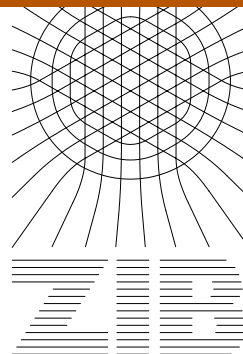

Konrad-Zuse-Zentrum
für Informationstechnik Berlin



Takustraße 7
D-14195 Berlin-Dahlem
Germany

PETER DEUFLHARD

**A Comparison of Related Concepts in
Computational Chemistry and Mathematics**

A Comparison of Related Concepts in Computational Chemistry and Mathematics¹

PETER DEUFLHARD

Abstract

This article studies the relation of the two scientific languages Chemistry and Mathematics via three selected comparisons: (a) QSSA versus dynamic ILDM in reaction kinetics, (b) lumping versus discrete Galerkin methods in polymer chemistry, and (c) geometrical conformations versus metastable conformations in drug design. The common clear message from these comparisons is that chemical intuition may pave the way for mathematical concepts just as chemical concepts may gain from mathematical precising. Along this line, significant improvements in chemical research and engineering have already been possible – and can be further expected in the future from the dialogue between the two scientific languages.

Key words: dimension reduction, model reduction, QSSA, ILDM, reaction kinetics, polyreaction kinetics, Hamilton systems, conformations, metastability, drug design

MSC (2000): 65–99, 65L05, 65L07, 65L99, 65P10

¹Invited talk, Leopoldina Symposium “Chemistry and Mathematics: Two Scientific Languages of the 21st Century”, MPI Göttingen

Introduction

When contemplating the role of Chemistry and Mathematics as two scientific languages of the 21st century, one may first associate the novel *The Glass Bead Game* (*Das Glasperlenspiel*) [9] from 1943. In this novel, HERMANN HESSE (1877–1962) had envisioned some unifying language of art and science, which many a reader might be tempted to identify as Mathematics. In the spirit of the Order of the Glass Bead Game players that language was only understood (and permitted!) to represent the close structural links between various up to all sciences and arts: linguistics versus history versus music versus physics versus mathematics – to be continued. Individuality of the scholars, even creativity, was definitely not meant to enter the Game. Through the personal development of the Magister Ludi Joseph Knecht, the main character of the novel, the poet’s genius left this narrow confinement opening up into the ‘real’ life. Is this a paradigm for the evolution of Mathematics after the invention of the computer, from the pure science of structures towards some key discipline, which intervenes in more and more parts of our real life? In fact, at the beginning of the 21st century, detailed mathematical models open new possibilities to master complexity and to explore smart technological options via modeling, simulation, and optimization. It is more than clear now that present Mathematics is not only a unifying language, but adds actual value in a joint interdisciplinary ‘game’ with nearly all fields of science and engineering – including, of course, Chemistry.

The present paper asks about the mutual role of the languages of Chemistry and Mathematics, restricted to computational chemistry and mathematics due to the author’s expertise. Rather than treating this kind of question in general, an answer is sought here via a synopsis of three comparable concepts from the two scientific disciplines. Of course, the selection is biased according to the author’s personal taste and experience. Section 1 deals with the development from the classical QSSA (quasistationary state approximation) to modern ILDM (intrinsic low-dimensional manifold) and its dynamics extension in reaction kinetics. Section 2 starts from lumping techniques in the numerical treatment of polyreaction kinetics and ends up with discrete Galerkin methods. Section 3 treats molecular conformations, from the essentially geometrical concept to the recent metastability concept.

