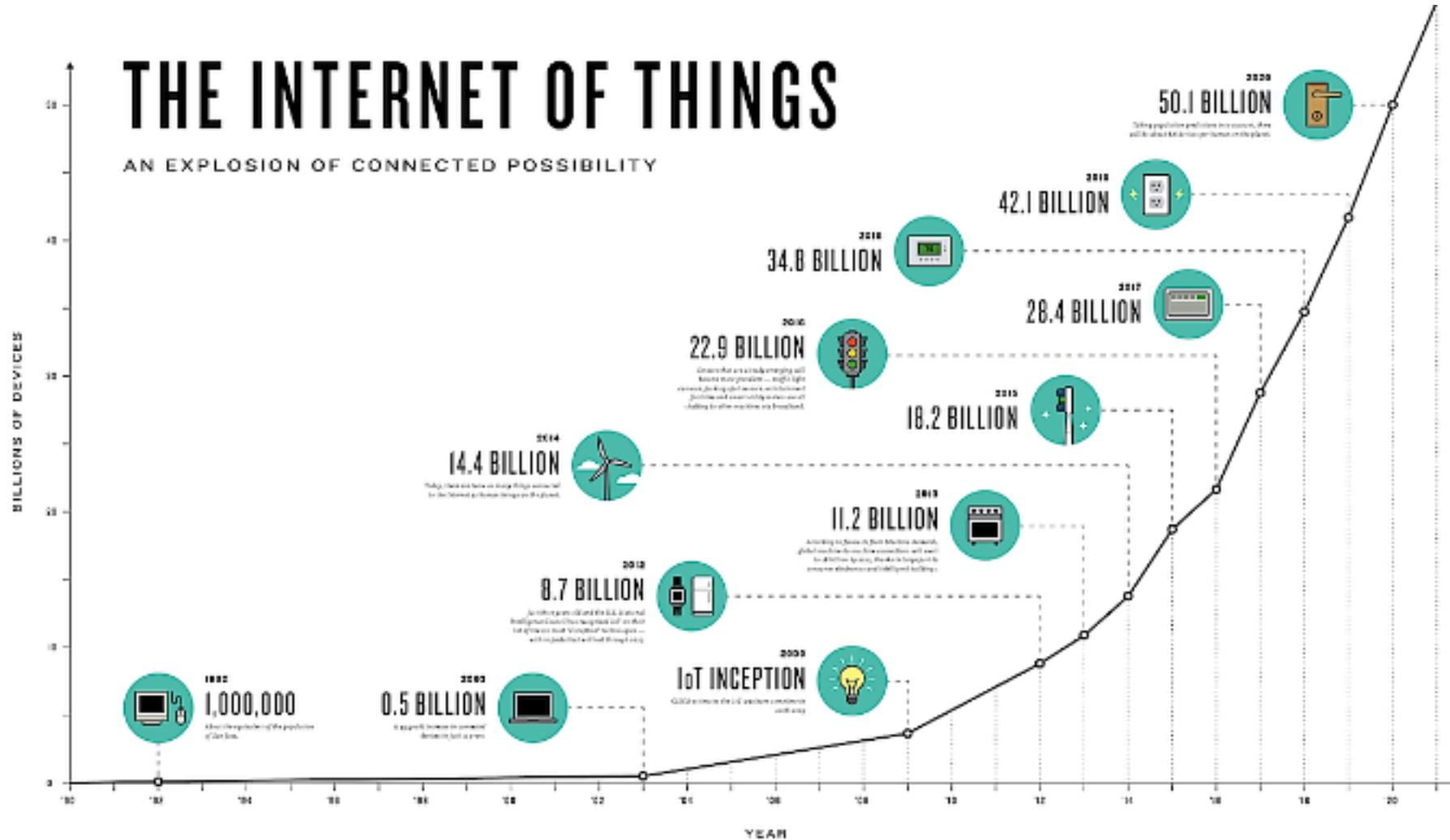
our profession is being disrupted

# Moore's Law

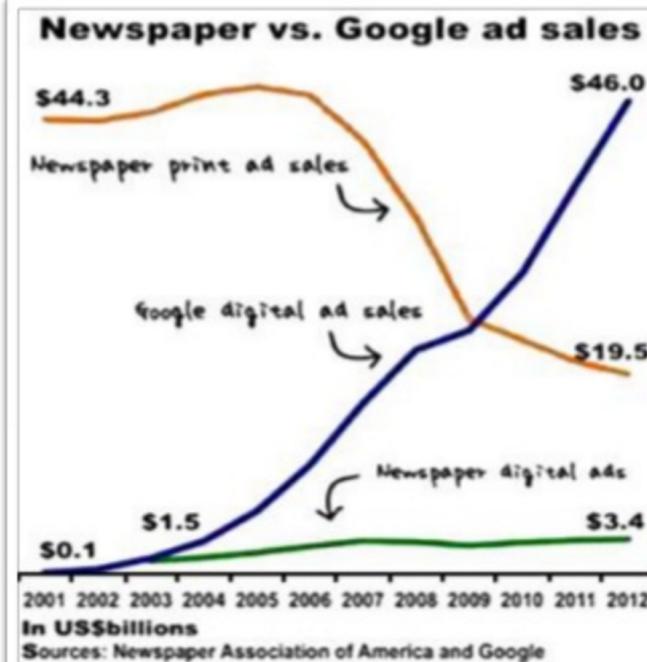
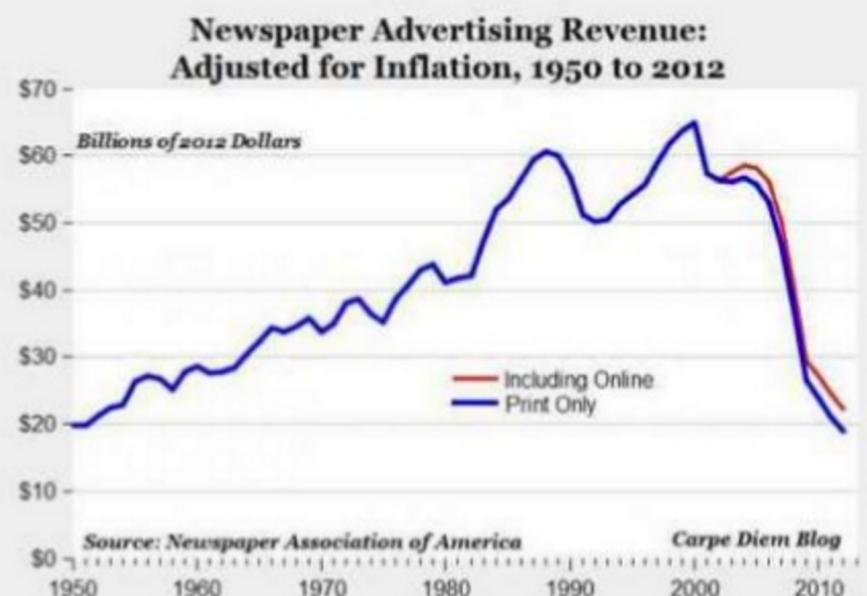


# THE INTERNET OF THINGS

AN EXPLOSION OF CONNECTED POSSIBILITY



# Newspapers



# LINEAR → EXPONENTIAL



1996

MarketCap: \$28B

Employees: 140,000



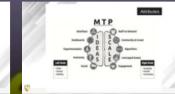
Employees: 7,000



April 2012

MarketCap: \$1B

Employees: 13



Requirements: A Kodak Moment

# The Second Machine Age and Water Management?

a case study in  
vegetation and hydraulic resistance

# Transformations .....



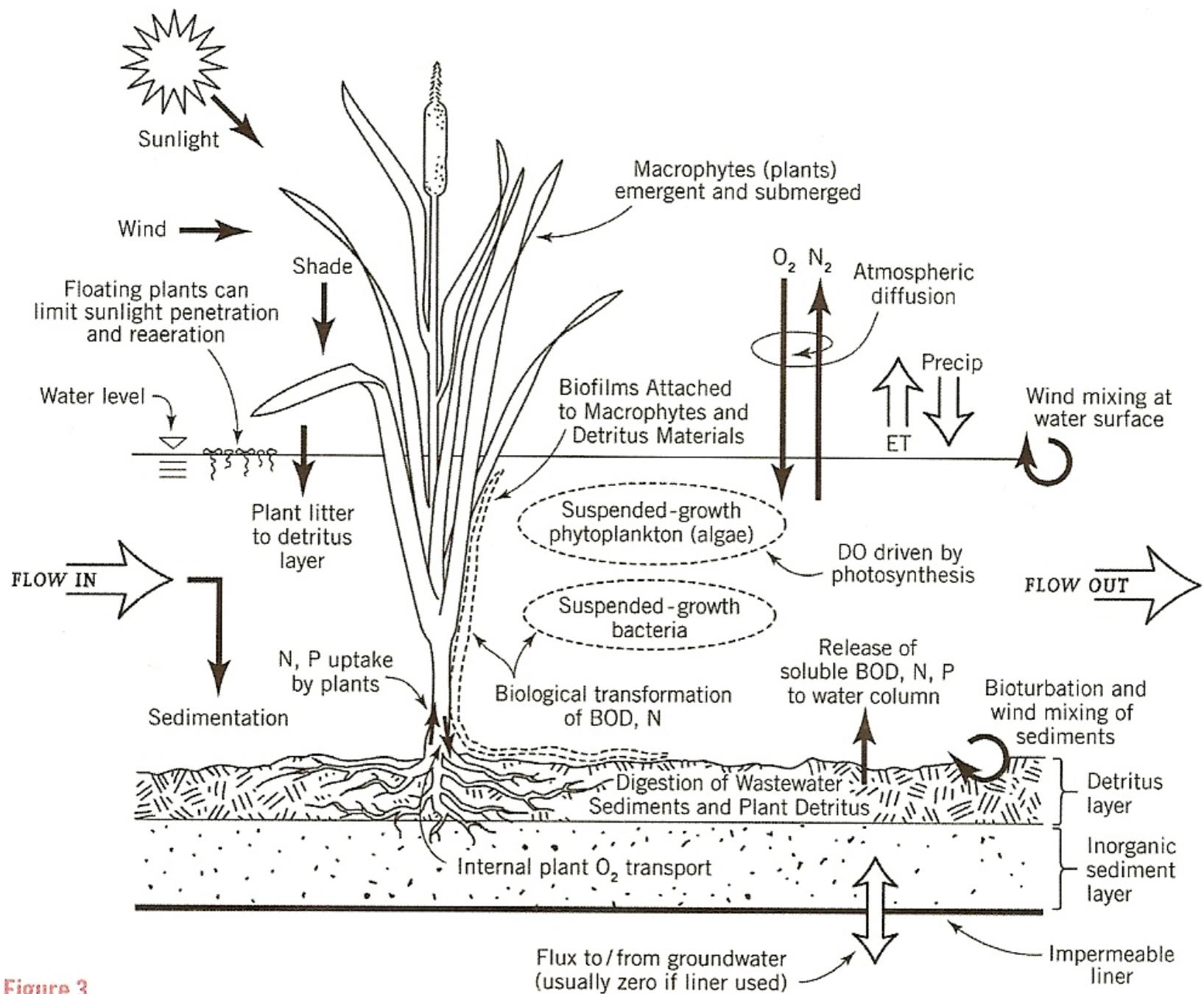


Figure 3

# Background

Antoine Chezy (1776)

