

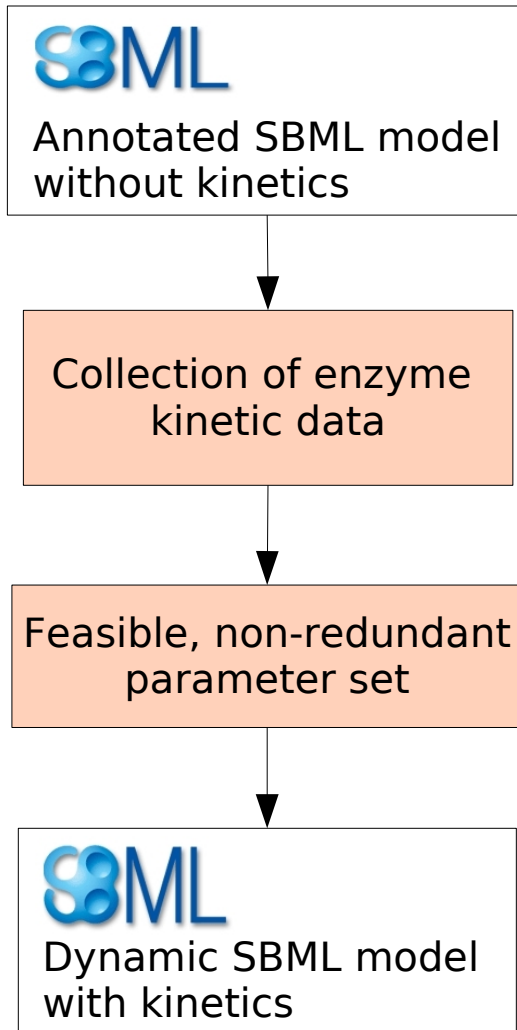
Integration of enzyme kinetic data from various sources

Storage and Annotation of Reaction Kinetics Data

EML, May 21-23, 2007

Wolfram Liebermeister, MPI-MG Berlin
Computational systems biology

Scheme for a modelling workflow

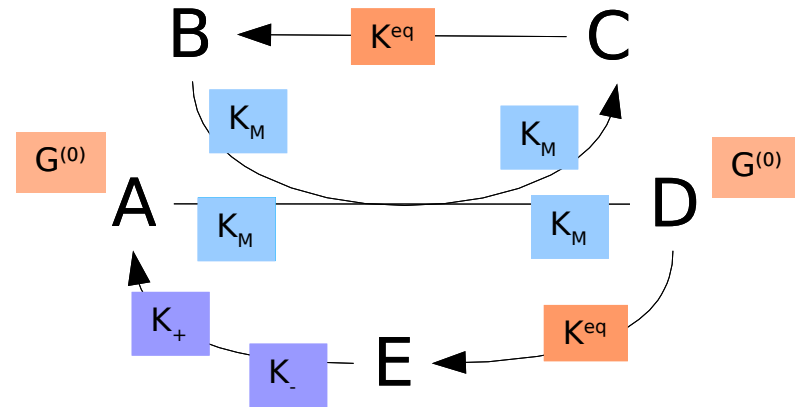


The screenshot shows a spreadsheet with the following data:

A	B	C	D	E	F	G	H
ParameterType	Unit	ReactionID	MetaboliteID	Value	StdDev	Organism	
1	G	nan	2-(alpha-Hydroxyethyl)thi	-122.21	435.68	nan	
2	G	nan	CO2	-472.49	2.0819	nan	
3	G	nan	Thiamin diphosphate	-256.38	435.84	nan	
4	G	nan	Pyruvate	-328.04	17.895	nan	
5	G	nan	ATP	-2236.1	65.149	nan	
6	G	nan	H2O	-212.47	35.217	nan	
7	G	nan	AMP	-565.82	65.417	nan	
8	G	nan	Phosphoenolpyruvate	-1139	27.775	nan	
9	G	nan	Orthophosphate	-1070.1	33.818	nan	
10	G	nan	ADP	-1405.1	59.227	nan	
11	G	nan	(S)-Malate	-712.34	27.956	nan	
12	G	nan	NAD+	1003.9	73.495	nan	
13	G	nan	NADH	1109.1	73.495	nan	
14	G	nan	NADP+	141.01	73.442	nan	
15	G	nan	NADPH	245.81	73.442	nan	
16	G	nan	Acetaldehyde	82.27	22.56	nan	
17	G	nan	CoA	-61.282	64.939	nan	
18	G	nan	Acetyl-CoA	-65.151	59.501	nan	
19	G	nan	H+	-37.07	15.263	nan	
20	G	nan	Acetyl phosphate	-1060.7	29.604	nan	
21	G	nan	Acetate	-245.08	23.373	nan	
22	G	nan	Pyrophosphate	-1915.3	50.826	nan	
23	G	nan	L-Glutamate	-349.09	76.675	nan	
24	G	nan	2-Oxoglutarate	-640.16	40.889	nan	
25	G	nan	NH3	48.945	76.675	nan	
26	G	nan	Oxalosuccinate	-1054.5	33.76	nan	
27	G	nan	UTP	-764.52	441.47	nan	
28	G	nan	D-Glucose 1-phosphate	-1290.9	41.99	nan	
29	G	nan	UDP-glucose	-134.41	437.32	nan	
30	G	nan	Oxaloacetate	-755.35	22.001	nan	
31	G	nan	Citrate	-1008.8	31.516	nan	
32	G	nan	Succinate	-483.99	45.467	nan	
33	G	nan	Succinyl-CoA	-301.13	64.248	nan	
34	G	nan	Acceptor	-287.58	542.35	nan	
35	G	nan	Fumarate	-494.6	34.448	nan	
36	G	nan	Reduced acceptor	-282.35	542.35	nan	
37	G	nan	Glyoxylate	-501.9	33.521	nan	
38	G	nan	Isocitrate	-999.09	31.516	nan	
39	G	nan	3-Carboxy-1-hydroxyprop	-413.7	436.27	nan	
40	G	nan	2-Phospho-D-glycerate	-1350.2	27.969	nan	
41	G	nan	(S)-Lactate	-282.34	24.726	nan	
42	G	nan	Ethanol	132.36	28.536	nan	
43	G	nan	alpha-D-Glucose	-433.04	50.639	nan	
44	G	nan					