1 Dynamic Dimension Reduction versus QSSA in Reaction Networks

Consider a singularly perturbed system of ordinary differential equations (ODEs) in the explicitly separated form

$$y' = f(y, z), \quad \epsilon z' = g(y, z), \quad (1.1)$$

for which the solution will be denoted by (y_ϵ, z_ϵ) . This might model the kinetics of a complex reaction network – see, e.g., the survey paper [7] of DEUFLHARD and NOWAK (1986), which is especially addressed to computational chemists. Typically, such an ODE system will be 'stiff' and nonlinear; for a deeper understanding of *stiffness*, the reader may refer to the recent textbook [3] by DEUFLHARD and BORNEMANN (2002). For our present purposes, we just recall that stiff systems asymptotically reach an equilibrium. In terms of the numerical treatment, this implies that information about the Jacobian matrix of the right hand side must enter – thus increasing the computational cost considerably compared to the solution of nonstiff systems. Typically, different components reach the equilibrium at rather different time scales, which gives rise to a separation of 'slow' modes y and 'fast' modes z . In complex reaction network models, the identification of 'fast' and 'slow' modes is often quite difficult. It is not for nothing that experts speak of the 'golden' ϵ that has to be found in each problem.

Given such an identification, the assumption of quasi-stationarity

$$\epsilon z' = 0 ,$$

then leads to the *quasi-stationary state approximation* (QSSA) (y_0, z_0) defined by the differential-algebraic equations (DAEs)

$$y' = f(y, z), \quad 0 = g(y, z) . \quad (1.2)$$

Mathematically speaking, the transition from the ODE model (1.1) to the DAE model (1.2) is justified under the assumption that the equilibrium point is unique and *attractive*. In this situation, an explicit local parametrization of the form

$$z = h(y)$$

will certainly exist. If, in addition, *consistent* initial values $y(0), z(0)$ are given, i.e. initial values satisfying

$$g(y(0), z(0)) = 0 ,$$

then system (1.2) may be replaced by the *reduced* ODE system

$$y' = f(y, h(y)) . \quad (1.3)$$

The transition from (1.1) to (1.3) is usually called *dimension reduction* or also *model reduction*.

Classical QSSA approach. In this approach, the selection of fast modes is based on chemical insight into the reaction network. This means that a chemist may have identified certain 'radicals' z exhibiting only a short-lived appearance in the chain of reactions. In the days before the advent of efficient numerical stiff integrators, the radical components would then be eliminated and an analytic expression $z = h(y)$ would be derived, which, when inserted into the reduced

model (1.3), could be cheaply integrated by some nonstiff numerical integrator. However, as it turned out, the thus obtained DAE systems are often not uniquely solvable – see, again, the book [3] or the survey in [7], where an example of this kind has been worked out.

Of course, modern computational chemists would just apply their favorite numerical stiff integrator DASSL, RADAU53, or LIMEX to the original unseparated system – see, e.g., the textbook [3] for details and references.

ILDm approach. This alternative approach (ILDm: *intrinsic low-dimensional manifold*) has been proposed by MAAS and POPE in 1992, see [13]. It avoids the occurrence of nonuniqueness in the derived DAE models, but requires a deeper mathematical understanding. In order to convey the main idea, we consider the general ODE model

$$x' = F(x), \quad x(0) = x_0 \tag{1.4}$$

and its local Jacobian

$$J = F_x(x(0)).$$

The corresponding linearized differential equation

$$\delta x' = J\delta x$$

has solution components with a growth behavior

$$\delta x(t) \sim \exp(\Re\lambda_i t) = \exp(-t/\tau_i).$$

specified by those eigenvalues λ_i of J , for which $\Re\lambda_i < 0$, or by the corresponding time scales

$$\tau_i = -\frac{1}{\Re\lambda_i} > 0. \tag{1.5}$$

As can be seen, the computation of the various time scales τ_i would require the solution of the corresponding eigenvalue problem – which, however, may be ill-conditioned for *nonsymmetric* Jacobian matrices, the usual case.