$$\bar{u} = C\sqrt{Ri}$$

Von Karman

$$u(z) = \frac{u_*}{\kappa} \ln \left( \frac{z}{z_0} \right)$$

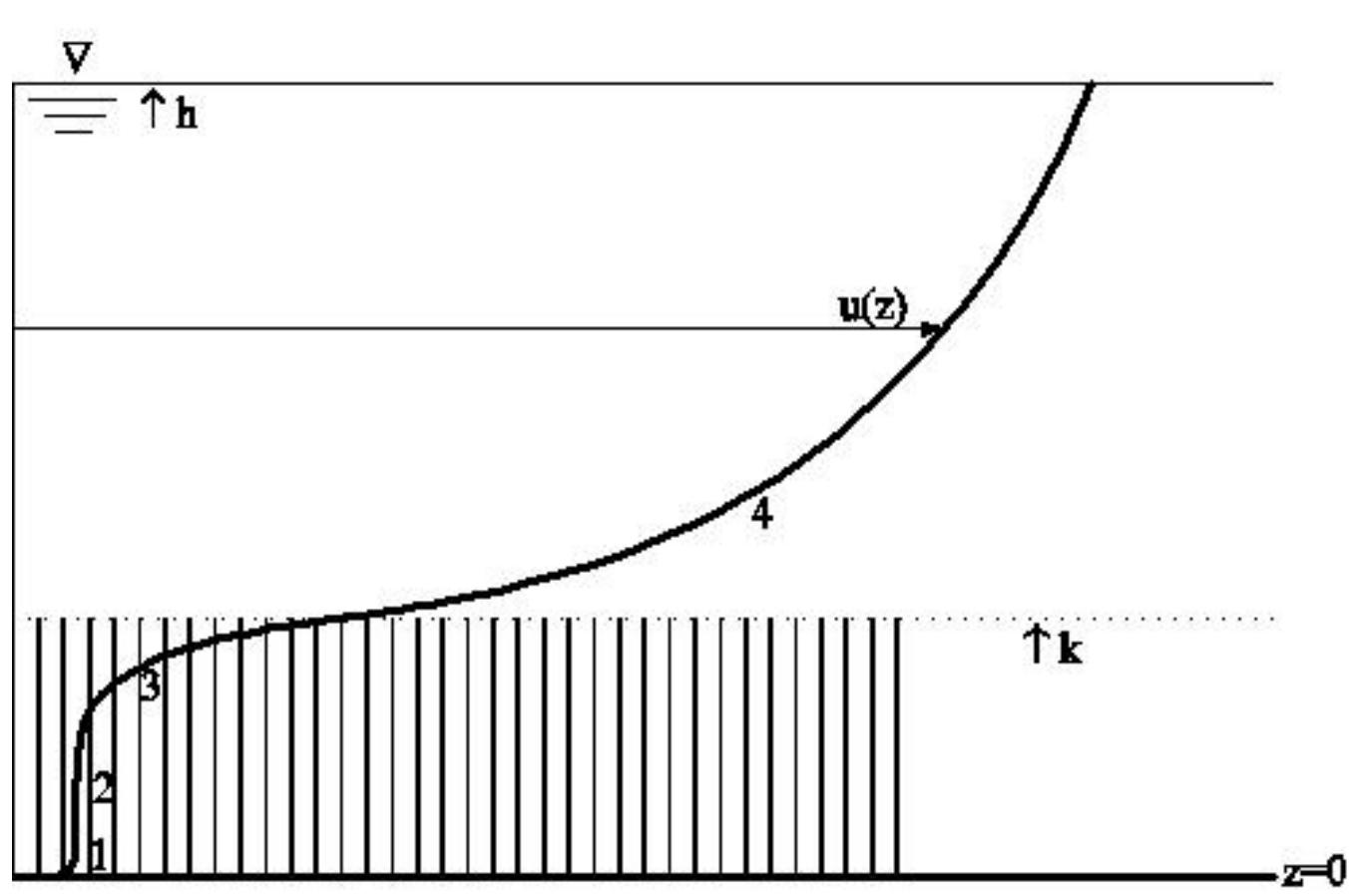
Nikuradze (1930)

$$z_0 = k_N / 30$$

White-Colebrook formula

$$C = 18 \log \left( \frac{12R}{k_N} \right)$$

# Vertical velocity profile



# Detailed Account of Roughness

1 DV Turbulence Model

Uittenbogaard (2003)

# Detailed description

## 1 DV Turbulence Model

$$\rho_0 \frac{\partial u(z)}{\partial t} + \frac{\partial p}{\partial x} = \frac{\rho_0}{1 - A_p(z)} \frac{\partial}{\partial z} \left( (1 - A_p(z)) (\nu + \nu_T(z)) \frac{\partial u(z)}{\partial z} \right) - \frac{F(z)}{1 - A_p(z)}$$

$$\frac{\partial u}{\partial x} = 0$$

$$\frac{\partial k}{\partial t} = \frac{1}{1 - A_p} \frac{\partial}{\partial z} \left( (1 - A_p) (\nu + \nu_T / \sigma_k) \frac{\partial k}{\partial z} \right) + T + P_k - B_k - \varepsilon$$

# Detailed description

## 1 DV Turbulence Model

Drag force

$$F(z) = \frac{1}{2} \rho_0 C_D(z) D(z) m(z) u(z) |u(z)|$$

Additional turbulence

$$T(z) = F(z) u(z)$$

Turbulence production

$$P_k = \nu_T \left( \frac{\partial u}{\partial z} \right)^2$$

Buoyancy

$$B_k = - \frac{\nu_T}{\sigma_k} \frac{g}{\rho_0} \frac{\partial \rho}{\partial z}$$

Dissipation rate

$$\frac{\partial \varepsilon}{\partial t} = \frac{1}{1 - A_p} \frac{\partial}{\partial z} \left( (1 - A_p) (\nu + \nu_T / \sigma_\varepsilon) \frac{\partial \varepsilon}{\partial z} \right) + P_\varepsilon - B_\varepsilon - \varepsilon_\varepsilon + c_{2\varepsilon} \frac{T}{\tau_{eff}}$$

# Analytical Formulations

Baptist (2005)

# Non-submerged vegetation

- Balance of horizontal momentum:

$$\tau_t = \tau_b + \tau_v$$

From which Chezy roughness is derived:

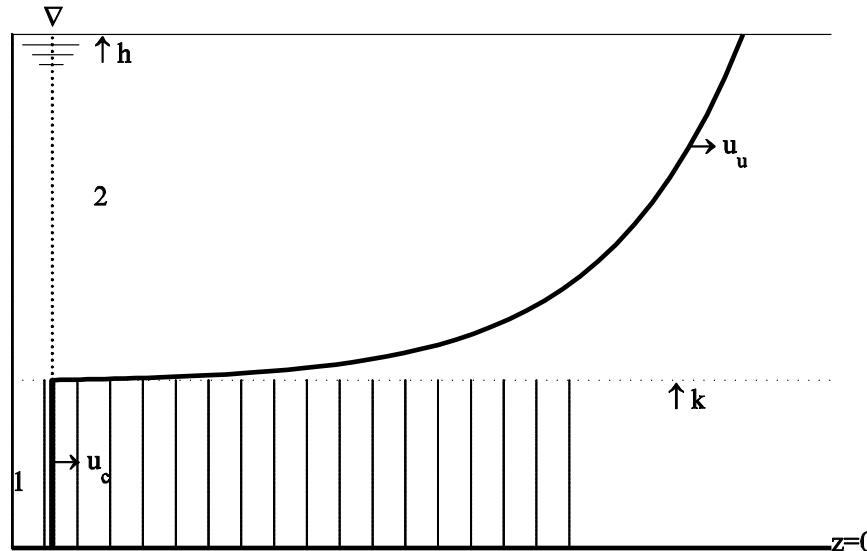
$$C_k = \sqrt{\frac{1}{\frac{1}{C_b^2} + \frac{C_D m D h}{2g}}}$$

or simplified

$$C_k = \sqrt{\frac{2g}{C_D m D h}}$$

# Submerged vegetation

Method of Effective Water Depth (Baptist, 2005)

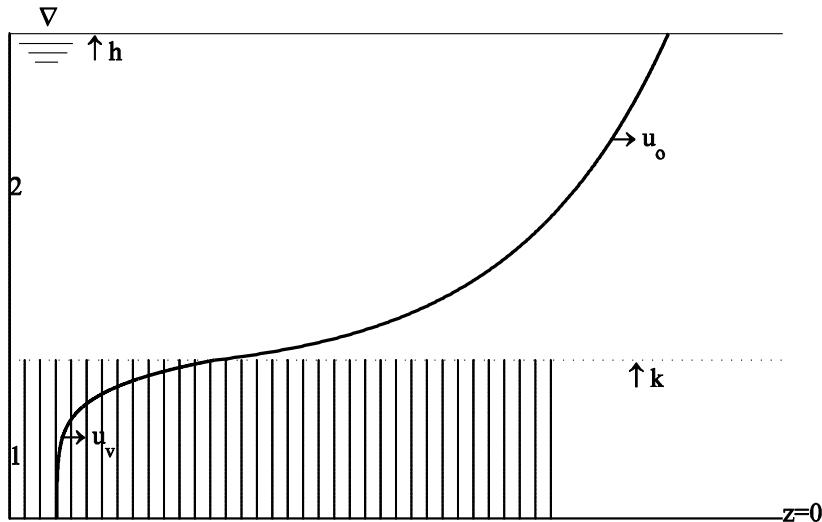


A

$$C_r = \sqrt{\frac{1}{C_b^{-2} + \frac{C_D m D k}{2g}}} + \frac{(h-k)^{3/2} \frac{\sqrt{g}}{\kappa} \ln\left(\frac{h-k}{e z_0}\right)}{h^{3/2}}$$

# Submerged vegetation

Analytical Method (Baptist, 2005)

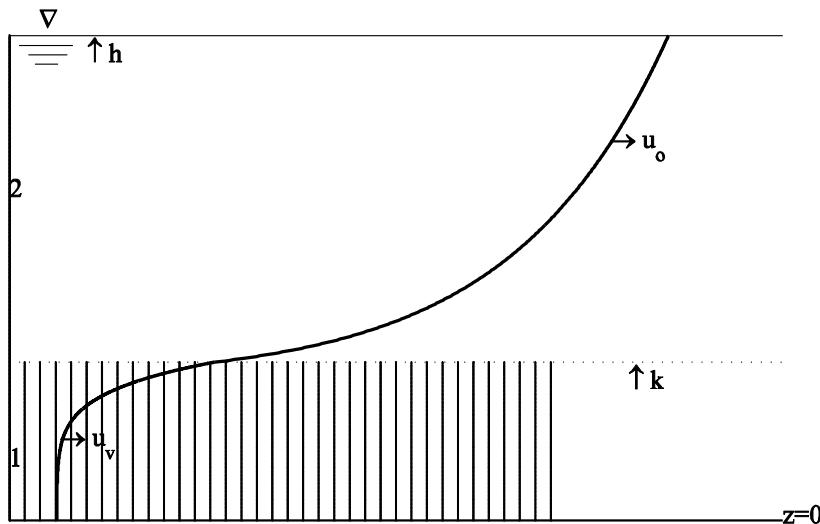


B

$$C_r = \frac{1}{h^{3/2}} \left\{ L \left[ 2 \left( u_{vk} - \sqrt{a_v + u_{v0}^2} \right) + u_{v0} \ln \left( \frac{(u_{vk} - u_{v0})(\sqrt{a_v + u_{v0}^2} + u_{v0})}{(u_{vk} + u_{v0})(\sqrt{a_v + u_{v0}^2} - u_{v0})} \right) \right] + \frac{\sqrt{g(h-k)}}{\kappa(h-k)} \left[ (h-d) \ln \left( \frac{h-d}{z_0} \right) - (k-d) \ln \left( \frac{k-d}{z_0} \right) - (h-k) \right] \right\}$$