How to combine enzyme kinetic data?

K_+	Turnover rate
K_M	K_M value
$G^{(0)}$	Gibbs free energy of formation
K^{eq}	Equilibrium constant



Kinetic parameters are dependent!

What is determined by the available data?

How can we account for the error widths?

Can we use prior knowledge to compensate for missing values?

Kinetic parameters are dependent



Reversible Michaelis-Menten kinetics:

$$v = E \frac{k_+^{\text{cat}} a / K_a - k_-^{\text{cat}} b / K_b}{1 + a / K_a + b / K_b}$$

Equilibrium constant K^{eq}

$$k^{\text{eq}} = \left(\frac{b}{a} \right)_{\text{eq}} = \frac{k_+^{\text{cat}} K_b}{k_-^{\text{cat}} K_a}$$

$$(1) \quad \ln k^{\text{eq}} = \ln k_+^{\text{cat}} - \ln k_-^{\text{cat}} + \ln K_b - \ln K_a$$

that means:

the parameters are dependent given the equilibrium constant

Gibbs free energy of formation $G^{(0)}$

$$k^{\text{eq}} = e^{-\Delta G^{(0)} / RT}$$

$$(2) \quad \ln k^{\text{eq}} = -\frac{1}{RT} N^T \mathbf{G}^{(0)}$$

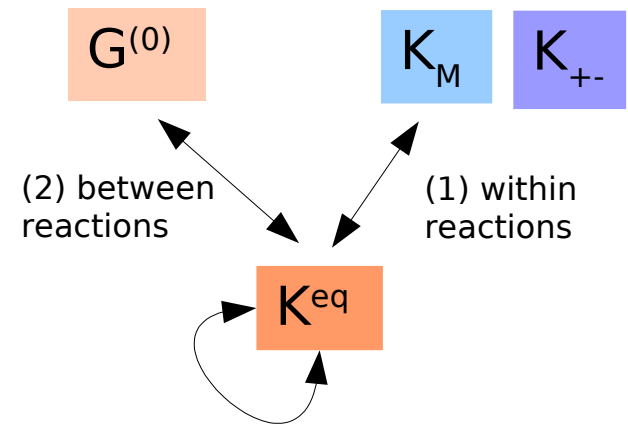
that means:

- the $G^{(0)}$ determine the K^{eq}
- the K^{eq} are dependent

Kinetic parameters are dependent

$$(1) \quad \ln k_l^{\text{eq}} = \ln k_{+l}^{\text{cat}} - \ln k_{-l}^{\text{cat}} + \sum_{\text{reactants } i} N_{il} \ln K_{il}^{\text{M}}$$

$$(2) \quad \ln k_l^{\text{eq}} = -\frac{1}{RT} \sum_i N_{il}^{\text{T}} G_i^{(0)}$$



Problem: complicated relationships between all parameters

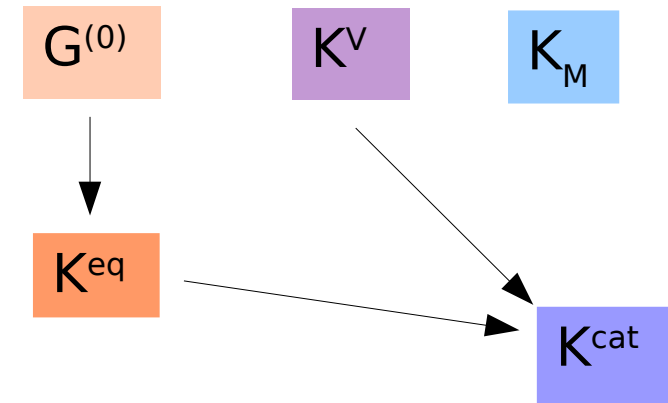
Solution: define a set of INDEPENDENT parameters