Fortunately, in the presence of a sufficiently large spectral gap between the corresponding eigenvalues and the rest of the spectrum, the computation of invariant eigenspaces is known to be a well-conditioned problem. For its solution we use two steps to transform the Jacobian matrix J . The first step consists of an orthogonal similarity transformation such that

$$Q^T J Q = \bar{S} = \begin{pmatrix} S_{11} & S_{12} \\ 0 & S_{22} \end{pmatrix}.$$

Here Q is an orthogonal matrix and \bar{S} (essentially) an upper triangular matrix. It is possible to arrange the diagonal elements of \bar{S} according to the order of magnitude of the real parts of the eigenvalues. If at least one of these real parts is negative, then, with a reduced dimension, say r , the above block decomposition

of \bar{S} can be defined in terms of an associated parameter $\mu_r < 0$ by the following condition:

$$\mu_r = \max_{\lambda \in S_{22}} \Re \lambda = \Re \lambda_{r+1} < 0 \quad \text{and} \quad \min_{\lambda \in S_{11}} \Re \lambda = \Re \lambda_r > \mu_r.$$

In the second step the coupling matrix S_{12} is eliminated (via a so-called Sylvester equation) so that a nonorthogonal similarity transformation

$$T_r^{-1} J T_r = S = \begin{pmatrix} S_{11} & 0 \\ 0 & S_{22} \end{pmatrix}$$

is realized. The transformation matrix and its inverse have the form

$$T_r = Q \left(I + \begin{pmatrix} 0 & C_r \\ 0 & 0 \end{pmatrix} \right) \quad \text{and} \quad T_r^{-1} = \left(I - \begin{pmatrix} 0 & C_r \\ 0 & 0 \end{pmatrix} \right) Q^T.$$

At this point, we are now able to identify 'fast' components z and 'slow' components y by virtue of

$$T_r^{-1} x = \begin{pmatrix} y \\ z \end{pmatrix}, \quad T_r^{-1} F = \begin{pmatrix} f \\ g \end{pmatrix}.$$

With this transformation we have obviously established the connection between our original full network model (1.4) and the singularly perturbed problem (1.1) – and found the 'golden'

$$\epsilon = \frac{1}{|\mu_r|} = \frac{1}{|\Re \lambda_{r+1}|}.$$

Obviously, the choice of the reduced dimension r is coupled with the selection of the perturbation parameter ϵ . A comparison with (1.5) shows that $\epsilon = \tau_{r+1}$ represents that time scale, below which a resolution of the system dynamics is ignored in the modeling.

In the language of chemistry, we have identified 'radicals' – but *without* any use of chemical insight, merely with the help of a numerical algorithm. The inevitable downer is that the thus obtained components z can no longer be associated with selected chemical species, they are just 'abstract radicals'. In summary, we see that *the intuitive, but mathematically deficient QSSA concept from chemistry has turned into the less intuitive ILDM concept based on more precise mathematical terms.*

Dynamic dimension reduction. On the above mathematical basis, we can now even go a step further and allow for a *time dependent* reduction of dimension. The reduced model will only be useful, when the differences $y_\epsilon - y_0$ and $z_\epsilon - z_0$ are 'sufficiently small'. Actually we are free to dispose about the initial values such that

$$z_\epsilon(0) = z_0(0) + \zeta_0(0), \quad \text{and} \quad y_\epsilon(0) = y_0(0).$$

Then, on the theoretical basis given in [5] by DEUFLHARD and HEROTH (1996), we may even derive a componentwise error criterion: Given a user prescribed error tolerance TOL, we just need to require that

$$\epsilon |f(y_\epsilon(0), z_\epsilon(0)) - f(y_\epsilon(0), z_0(0))| \leq \text{TOL}. \quad (1.6)$$

This criterion actually permits a time varying definition of 'abstract radicals', i.e. we obtain a reduced dimension $r(t)$ – which explains the name *dynamic* dimension reduction. Note that the actual evaluation of the criterion requires consistent initial values $(y_0(0), z_0(0))$ as input arguments – for more details we refer again to the textbook [3].

In passing we want to mention that an actual *elimination* of all the fast modes z succeeds only, if the algebraic equations $g = 0$ need not be solved explicitly, but through pointwise evaluation of the fast variables z via simple *table lookups*.

Illustrative example: Oxyhydrogen combustion. We base our discussion on the chemical model given in [13], which involves 37 elementary chemical reactions for 8 chemical species. This leads to a system of 8 differential equations. A reduced version of this initial value problem including only 7 chemical reactions for 7 species has been given in [10] by HOPPENSTEADT ET AL. (1981) as an example of the failure of the classical QSSA approach: there an elaborate analytical treatment showed that *before*, *during*, and *after* the combustion rather different 'golden' ϵ 's had to be defined.