# Submerged vegetation

Analytical Method (Baptist, 2005)



Where:

$$a = \frac{2Lg(h-k)i}{c_p l \exp\left(\frac{k}{L}\right)}$$

$$d = k - \int_0^k \frac{\exp\left(\frac{z}{L}\right)}{\exp\left(\frac{k}{L}\right)} dz = k - L \left(1 - \exp\left(-\frac{k}{L}\right)\right)$$

$$z_0 = (k - d) \exp\left(-\kappa \sqrt{\frac{2L}{c_p l} \left(1 + \frac{L}{h-k}\right)}\right)$$

$$c_p = \frac{1}{20} \frac{h-k}{l}$$

# Symbolic Regression

Knowledge Discovery in Petabyte Age

# The Age of Information Overflow

- Development of technologies for data acquisition, storage and transmission continues at a breakneck pace
- The same cannot be asserted about extraction of useful information and knowledge from data
- *What is to be done with all these data?*

# Data Mining for Scientific Discovery

- Philanthropic activity (without major material satisfaction)
- Every scientific discipline holds a vast amount of knowledge that should not be just ignored (black boxes are not tolerated)
- Combination of theory-driven (understanding-rich) and data-driven discovery process

# 3rd Kepler's Law of Planetary Motion

Tycho Brahe collected observables

- READINGS OF THE BOOK OF NATURE ITSELF

Kepler analysed the data

- MODEL OF THE PROCESS - ECONOMY OF THOUGHT

# Evolutionary Algorithms

- ★ Loosely based on Darwin's theory of evolution
- ★ Survival of the fittest
- ★ Genetic operators
  - ★ Reproduction (crossover)
  - ★ Mutation
  - ★ Selection

# Evolutionary Criteria

- Criterion of Heredity

Offspring are similar to their parents: the copying process maintains high fidelity

- Criterion of Variability

Offspring are not exactly the same as their parents or each other: the copying process is not perfect

- Criterion of Fecundity

Variants leave different number of offspring: specific variations have an effect on behaviour, and behaviour has an effect on reproductive success

# Evolutionary Algorithms

## Crossover

- Select two parental strings according to the fitness criterion
- Randomly select crossover site
- Swap the two sub-strings

ATTCGGCTTAAG

ATACGCGCTACGC

ATACGC CTTAAG  
ATT CGG CTACGC

ATACGC CTTAAG

ATT CGG CTACGC

# Evolutionary Algorithms

## Crossover (1)

THIS IS THE WEST MESSAGE

THIN IS THE BEST MESSAGE

# Evolutionary Algorithms

## Crossover (2)

THIS IS THE WEST MESSAGE

THIN IS THE BEST MESSAGE

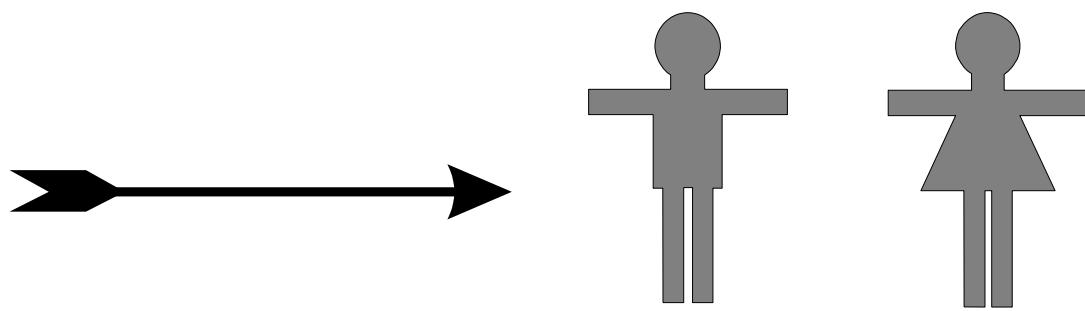
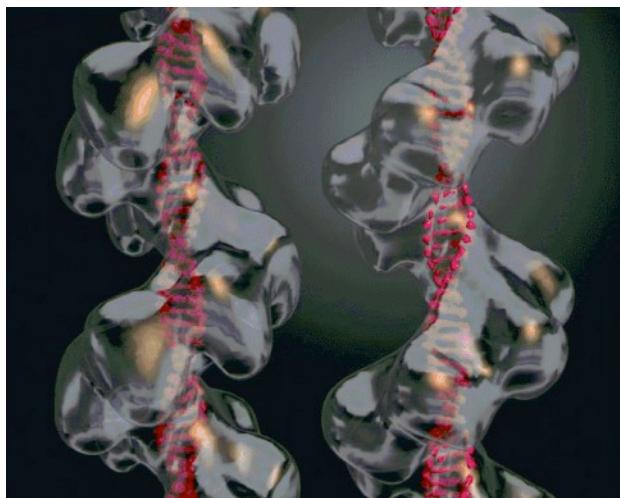
# Evolutionary Algorithms

## Crossover (3)

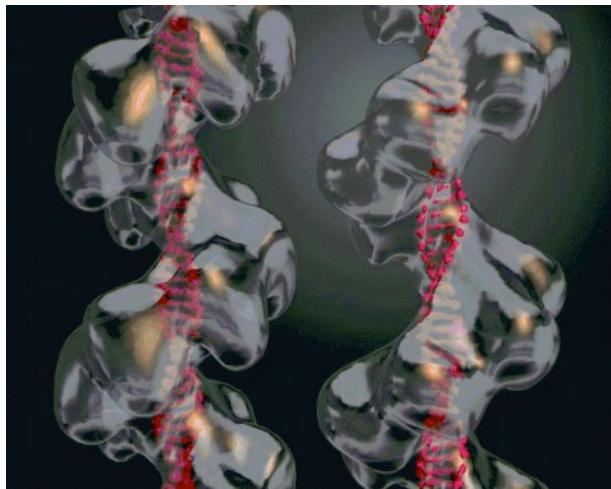
THIS IS THE NEWEST MESSAGE  
THIN IS THE BEST MESSAGE

THIS IS THE BEST MESSAGE  
THIN IS THE NEWEST MESSAGE

# Natural Evolution



# Genetic Programming



$$\frac{\partial x}{\partial t} + \frac{\partial y}{\partial t} + \frac{\partial z}{\partial t} = 0$$

$$\sin \theta + \ln \lambda - \sum_{t=0}^{\rho} x^t$$

# Genetic Programming

## Crossover (1)

$$\ln(x) + \sin y / z$$
$$(\cos(x) + \sqrt{z}) / 3$$

# Genetic Programming Crossover (2)

$$\ln(x) + \sin(y) / z$$
$$(\cos(x) + \sqrt{z}) / 3$$

# Genetic Programming

## Crossover (3)

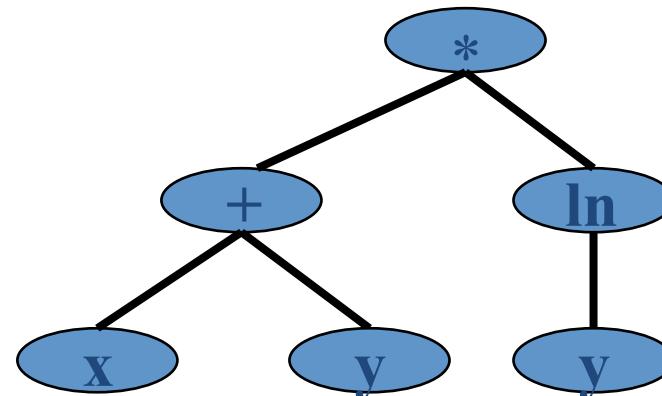
$$\ln(x) + \sin y / z$$
$$(\cos(x) + \sqrt{z}) / 3$$

$$\ln(x) (x + \sqrt{z}) / 3$$
$$(\cos + \sin y / z$$

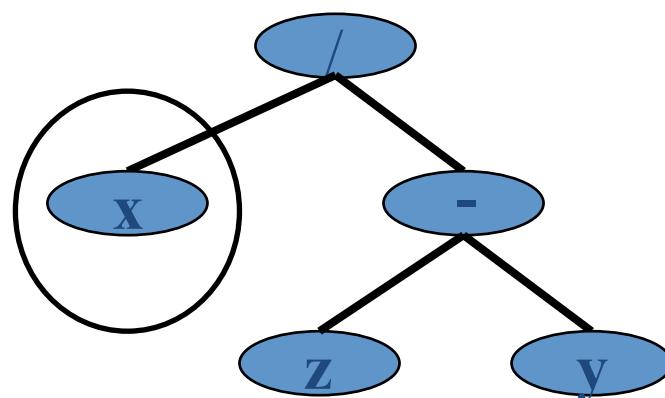
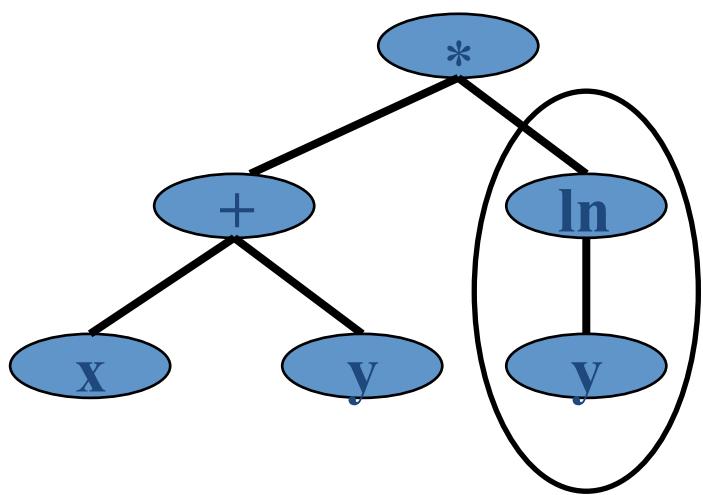
# Genetic Programming

- Parse Trees
  - Function Set
  - Terminal Set
  - Constants
- Symbolic Regression
  - Find a relationship between:
  - dependent
  - independent variables