Independent system parameters

Velocity constant $k^V = \sqrt{k_+^{\text{cat}} k_-^{\text{cat}}}$

$$\ln k_l^{\text{eq}} = -\frac{1}{RT} \sum_i N_{il}^T G_i^{(0)}$$

$$\ln k_{\pm}^{\text{cat}} = \ln k^V \pm \frac{1}{2} \ln k^{\text{eq}}$$

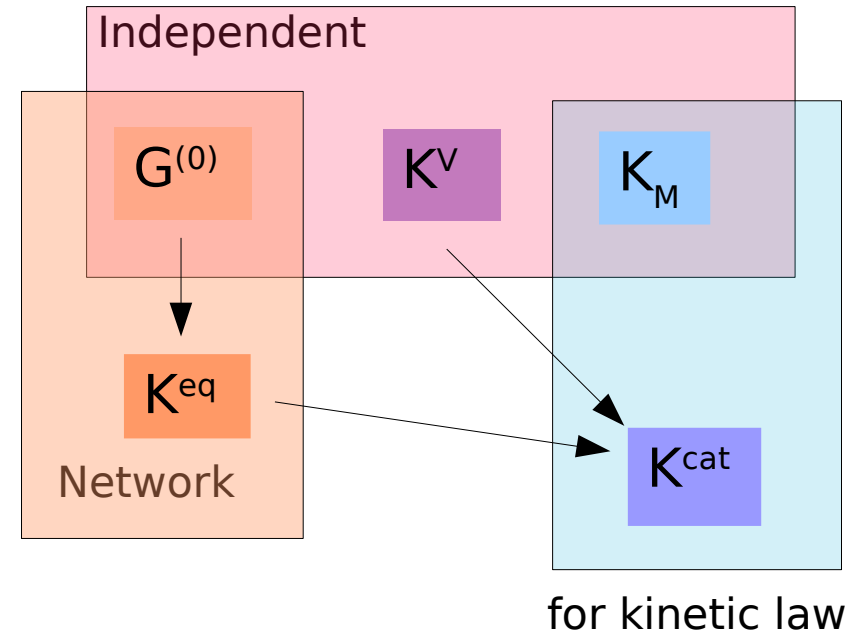


Independent system parameters

Velocity constant $k^V = \sqrt{k_+^{\text{cat}} k_-^{\text{cat}}}$

$$\ln k_l^{\text{eq}} = -\frac{1}{RT} \sum_i N_{il}^T G_i^{(0)}$$

$$\ln k_{\pm}^{\text{cat}} = \ln k^V \pm \frac{1}{2} \ln k^{\text{eq}}$$

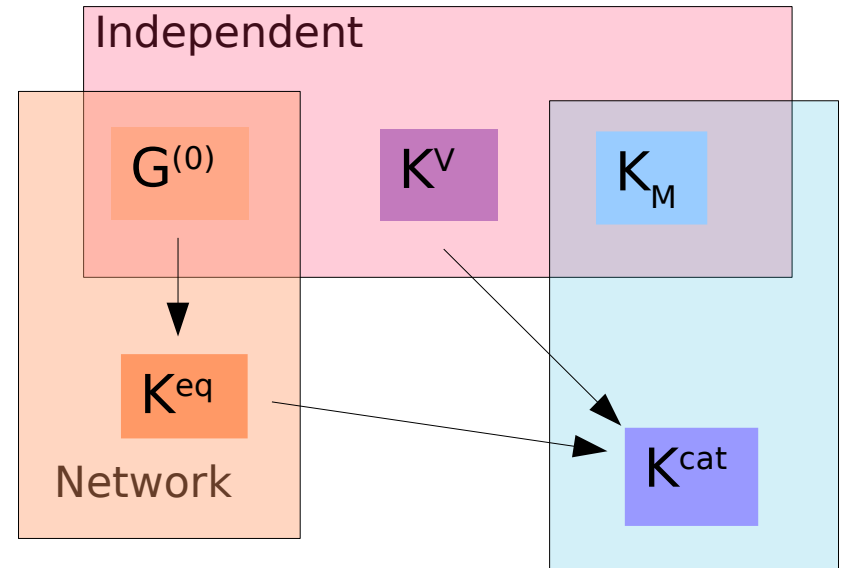


Independent system parameters

Velocity constant $k^V = \sqrt{k_+^{\text{cat}} k_-^{\text{cat}}}$

$$\ln k_l^{\text{eq}} = -\frac{1}{RT} \sum_i N_{il}^T G_i^{(0)}$$

$$\ln k_{\pm}^{\text{cat}} = \ln k^V \pm \frac{1}{2} \ln k^{\text{eq}}$$



Linear dependencies for logarithms:

$$x = R p$$

log. observable parameters

log independent parameters

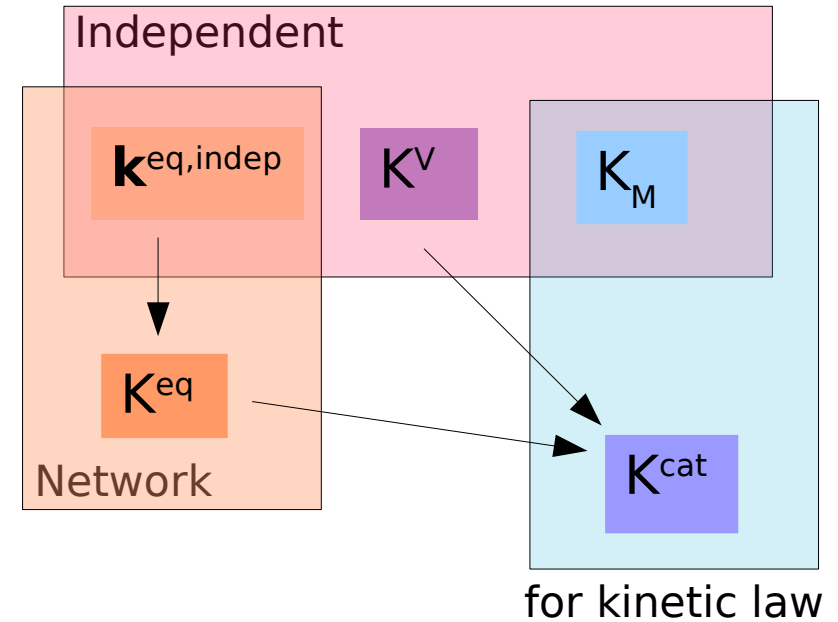
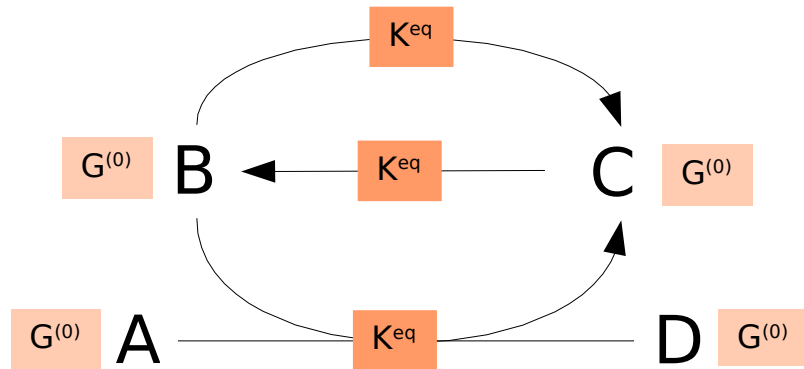
Dependence matrix