The computational results presented in Fig. 1 were obtained by means of the differential–algebraic numerical integrator LIMEX with adaptive control of order and stepsize; the prescribed error tolerance was $\text{TOL} = 10^{-2}$ in (1.6). Prior to the treatment as a singular perturbation problem, two dynamical invariants of the system were eliminated (since these induce associated zero eigenvalues). Consequently, the dimension $n = 8$ was reduced in advance to an effective system dimension $n_{\text{eff}} = 6$. The upper part of Fig. 1 shows the computed solution for the species H_2 , O_2 and H_2O as a function of time t . Within the required error tolerance, these numerical results agree with those for the full model. The lower part of Fig. 1 shows that the dimension reduces to $r = 2$ *before*, to $r = 1$ *during*, and even to $r = 0$ (equilibrium) *after* the combustion.

This approach only pays off in the context of partial differential equations, especially of the reactive flow type – see the (German) habilitation thesis by MAAS (1993) as quoted in the English survey paper [12].

2 Discrete Galerkin Methods versus Lumping in Polyreaction Kinetics

Polymers are known to be chains of typically ten thousand up to ten million monomers, which are simple molecules or molecular groups. Corresponding

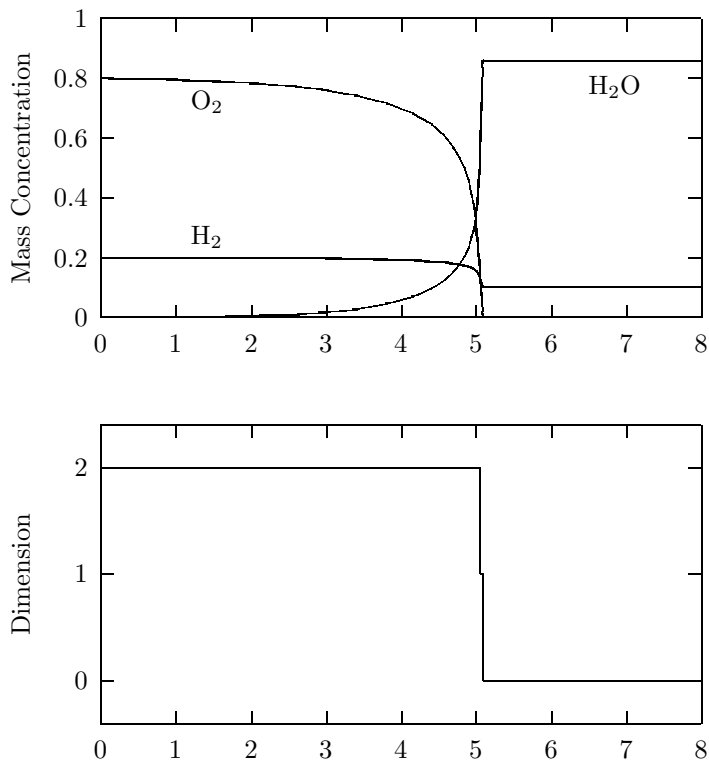


Figure 1: Oxyhydrogen combustion. *Top*: Chemical species behavior as a function of time t . *Bottom*: Reduced dimension $r(t)$.

mathematical models of polyreaction kinetics involve the same number of ODEs, usually nonlinear and stiff. For quite a while, this problem class had been among the real challenges in both computational mathematics and computational chemistry. In what follows we denote by $P_s(t)$ the concentration of some polymer of *chain length* s at time t . For ease of writing, we will not distinguish between the chemical species P_s , its concentration $P_s(t)$, and the *chain length distribution* $\{P_s(t)\}_{s=1,2,\dots}$, but just rely on the context. In order to convey an impression of the CODE problem class, we start with a real life example.