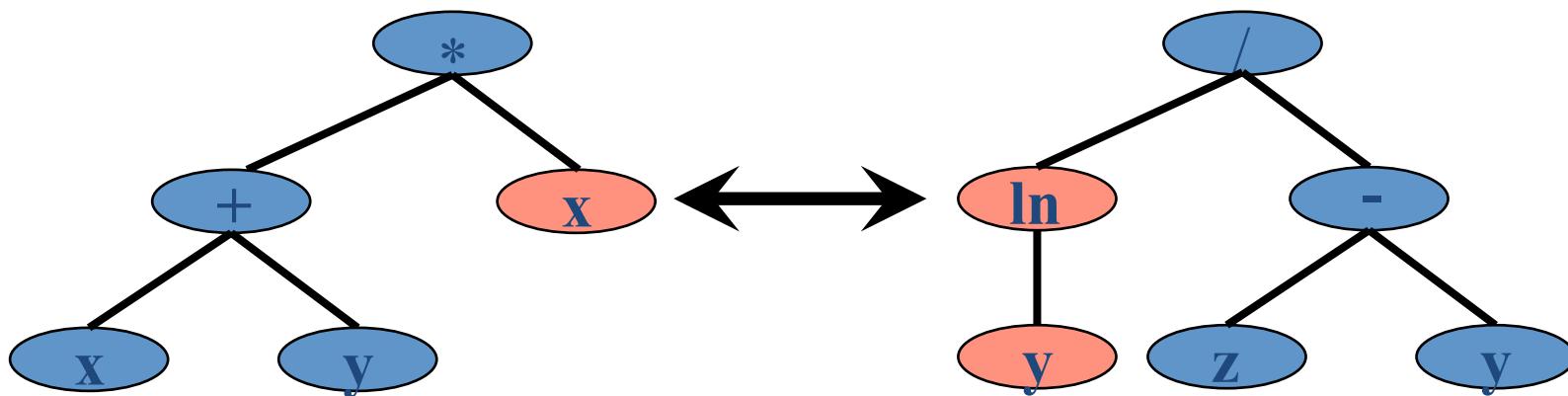
$$(x + y) \ln(y)$$



# Breeding: Crossover



# Breeding: Crossover



# Evolutionary Algorithms

## Mutation

- Randomly select parental string
- Randomly alter one of the symbols

THISIS THE WES T MESSAGE

THISIS THE WES R MESSAGE

# Genetic Programming

## Mutation

### ① Branch-mutation

- » a complete sub-tree is replaced with another arbitrary sub-tree

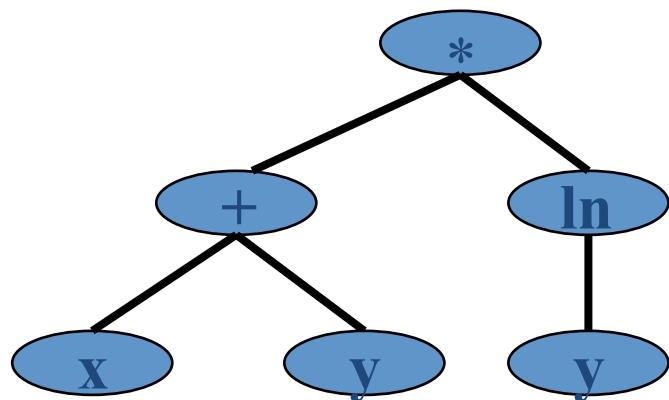
### ② Node-mutation

- » a random change to a single node with another, arbitrary node of a same arity

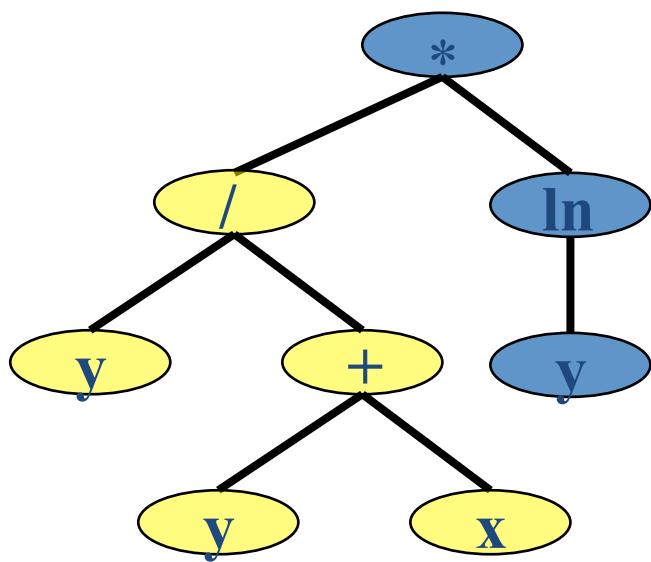
### ③ Inversion-mutation

- » inversion of an order in which operands are ordered in an expression, for example  $f(x, y, z)$  may become  $f(z, x, y)$

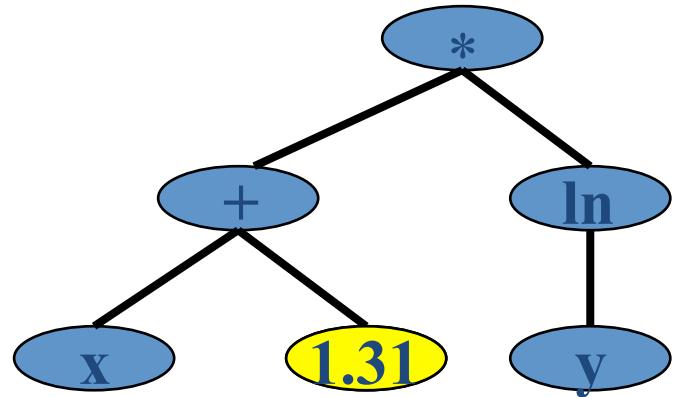
# Breeding: Branch Mutation



# Breeding: Branch Mutation



# Breeding: Constant Mutation



- White noise
- Small uniform steps
- Occasional large jumps
- Real vector optimization (gauss-Newton, evol. Strategies)

# Genetic Programming Selection

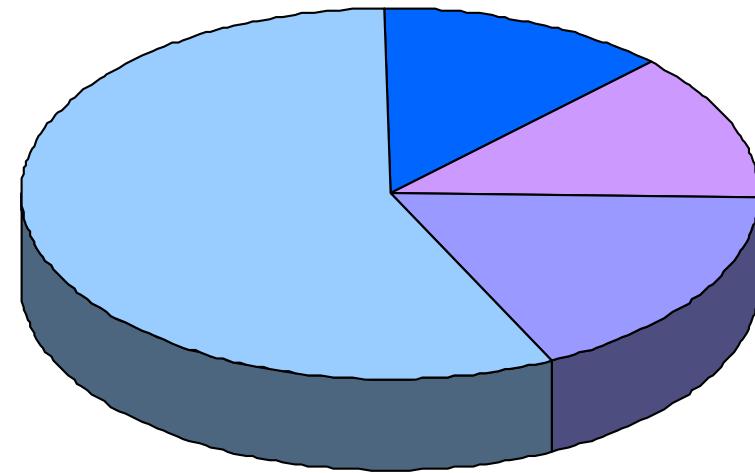
Gives direction to evolutionary search process

Favours more fit entities by selecting them more often

‘Survival of the fittest’

# Genetic Programming Selection

- Biased Roulette Wheel
- Areas are proportional to individual fitnesses
- Selection is random, but still promoting more fit individuals
- Balance between homogenising selection and generation of variation



# Genetic Programming

## Fitness

x	y	z	$z = \frac{x^2}{y}$	$z = x^3 + y^2$	$z = x^2 + y$	Calculated Z
15.0	45.0	1289.0	5.0	5400.0	2250.0	
35.0	32.0	354.0	38.3	43899.0	2249.0	
0.0	8.00	51381.0	0.0	64.0	64.00	
5.0	51.0	518.0	0.49	2726.0	2626.0	

ERRORS:

21837.0 101181.0 56281.0

FITNESSES:

6.67E-06 1.66E-07 6.12E-07

‘Fitness’ can be any measure (RMSE, NMSE, Classification Accuracy)

# Evolutionary Algorithms

## pseudo code

- 1.) Randomly create an initial population
- 2.) Iteratively perform the following:
  - (a) evaluate fitness of each individual
  - (b) select ‘parents’ on the basis of fitness
  - (c) perform (sexual) recombination
  - (d) mutate
- 3.) Best so-far individual

# Why Symbolic Regression?

- Approximate a relationship  $y = f(x)$  by minimizing error where  $f$  is built of a number of primitive functions
- Other regression methods require pre-specified functional form
  - Linear regression
  - Nonlinear regression (including ANN and SVM)
  - Coefficient fitting
- GP manipulates symbolic information
  - Hypothesis: GP can lead to *interpretable* equations

# Caveats

- When uncontrolled
  - Tree size grows
  - Referred to as bloat
- Possible solutions
  - Parsimony Pressure
  - Short formulations higher fitness than long
  - Length as penalty
  - Length as second objective
  - ...
- (((((0.18916 / 0.098799) -  
sqrt((((x\_1 / 0.22097) -  
sqrt(0.22097)) -  
sqrt(sqrt((((x\_1 / 0.22097) -  
sqrt(0.019076)) - sqrt(0.28856)) +  
x\_1) - sqrt(0.22097)))) + 0.22097)  
- sqrt(0.22097))) + -0.034116) -  
sqrt(x\_1)) + ((((((0.017471 -  
0.071534) - sqrt((((x\_1 / 0.22097)  
- sqrt(0.22097)) - sqrt(0.28856)) +  
x\_1) - sqrt(0.22097)))) + x\_1) -  
sqrt(x\_1)) + (((((x\_1 / 0.22097) -  
sqrt(x\_1)) - sqrt(x\_1)) + x\_1) -  
sqrt(0.017471))) - sqrt((((x\_1 /  
0.22097) - sqrt(0.019076)) -  
sqrt(0.28856)) + x\_1) -  
sqrt(0.22097))) - sqrt(x\_1)) + x\_1)  
- 0.017471))
- Is there an approach one can follow in order control this and similar problems?

# Genetic Programming in Natural Sciences

Can we utilise already existing knowledge in order to enhance interpretability and avoid bloat?

# Genetic Programming

- Programming computers by telling them what to do, instead telling them how to do it
- But, how do we say to computer what we want?
- And what is that we want anyway?

# What is that we want from a learning machine?

Find a good fit, do not worry about inconsistencies

- Transform the problem to the one without units of measurements (Buckingham's Pi theorem) and then find a good fit

Find a good fit, do not accept incorrect or inconsistent solutions

Find a good fit with least possible inconsistencies

# Find a good fit only

- Standard GP, neural networks, polynomial fits, splines, SVM, RVM...
- Problem is posed in either raw or dimensionless form
- If we are lucky we can interpret GP expressions
- In similar sense one can attempt to interpret other approximations

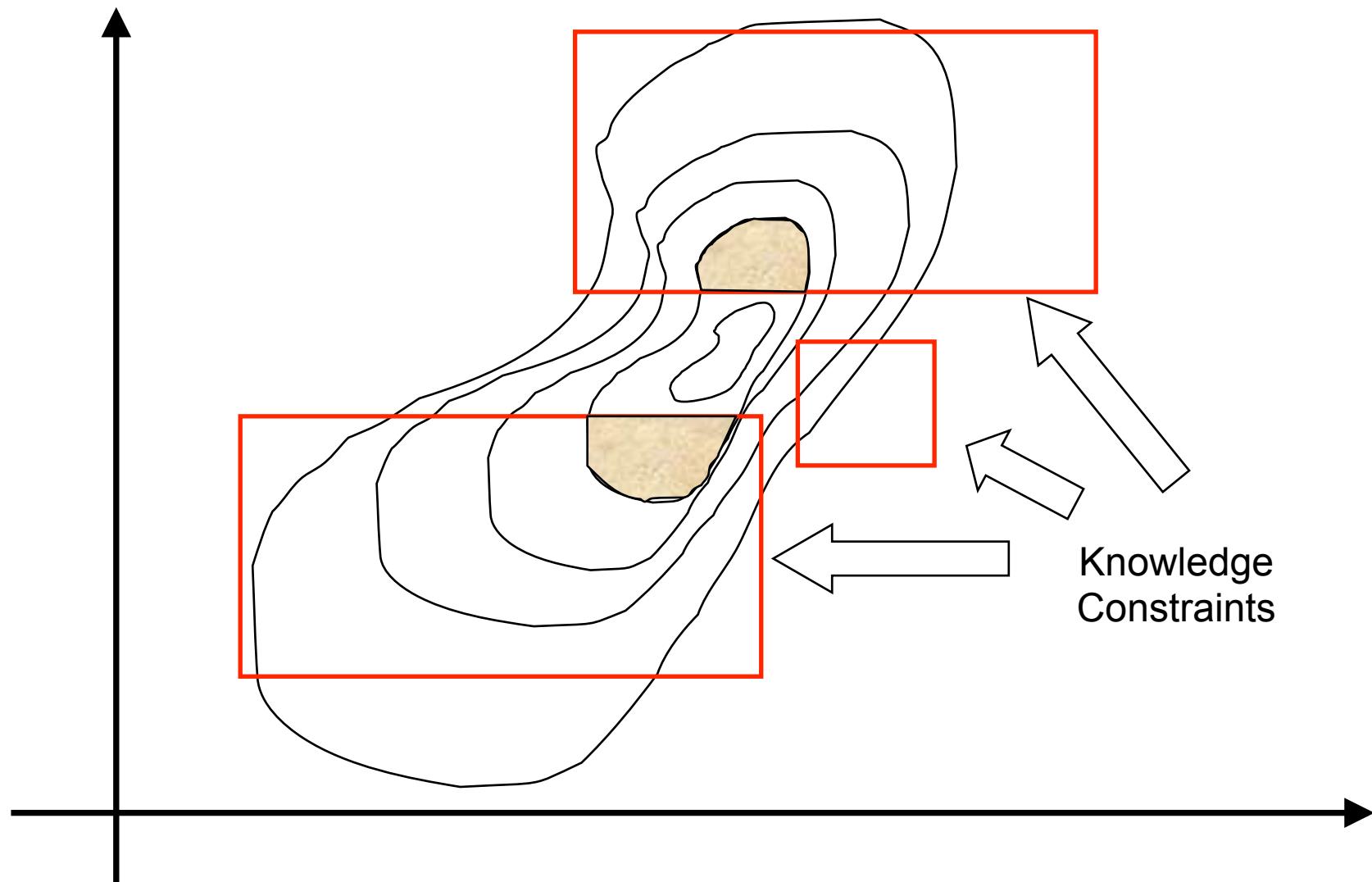
# Check for inconsistencies using background knowledge - Bias