$$x = \begin{pmatrix} G^{(0)} \\ \ln k^V \\ \ln K^M \\ \ln k^{\text{eq}} \\ \ln k_+^{\text{cat}} \\ \ln k_-^{\text{cat}} \end{pmatrix}$$

$$p = \begin{pmatrix} G^{(0)} \\ \ln k^V \\ \ln K^M \end{pmatrix}$$

$$R = \begin{pmatrix} I & 0 & 0 \\ 0 & I & 0 \\ 0 & 0 & I \\ -\frac{1}{RT} N^T & 0 & 0 \\ -\frac{1}{2} \frac{1}{RT} N^T & I & 0 \\ +\frac{1}{2} \frac{1}{RT} N^T & I & 0 \end{pmatrix}$$

An alternative: independent equilibrium constants



Define independent equilibrium constants:

Set $N = N_0 L$ (N_0 has independent columns)

$$\begin{aligned} \ln k^{eq} &= -\frac{1}{RT} N^T G^{(0)} \\ &= -\frac{1}{RT} L^T N_0^T G^{(0)} = L^T \left(-\frac{1}{RT} N_0^T G^{(0)} \right) = L^T k^{eq, indep} \end{aligned}$$

Exploiting the independent parameters

Are the parameter data feasible?

given data x^* , check whether

$$\exists p: x^* = R p$$

Balance contradicting data

from redundant, contradictory x^* ,
obtain complete feasible x :

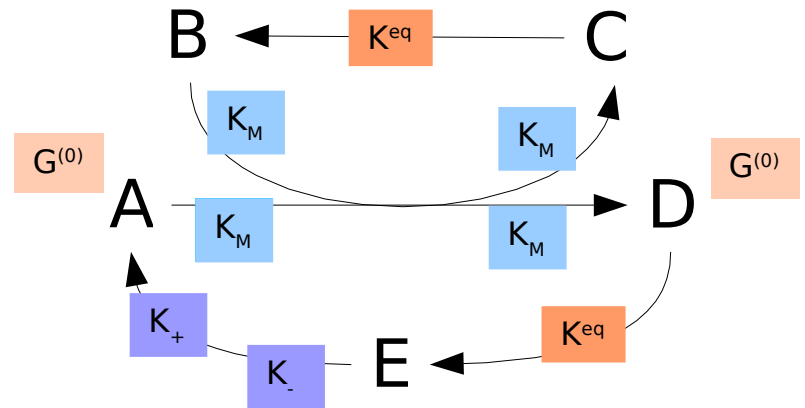
- solve $x^* = R^* p$ by least squares,
- set $x = R p$

But if there are too few data?