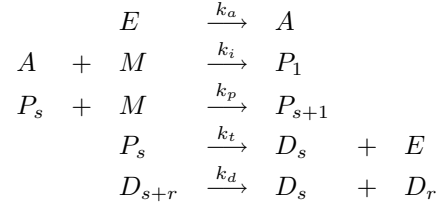
Example: Biopolymerization. This problem deals with an attempt to recycle waste of synthetic materials in an ecologically satisfactory way – which is certainly an important problem of modern industrial societies. An attractive idea in this context is to look out for synthetic materials that are both produced and eaten by bacteria – under different environmental conditions, of course. A schematic illustration of the production process within such a bacterial recycling

is given in Fig. 2: there certain bacteria use a monomer, fructose, as chemical input to produce a polymer, polyester (PHB), as chemical output.



Figure 2: Biopolymerization: bacteria eat sugar and produce polyester. White areas: polyester granules within bacteria cells.

The macromolecular reaction steps of production and degradation of the polymer can be summarized in the following *chemical model*



with $s, r = 1, 2, \dots, s_{max}$. Herein M denotes the monomer, E an enzyme, A the activated enzyme, P_s the 'living' and D_s the 'dead' polymer.

The corresponding *mathematical model* for the above process reads

$$\begin{aligned}
 E' &= -k_a E + k_t \sum_{r=1}^{s_{max}} P_r \\
 A' &= +k_a E - k_i A M \\
 M' &= -k_p M \sum_{r=1}^{s_{max}} P_r - k_i A M \\
 P_1' &= -k_p M P_1 + k_i A M - k_t P_1 \\
 P_s' &= -k_p M (P_s - P_{s-1}) - k_d P_s, \quad s = 2, 3, \dots, s_{max} \\
 D_s' &= +k_t P_s - k_d (s-1) D_s + 2k_d \sum_{r=s+1}^{s_{max}} D_r, \quad s = 1, 2, \dots, s_{max}.
 \end{aligned}$$

The truncation index s_{max} is not known a priori; practical considerations suggest $s_{max} = 50.000$ or so – which means that the above system consists of 100.000 ODEs, each of which has about the same number of terms (!) in the right hand side.

Lumping techniques. In the chemical literature, certain linear combinations of polymer components are defined to yield a much smaller number of ODEs involving the corresponding linear combinations of the right hand sides. The thus constructed smaller systems are then numerically tractable – see KUO and WEI (1969) in [11]. Clearly, a proper collection of components will require detailed a-priori insight into the process under consideration. In some cases, a logarithmic mesh equilibration is imposed on the basis of chemical intuition, as privately communicated to the author by EDERER (1984): the larger s , the more components are ‘lumped’ together. However, even though this technique is reported to work satisfactorily in some *linear* ODE models, it is certainly unreliable for *nonlinear* models – which actually represent the bulk of industrially relevant models.

Discrete Galerkin methods. This approach dates back to suggestions of DEUFLHARD and WULKOW(1989) in [8]. Facing such a huge number of ODEs, they suggest to expand the problem to infinitely many ODEs, more precisely: to Countably many ODEs – in short named CODEs. In fact, such an approach appears to be natural for a mathematician! The main idea of the approach is to construct a special discrete Hilbert space, a so-called sequence space, and an associated Galerkin method. The key to these so-called discrete Galerkin methods is the introduction of a *discrete inner product*

$$(f, g) := \sum_{s=1}^{\infty} f(s)g(s)\Psi(s) \quad (2.1)$$

in terms of some prescribed positive decaying *weighting function* Ψ , which takes care of the regularity of the infinite sum. This product induces a set of *orthogonal polynomials*, say $\{l_j\}, j = 1, 2, \dots$, for the discrete variable s such that

$$(l_j, l_k) = \gamma_j \delta_{jk} \quad , \quad \gamma_j > 0 \quad j, k = 0, 1, 2, \dots \quad (2.2)$$