- Bias refers to any kind of basis for choosing one generalisation over another, other than strict consistency with the observed training examples
- Bias is based on knowledge about the domain
- Bias can be provided in strong (declarative) and weak (preferential) form

# Declarative Bias

- Strongly Typed GP (STGP)
- Declarative bias takes form of constraints and strict adherence with constraints is fully enforced
- Violation of constraints eliminates proposed equation altogether (all or nothing)
- Indispensable for well defined problems
- Dramatically reduces search space
- But, it also breaks down search space resulting in a loss of diversity

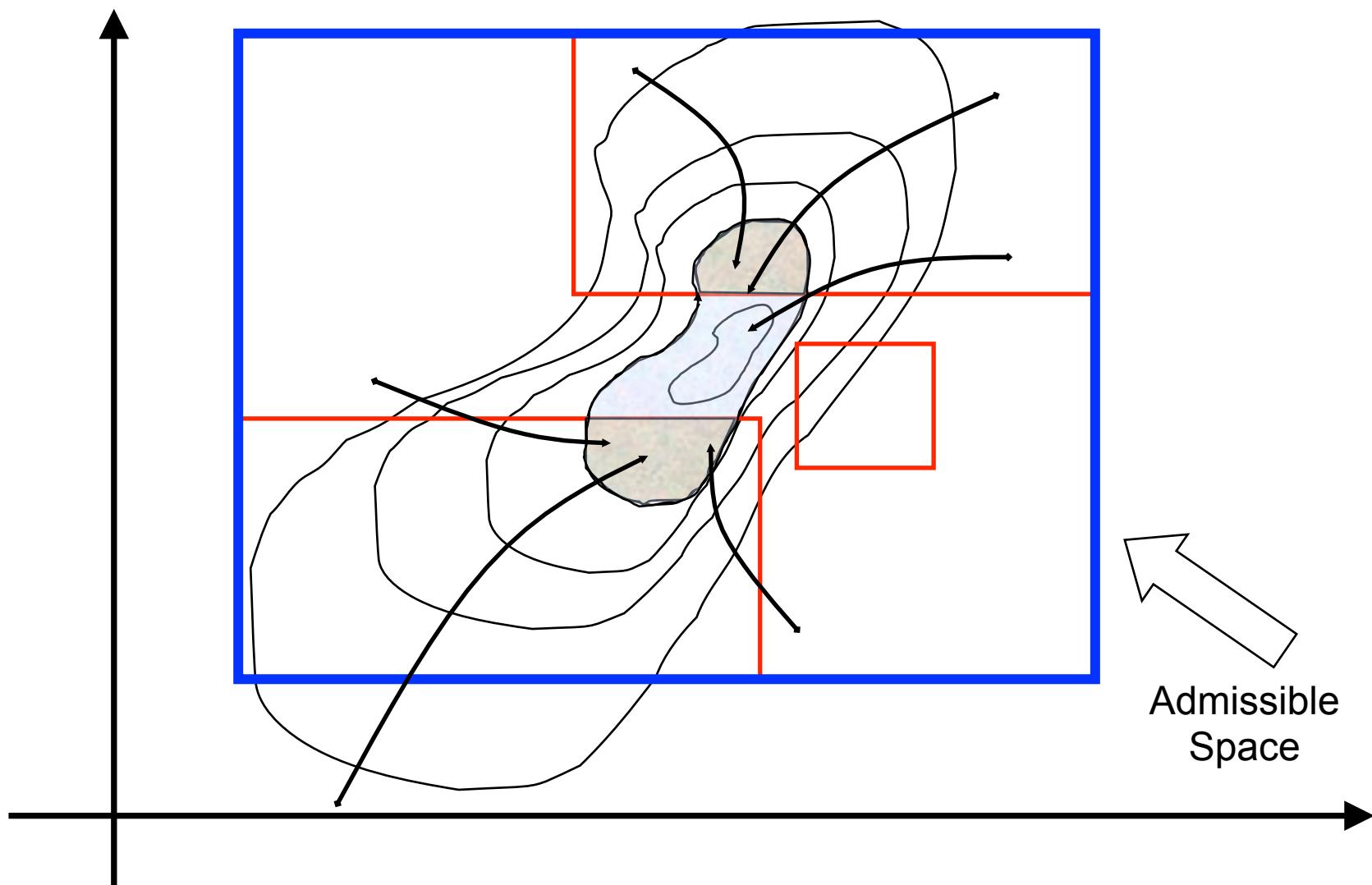
# Broken ergodicity



# Preferential bias

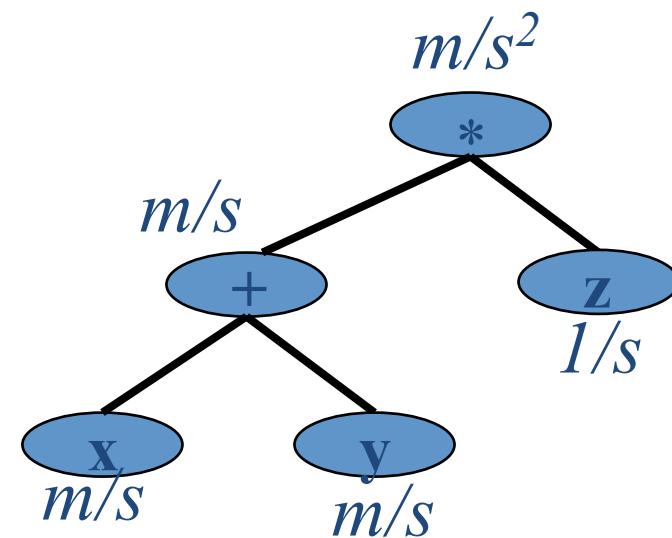
- Introduces knowledge in a form of objective function(s)
- It does not guarantee strict adherence with background knowledge
- It only defines preference towards equations which more closely satisfy requirements imposed by background knowledge

# Preferential bias

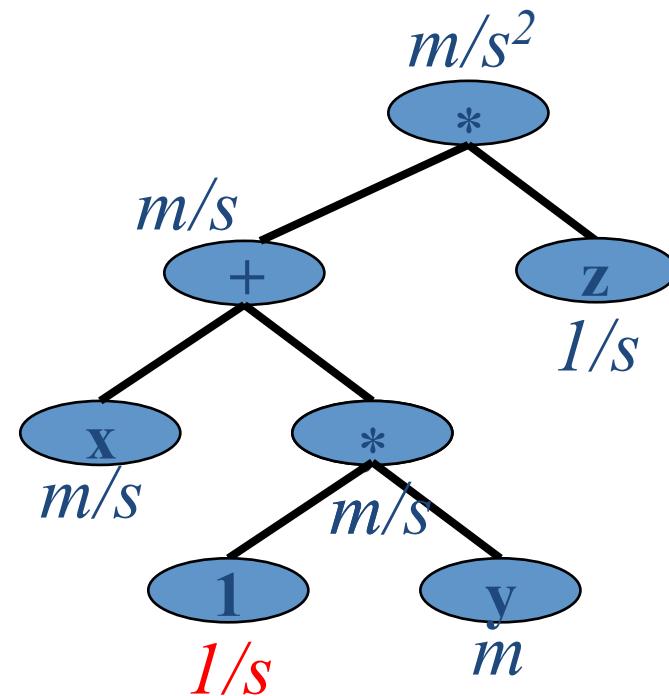


# Bias in the form of units of measurements

- Why restrict to dimensionless values?