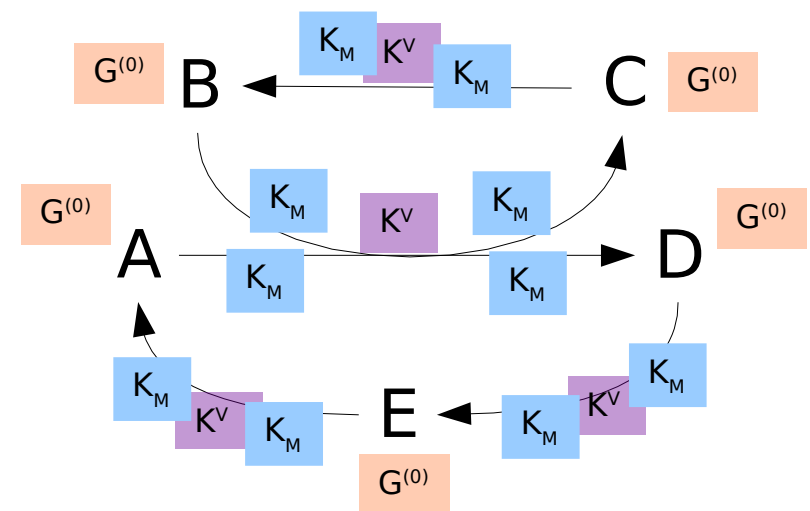
-> Values remain undetermined

Use plausible parameter ranges
as priors in Bayesian estimation

log data x^*

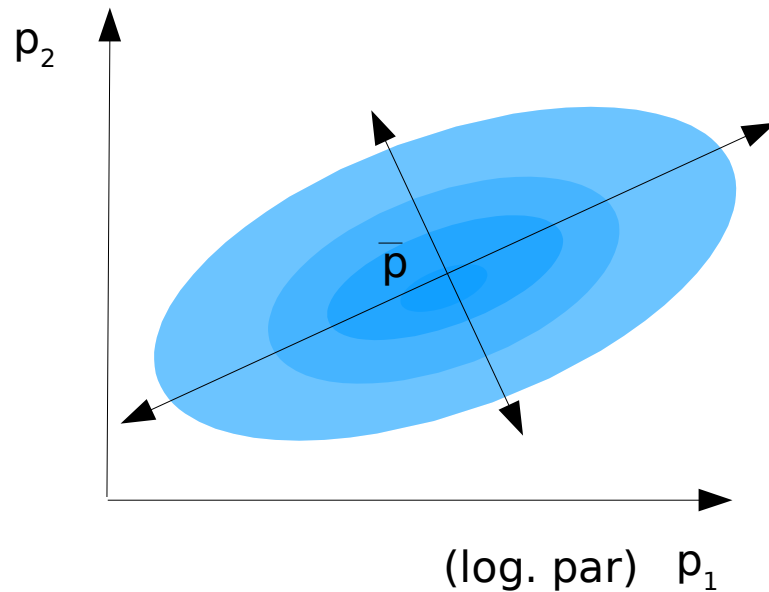


log. independent parameters p

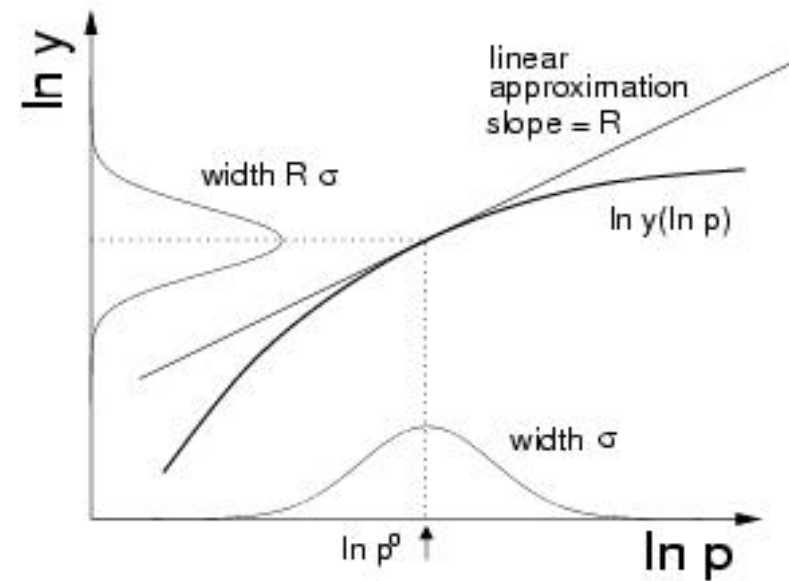


Uncertain parameters: distributions

Linear relationships between **log** parameters
-> **Gaussian** distributions of **log** parameters



Propagation of uncertainties

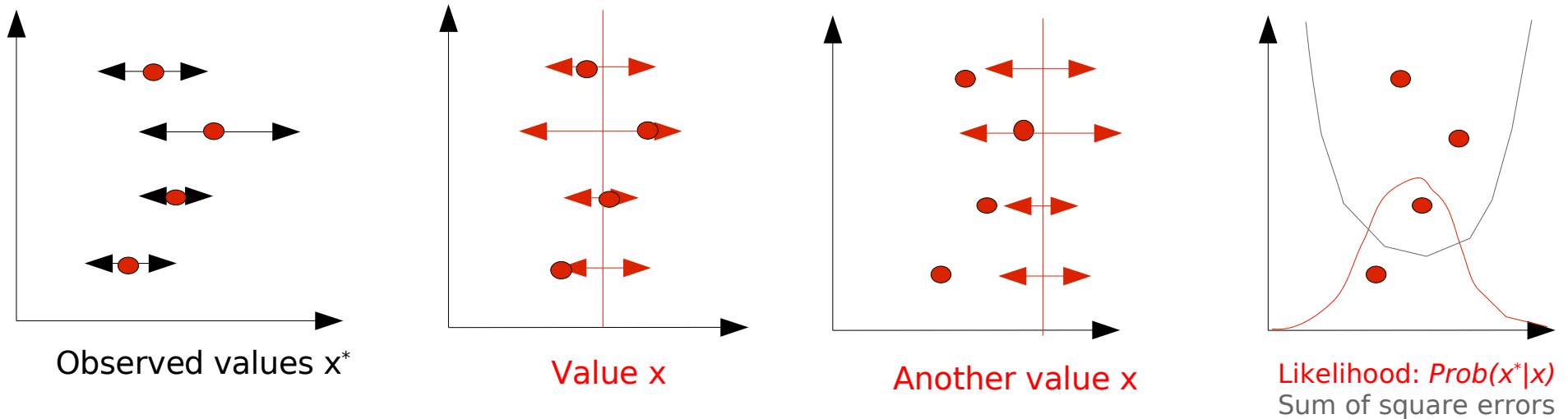


Linear functions $a = R b$

a multivariate Gaussian -> b multivariate Gaussian

$$\begin{aligned}\bar{a} &= R \bar{b} \\ \text{cov}(a) &= R \text{cov}(b) R^T\end{aligned}$$

What was maximum likelihood again?