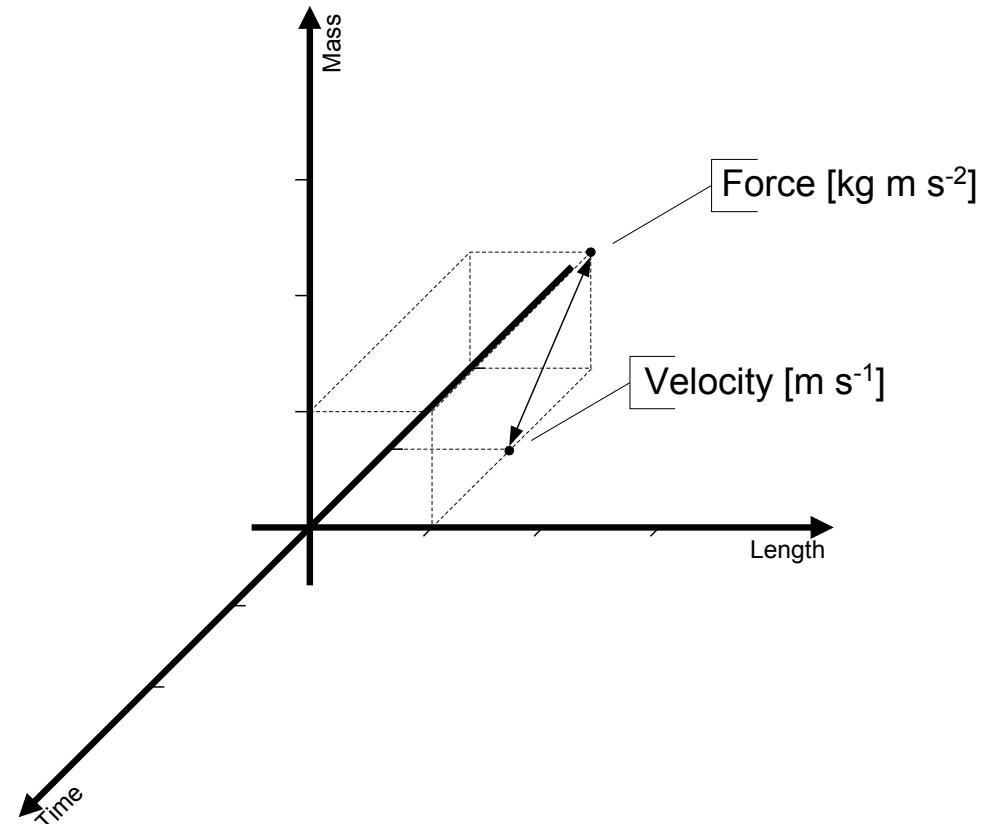
# Coercion



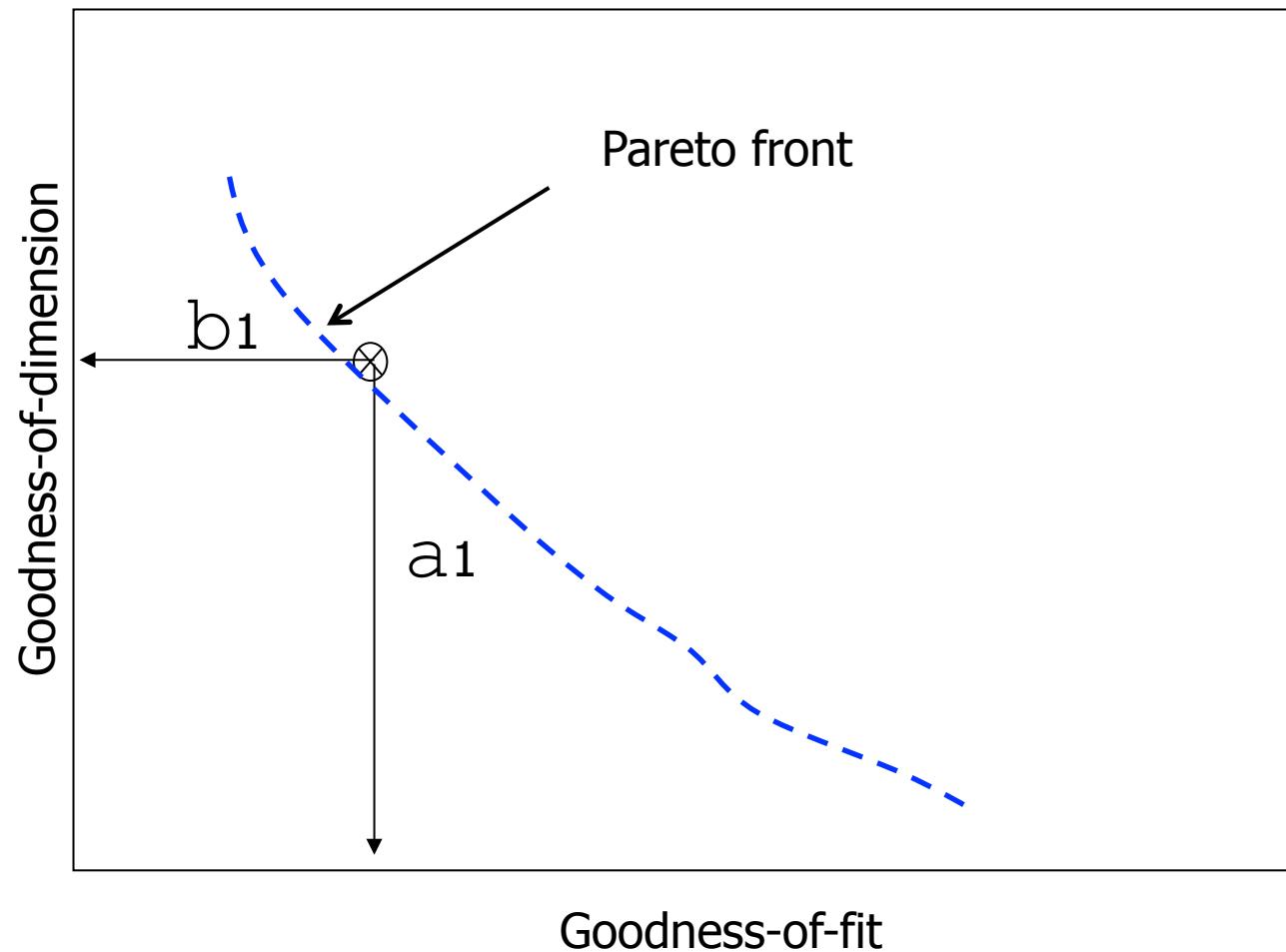
# Dimensionally Aware GP

(Babovic & Keijzer 1998)

- Counts number of *coercions* needed
  - multiplication with unity of needed transformation
- Pareto Balance
  - Goodness of fit
  - Goodness of dimension



# Dimensionally Aware GP: the idea



# What You Get

- Judge which hypothesis balances the objectives best
  - Rearrange terms
  - Set up ‘little theory’
  - Use results to gain insight

# Back to Roughness - Data

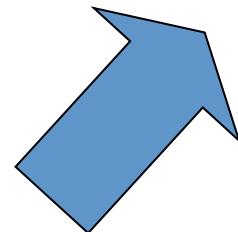
990 'observations' generated by the 1 DV model turbulence model  
(75% training, 25% validation)

<i>Input</i>	<i>Dimension</i>	Description
<b>D</b>	<b>L</b>	<b>Diameter of the stems.</b>
<b>m</b>	<b>L<sup>-2</sup></b>	<b>Number of stems per square meter</b>
<b>k</b>	<b>L</b>	<b>Vegetation height.</b>
<b>C<sub>D</sub></b>	<b>-</b>	<b>Drag coefficient of a single stem.</b>
<b>C<sub>b</sub></b>	<b>L<sup>0.5/T</sup></b>	<b>Bed Chézy resistance coefficient.</b>
<b>h</b>	<b>L</b>	<b>Water depth.</b>

177 flume experiments data points (out of sample validation)

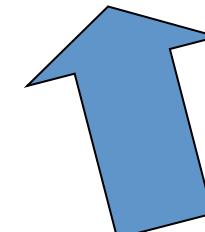
# GP Result

$$C_r = \sqrt{\frac{2g}{c_D m D k}} + 2\sqrt{g} \ln\left(\frac{h}{k}\right)$$



Non submerged  
conditions

C



Logarithmic term

# The fun is only starting

$$C_r = \sqrt{\frac{2g}{c_D m D k}} + 2\sqrt{g} \ln\left(\frac{h}{k}\right) \quad \mathbf{D}$$

$$C_r = \sqrt{\frac{1}{\frac{1}{C_b^2} + \frac{1}{2g} C_D m D k}} + \frac{\sqrt{g}}{\kappa} \ln\left(\frac{h}{k}\right)$$

# Differential Form

$$C_r = \sqrt{\frac{1}{\frac{1}{C_b^2} + \frac{1}{2g} C_D m D k}} + \frac{\sqrt{g}}{\kappa} \ln\left(\frac{h}{k}\right)$$

$$C_r = C_1 + \frac{\sqrt{g}}{\kappa} \ln\left(\frac{h}{k}\right)$$

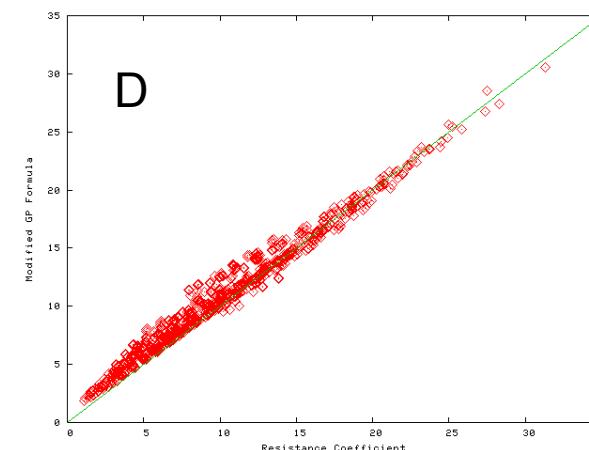
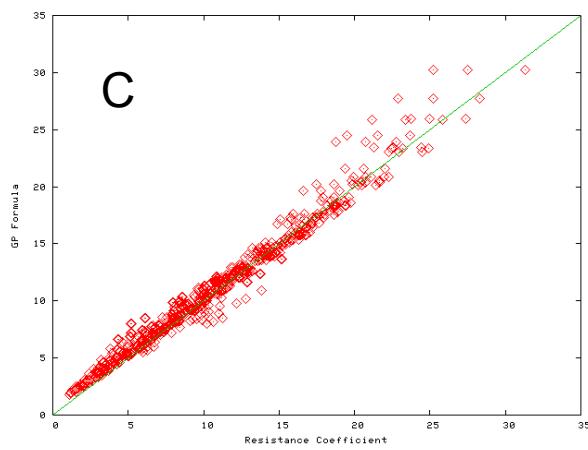
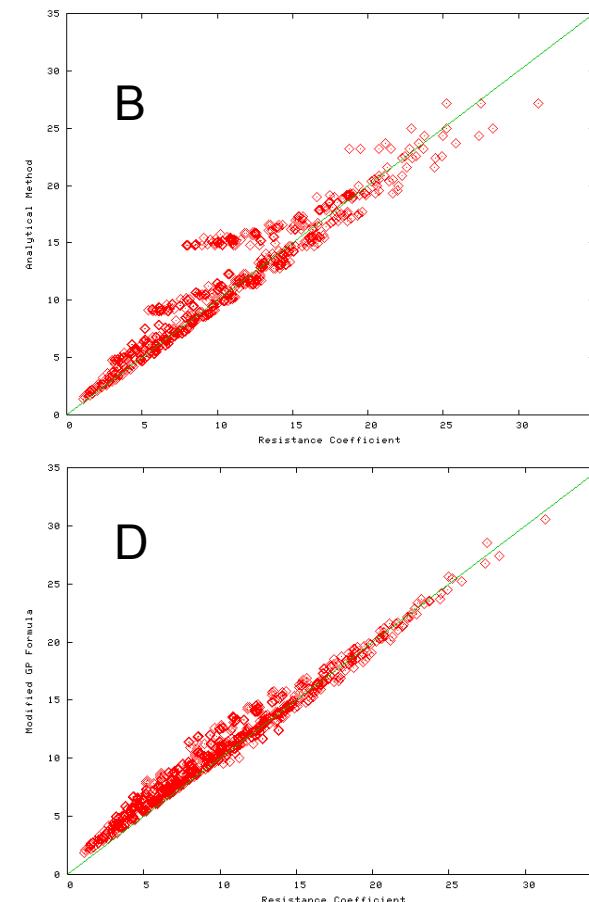
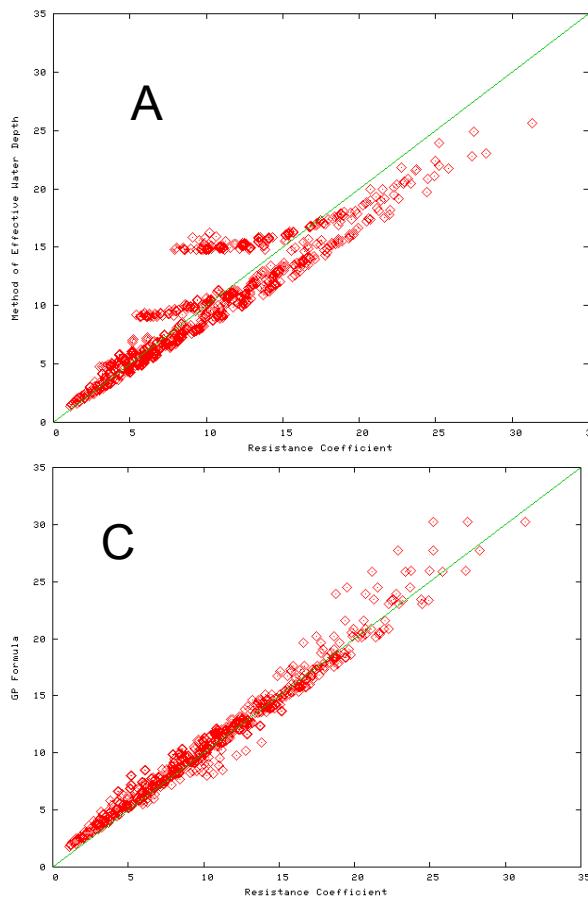
$$\frac{\partial C_r}{\partial h} = \frac{\sqrt{g}}{\kappa h}$$