Maximum likelihood estimator:

$$x_{\text{est}} = \operatorname{argmax}_x \text{Prob}(\text{data } x^* | \text{value } x)$$

Observations with Gaussian errors:

$$x_{\text{est}} = \operatorname{argmin}_x \sum_i \frac{(x_i^* - x)^2}{\sigma_i^2}$$

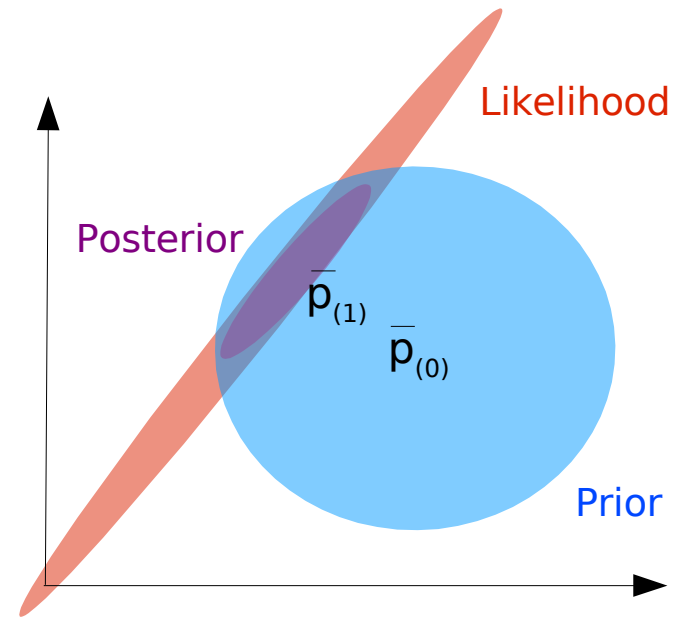
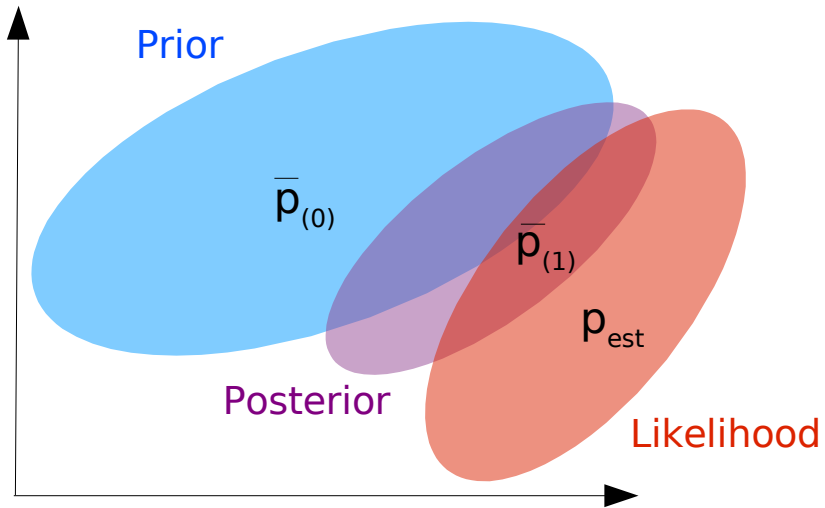
Estimate vector p from observations of $x = R p$:

$$x^* = \mathcal{N}(x = R p, C_x)$$

$$p_{\text{est}} = \operatorname{argmax}_p \text{Prob}(\text{data } x^* | \text{value } p)$$

$$\begin{aligned} p_{\text{est}} &= \operatorname{argmin}_p (x^* - R p)^T C_x^{-1} (x^* - R p) \\ &= (R^T C_x R)^{-1} R^T C_x^{-1} x^* \end{aligned}$$

Bayesian parameter estimation



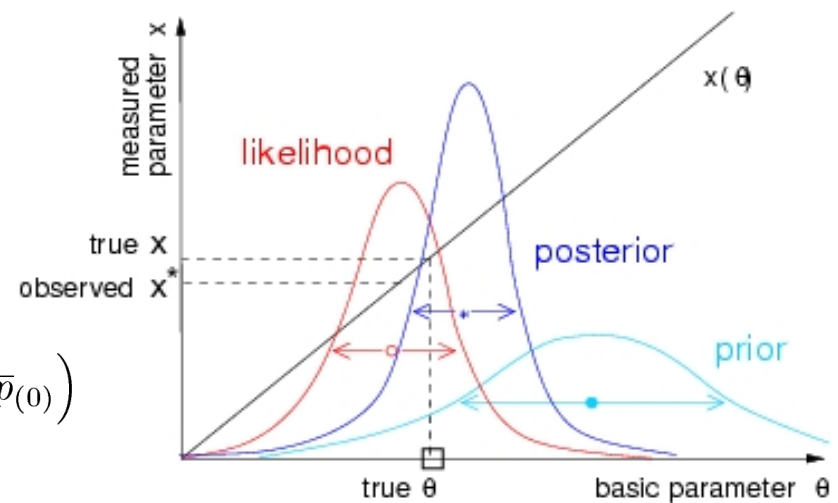
Additional prior knowledge in form of a distribution $Prob(p)$:

$$Prob(p | x^*) \sim Prob(x^* | p) Prob(p)$$

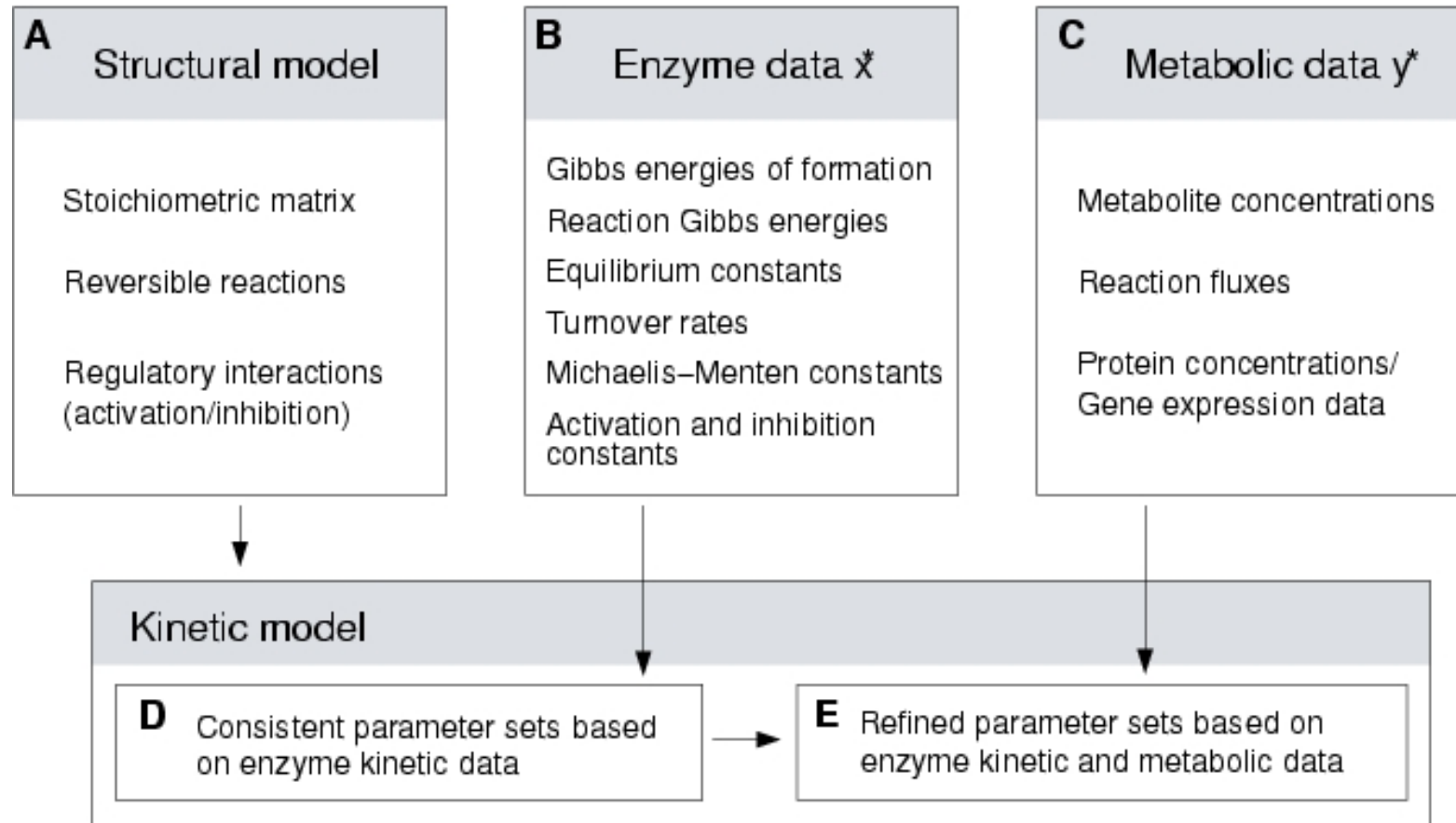
Given data $x^* = R p$:
mean and covariance matrix

$$\bar{p}_{(1)} = \left(C_{(0)}^{-1} + (R)^T C_x^{-1} R \right)^{-1} \left((R)^T C_x^{-1} x^* + C_{(0)}^{-1} \bar{p}_{(0)} \right)$$

$$C_{(1)} = \left(C_{(0)}^{-1} + (R)^T C_x^{-1} R \right)^{-1}$$



Improving the parameter estimates: integration of metabolic data



Important points in data collection

Each data point needs an error range

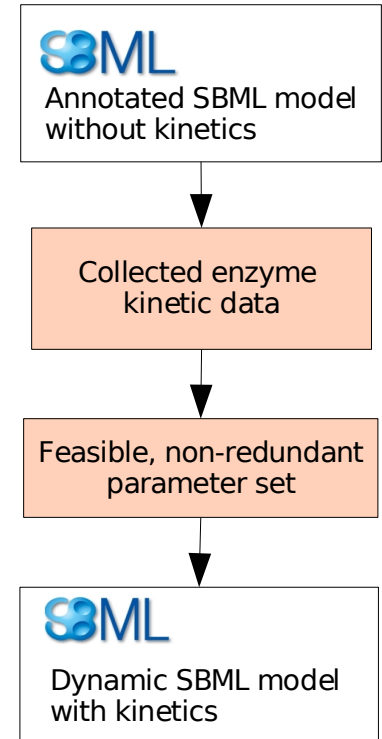
- Must be realistic
- How to compare error ranges for different kinds of data?
- Different error sources

Data must stem from independent sources

- never use the same database entry twice
- never reuse data that were already used in the computations

It is important to keep track of correlations

- “posterior is the new prior”
- some knowledge about equilibrium constants is stored in $G^{(0)}$ correlations



Aims and weaknesses of the workflow:

How can I use it?