# How Good is All This?

<i>Equation</i>	<i>RMSE 1-DV Data [m<sup>1/2</sup>s<sup>-1</sup>]</i>	<i>CoD 1-DV Data[-]</i>
Method of effective water depth (A)	2.2	0.831
Analytical solution method (B)	2.06	0.892
Original GP-formula (C)	0.98	0.971
Modified GP-formula (D)	1.13	0.977
1-DV numerical model	0	1.00

1 DV validation data

# How Good is All This?



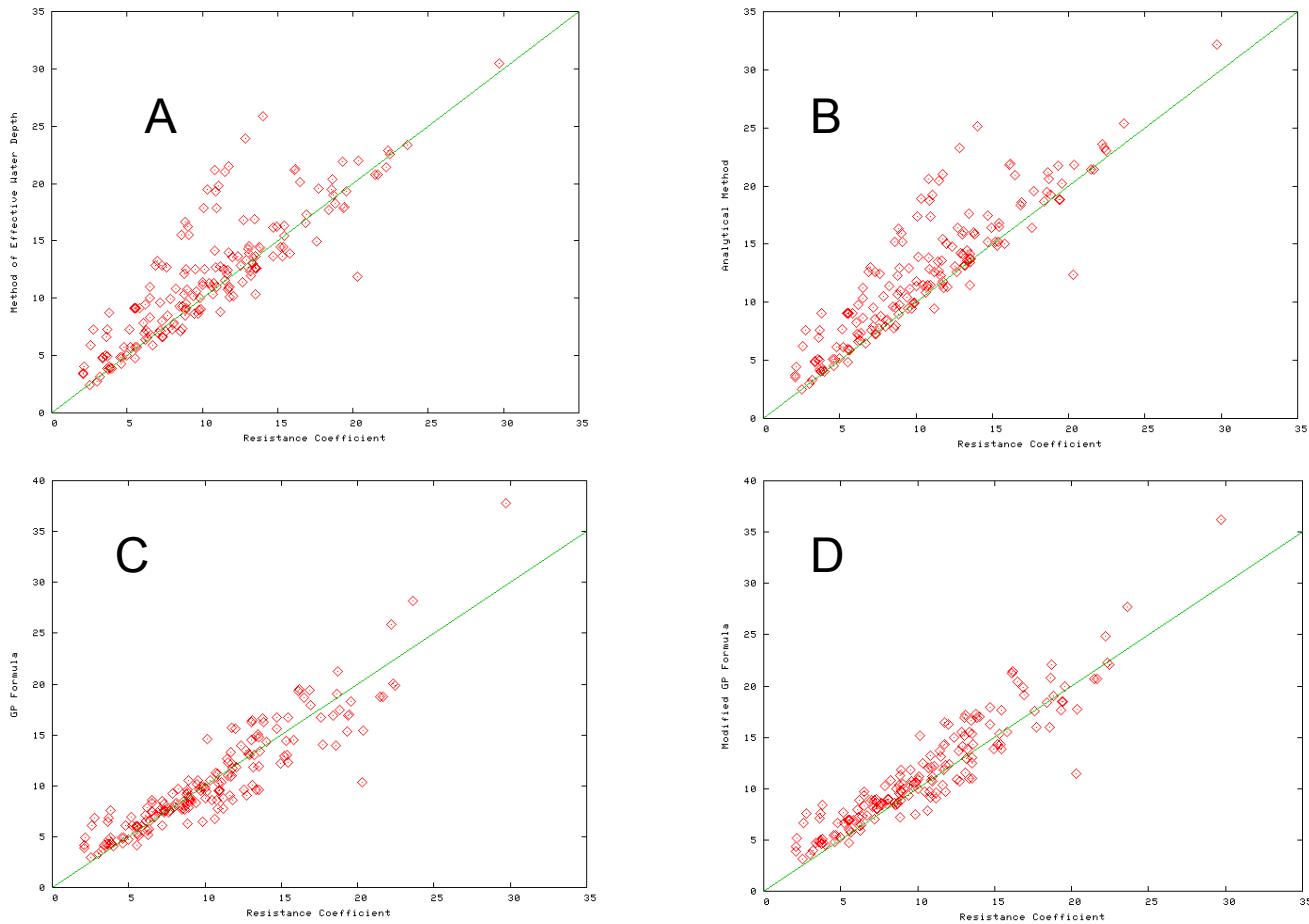
1 DV validation data

# How Good is All This?

<i>Equation</i>	<i>RMSE 1-DV Data [m<sup>1/2</sup>s<sup>-1</sup>]</i>	<i>CoD 1-DV Data[-]</i>
Method of effective water depth (A)	3.17	0.734
Analytical solution method (B)	3.15	0.787
Original GP-formula (C)	2.11	0.838
Modified GP-formula (D)	2.11	0.872
1-DV numerical model	1.86	0.873

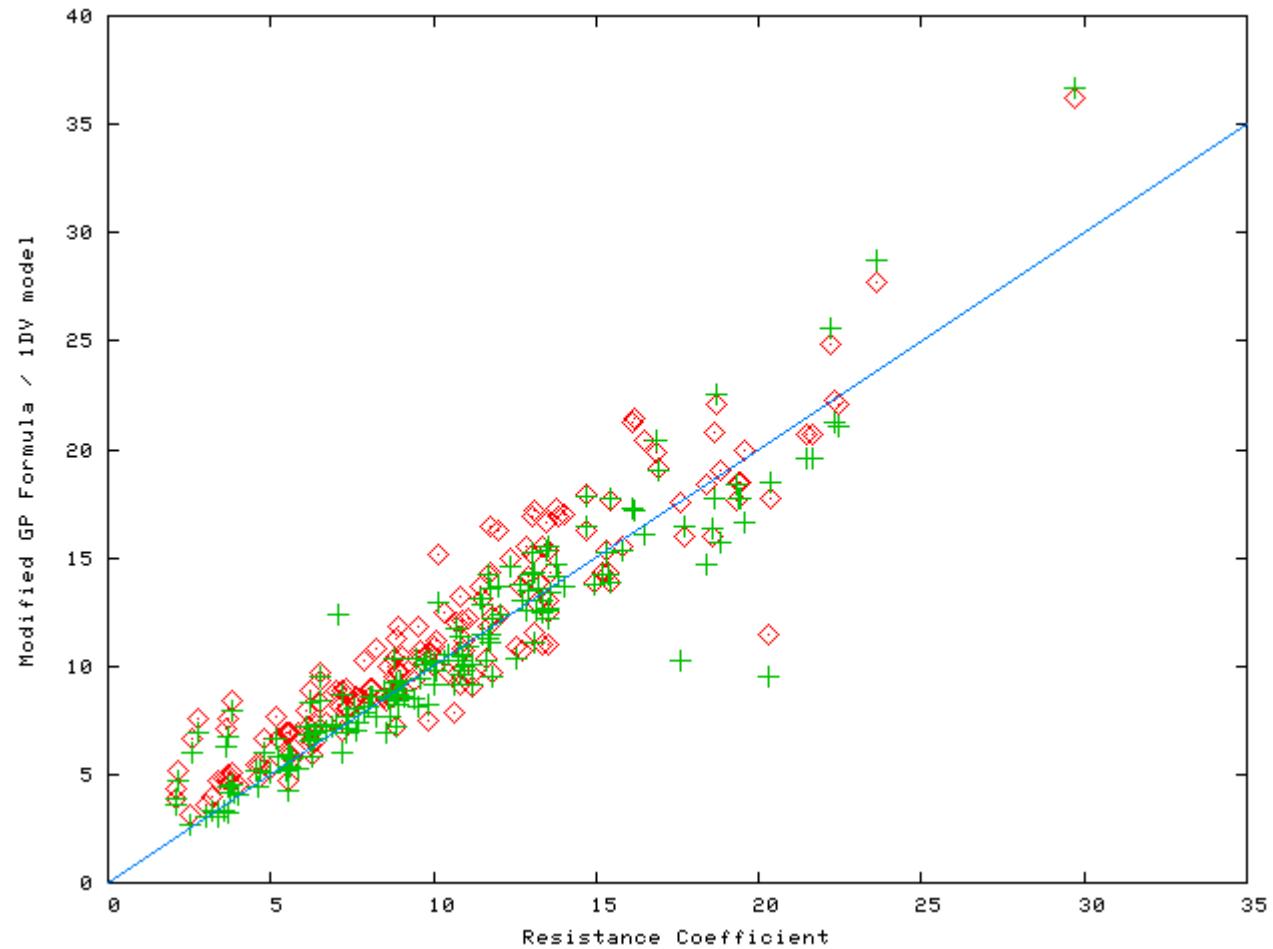
Flume data

# How Good is All This?

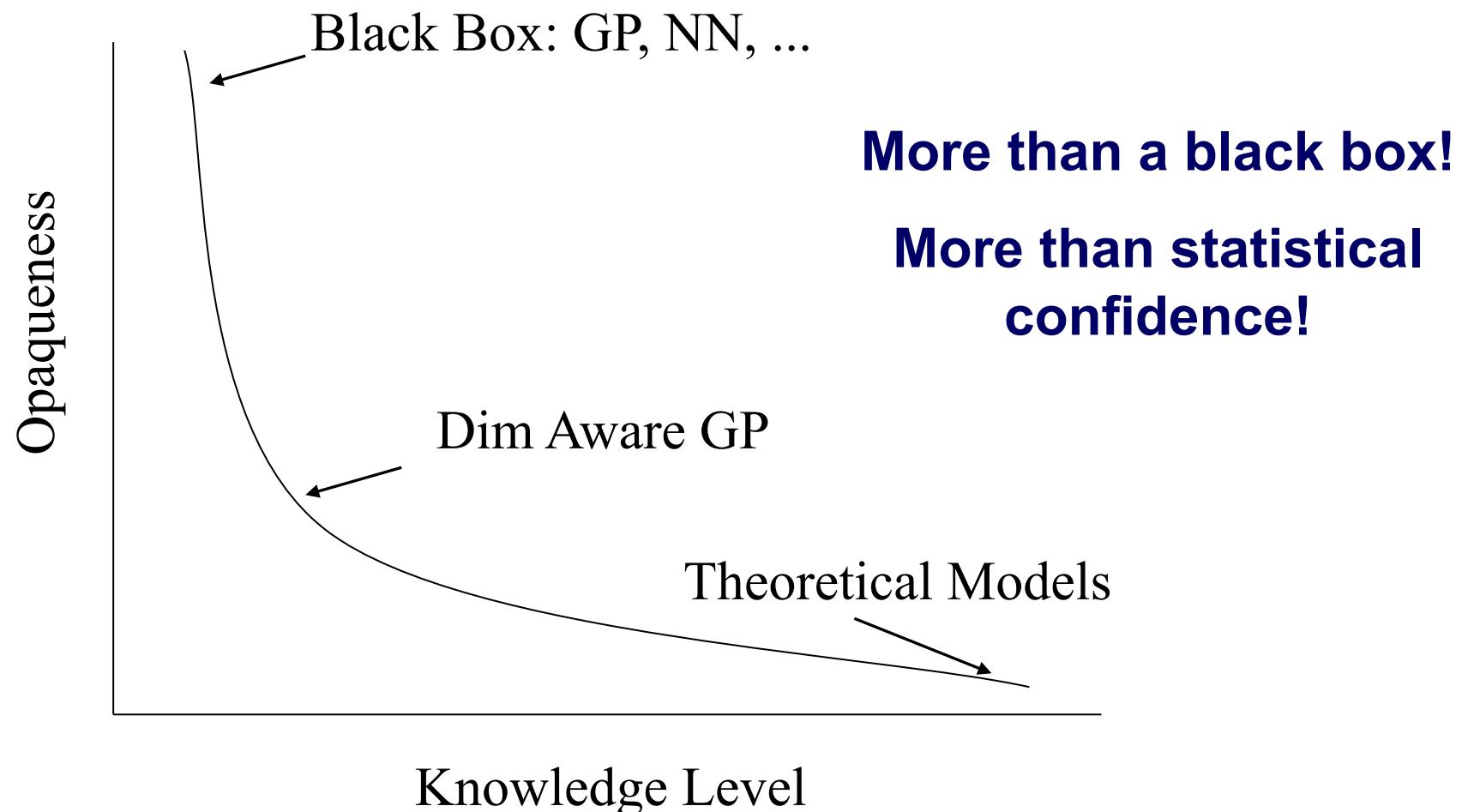


Flume data

# Eq. D vs. 1DV



# Model Induction Trade-Off



# Computational Scientific Discovery?

- Scientific Discovery represents the pinnacle of human creative thought
- Is it likely candidate for automation by computer?
- Certainly yes!!!