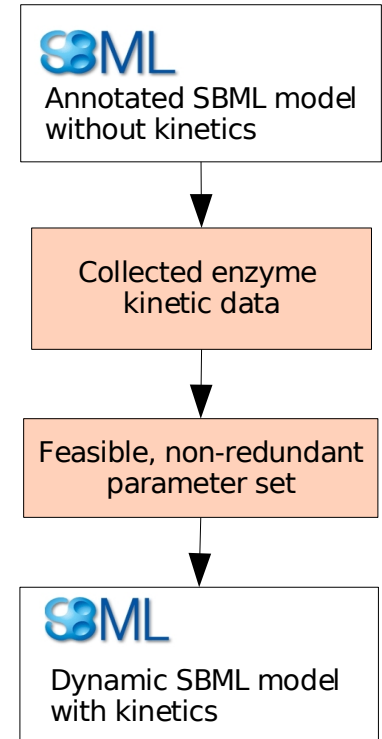
- model-driven data collection and mapping of relevant data
- assess the possible ranges and correlations of parameters;
- provide parameter priors for model fitting

What's the use of low-fi models??

- initial model as starting point for detailed models
- complete high quality models, fill gaps (better than nothing)
- determine which additional data are most needed

What are their weaknesses?

- convenience kinetics differs from true laws
- definition of parameters may differ
- computed, in-vitro, or inferred values differ from the ones needed in the model
- automatic data mapping requires detailed and reliable annotations
- general weaknesses of bottom-up-approaches



What do we need to make the workflow useful?

Data collection and mapping

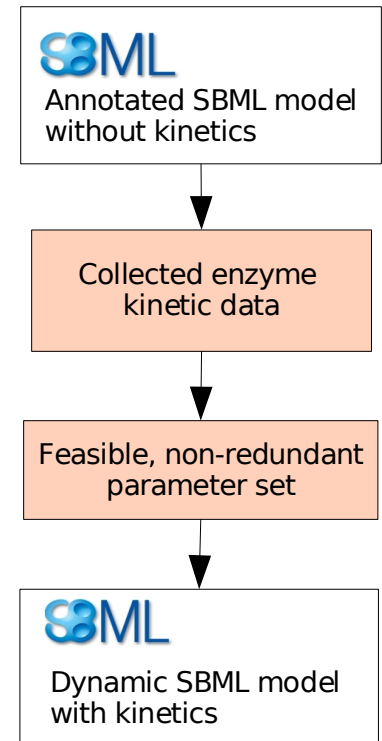
- easy access to databases
- precise annotations in models and kinetic data
- standard exchange formats for enzymatic data

Additional data sources

- databases
- statistical learning methods
- ab-initio calculations

Treatment of uncertainties & error propagation

- all data should come with error bars
- standards to describe uncertainties in SBML and enzyme data
- standards to describe relationships between parameters
- standards to quantify parameter correlations
- users: model fitting tools that use parameter priors



Thanks to ...



MPI for Molecular Genetics

Simon Borger
Jannis Uhlendorf
Anselm Helbig
Edda Klipp

Dietmar Schomburg lab

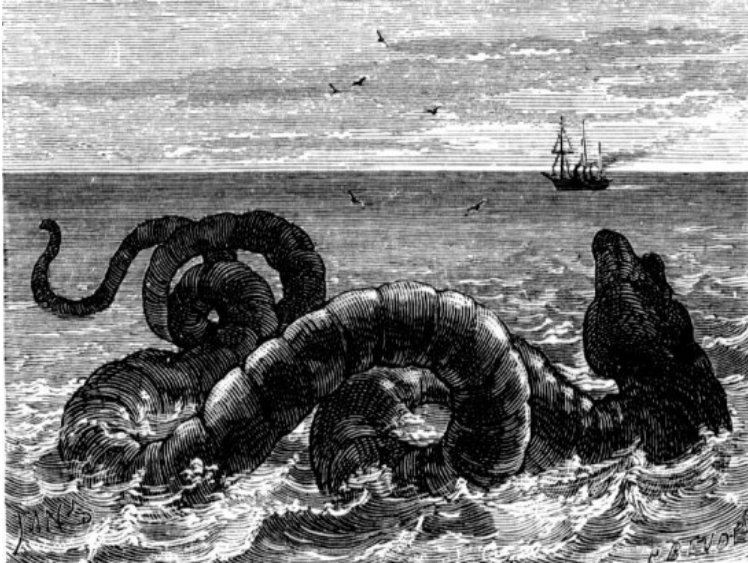
Sebastian Buchinger
Eva Vaylann

Genoscope, Evry

Vincent Schachter
Maxime Durot

... and to you !!!

Bottom-up and top-down model building



Scylla (“bottom-up”)
in-vitro data may be wrong ...
data may not be transferable ...
models SHOULD not fit all data ...



Charybdis (“top-down”)
the fitted model may work,
but the parameters have an
unclear interpretation ...