# How far can we push...

- Evolution of Systems of Coupled Ordinary Differential Equations

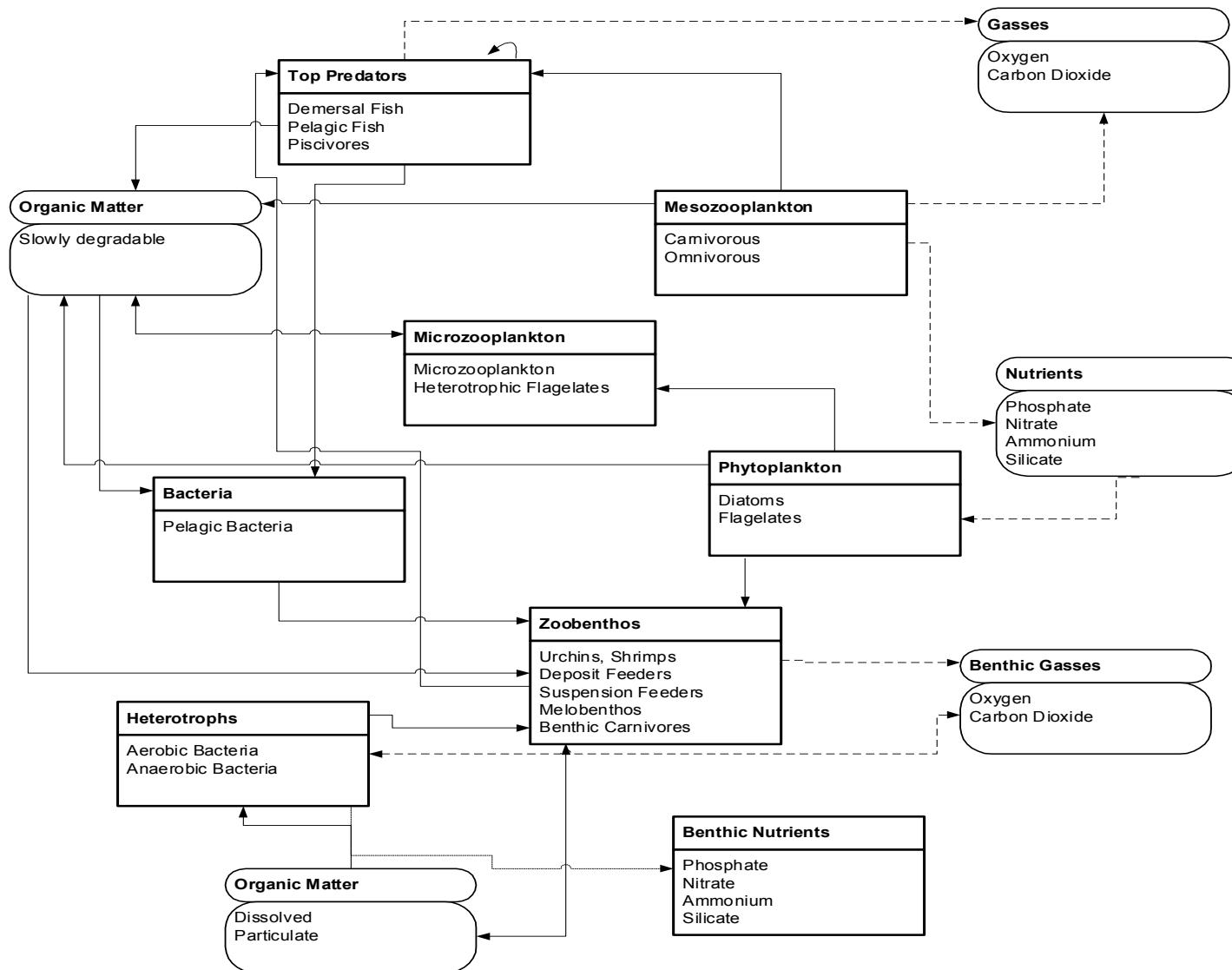
$$\frac{dx_1}{dt} = f_1(t, x_1, x_2, \dots, x_n)$$

$$\frac{dx_2}{dt} = f_1(t, x_1, x_2, \dots, x_n)$$

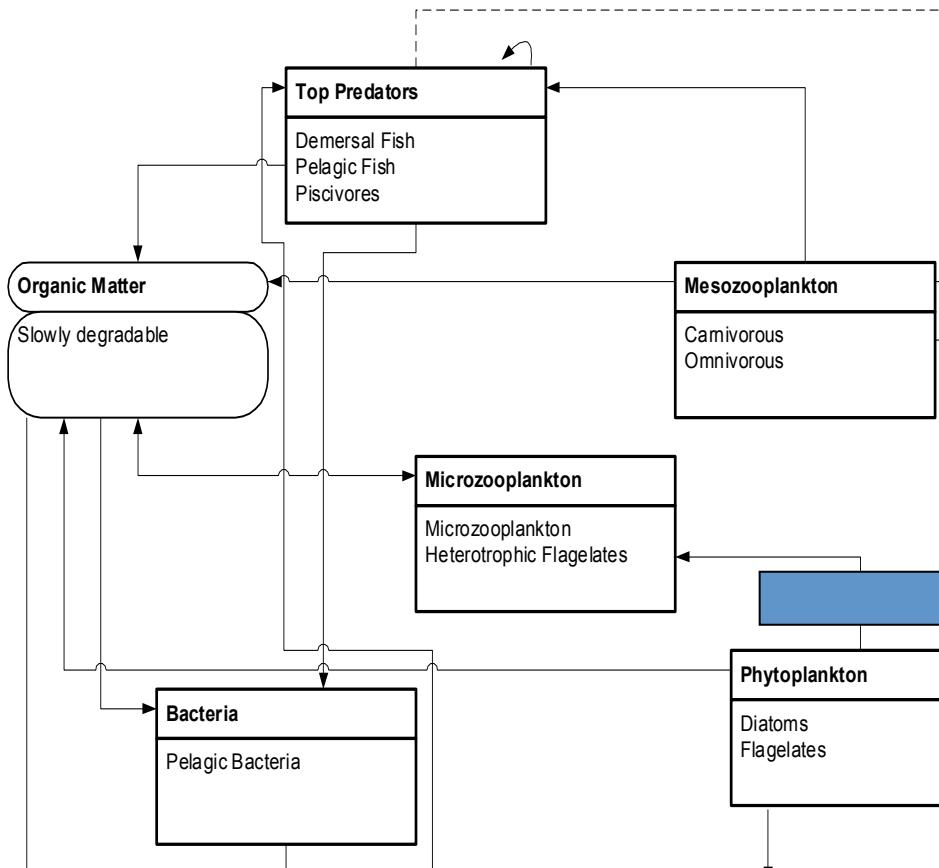
⋮

$$\frac{dx_n}{dt} = f_n(t, x_1, x_2, \dots, x_n)$$

# Process Flow Diagram



# Process Library



```

F:\pldata\BOD_DECAY.f90
=====
0001 ! * * * Top of File * * *
0002 SUBROUTINE ( PMSA , FL , IPPOINT , INCREM , NOSEG ,
0003 + NOFLUX , IEXPNT , IKNMRR , NOQ1 , NOQ2 ,
0004 + NOQ3 , NOQ4 )
0005 ****
0006 C
0007 C IMPLICIT NONE
0008 C
0009 C Type Name I/O Description
0010 C
0011 REAL(4) PMSA(*) ! I/O Process Manager System Array, window of routine to process library
0012 REAL(4) FL(*) ! O Array of fluxes made by this process in mass/volume/time
0013 INTEGER IPINT( 2 ) ! I Array of pointers in PMSA to get and store the data
0014 INTEGER INCREM( 2 ) ! I Incremental, IPINT for element loop, =constant, 1=spatially varying
0015 INTEGER NOSEG ! I Number of computational elements in the model
0016 INTEGER NOFLUX ! I Number of fluxes increment in the FL array
0017 INTEGER IEXPNT ! I From, To, From-1 and To+1 segment numbers of the exchange surfaces
0018 INTEGER IKNMRR ! I Active-Inactive, Surface-water-bottom, see manual for use
0019 INTEGER NOQ1 ! I Nr of exchanges in 1st direction, only horizontal dir if irregular mesh
0020 INTEGER NOQ2 ! I Nr of exchanges in 2nd direction, NOQ1+NOQ2 gives hor. dir. reg. grid
0021 INTEGER NOQ3 ! I Nr of exchanges in 3rd direction, vertical direction, pos. downward
0022 INTEGER NOQ4 ! I Nr of exchanges in the bottom (bottom layers, specialist use only)
0023 INTEGER INT( 2 ) ! Local work array for the pointerring
0024 INTEGER ISEG ! Local loop counter for computational element loop
0025 C ****
0026 C
0027 C Type Name I/O Description Unit
0028 C
0029 C
0030 REAL(4) BOD_D_Rate ! I Decay rate for Biochemical Oxygen Demand 1 / day
0031 REAL(4) BOD ! I Biochemical Oxygen Demand g O2 / m3
0032 REAL(4) BOD_D_flux ! F Decay flux of Biochemical Oxygen Demand g O2 / m3 / day
0033 INTEGER IBOD_D_flux ! Pointer to the Decay flux of Biochemical Oxygen Demand
0034 C ****
0035 C
0036 C IPNT = IPINT
0037 C IBOD_D_flux = 1
0038 C
0039 C DO 9000 ISEG = 1 , NOSEG
0040 C
0041 C BOD_D_Rate = PMSA( IPNT( 1 ) )
0042 C BOD = PMSA( IPNT( 2 ) )
0043 C **** Insert your code here ****
0044 C
0045 C BOD_D_flux = ??????
0046 C
0047 C **** End of your code ****
0048 C
0049 C
0050 C FL ( IBOD_D_flux ) = BOD_D_flux
0051 C
0052 C IBOD_D_flux = IBOD_D_flux + NOFLUX
0053 C IPNT = IPNT + INCREM
0054 C
0055 C 9000 CONTINUE
0056 C
0057 C RETURN
0058 C END
0059 C
0060 *** End of File ***
  
```

# The Role of Scientist

- Problem formulation
- Effective representation (representational engineering)
- Data manipulation
- Algorithm manipulation
- Post-processing

# Innovation Engine

- Logic is not the key to invention! What human's bring to invention process is not their logic, but their illogic!
- Logic should be used as a constraining factor (*i.e.* dimensional correctness, mass and momentum conservation)
- For the first time in history, we can tell to computer what to do, and not how to do it!

# Conclusions

- Knowledge Discovery is coming to a theater near you
- The best results are generated ONLY when a tight interaction between domain and method knowledge is secured
- This is MAINLY Hypothesis generation!