A first attempt of this kind has been based on the discrete weight function

$$\Psi(s) = \rho^s, \quad \rho < 1 \quad (2.3)$$

The kernel $\rho = \exp(-\beta)$ can be interpreted as a uniform discretization of the decaying exponential $\exp(-\beta t)$ in terms of some $\beta > 0$. The free choice of β corresponds to the free choice of ρ . Since the exponential generates the classical LAGUERRE polynomials $L_k(t)$ for the continuous variable t , the $l_k(s)$ are here called *discrete* LAGUERRE polynomials. For the actual adaptation of ρ see the “moving weight function” concept as worked out in [8]. For the solution P of the CODE, discrete Galerkin methods try the corresponding separation ansatz

$$P(s, t) = \Psi(s) \sum_{k=0}^{\infty} a_k(t) l_k(s). \quad (2.4)$$

In this setting, a *linear* CODE would be written as

$$\frac{\partial P}{\partial t} = \mathcal{A}P,$$

where \mathcal{A} denotes some discrete operator, bounded or unbounded in dependence on the modeled polyreaction mechanisms. Upon insertion of the expansion (2.4), multiplication by the test function $l_j(s)$, summation over s , change of the summation order, and use of the above orthogonality relations, we again end up with a CODE, this time for the Galerkin coefficients

$$\gamma_j a_j'(t) = \sum_{k=0}^{\infty} a_k(t) (l_j, \mathcal{A}l_k) \quad j = 0, 1, \dots \quad (2.5)$$

The attractive qualitatively new feature of this CODE compared to the original one is that the Galerkin coefficients are known to decay, if only the desired solution actually lies in the prescribed discrete Hilbert space – an assumption to be carefully taken into account using reasonable approximation arguments (estimation of truncation errors). Under this assumption we are then able to truncate the above k -sum, avoiding any artificial 'closure assumptions'.

There still remain infinite sums to be evaluated in the right hand side, which have the structure

$$(l_j, \mathcal{A}l_k) = \sum_{s=1}^{\infty} l_j(s) \mathcal{A}(s) l_k(s) \Psi(s).$$

These sums also must be approximated somehow by a finite number of terms (to prescribed accuracy TOL). In most of the relevant cases, analytical expressions are not available so that numerical approximations are the only choice. Today, the most efficient approximation is done via a discrete GAUSS-CHRISTOFFEL quadrature, which is directly based on the selected weight function Ψ . For $\Psi = \rho^s$, let τ_{in} , $i = 0, \dots, n$ denote the zeroes of the discrete Laguerre polynomials l_{n+1} and λ_{in} the associated weights of the discrete Gauss-Laguerre quadrature rule. Then we obtain the approximation

$$(l_j, \mathcal{A}l_k) \approx \sum_{i=0}^n \lambda_{in} l_j(\tau_{in}) \mathcal{A}(\tau_{in}) l_k(\tau_{in}).$$

The sum might even be exact, if a proper choice of finite n can be made that corresponds with the polynomial order of the discrete integrand. Obviously, this is nothing else than a special kind of 'lumping', this time based on subtle mathematical tools. In summary, we see that *the intuitive, but deficient chemical concept of lumping has turned into the less intuitive, but algorithmically more efficient mathematical concept of discrete Galerkin methods.*

This new algorithmic view of polyreaction models opened the door to a much improved version, a discrete $h - p$ finite element method, which is now the

basis for an efficient treatment of challenging polymer problems in industry. For further references see, e.g., the survey [16] by WULKOW (1996), especially addressed to computational chemists, or chapter 3 in the survey article [1] of DEUFLHARD (2000), especially addressed to mathematicians.

3 Metastable versus Geometrical Conformations in Drug Design

In computational biotechnology, algorithms from discrete mathematics or computer science already have played a publicly visible role for some time already – for example, in the decoding of the human and other genomes. These approaches primarily aim at a clarification of the *secondary structure* of biomolecular systems. However, in most cases, the key to an understanding of *biomolecular function* is the *tertiary structure*, i.e. the *geometrical shape* in 3D. On top of it, molecular *dynamics* rules biological function, which makes computational drug design an extremely challenging task.

The whole situation is basically characterized by the spreading between real times of pharmaceutical interest (in the region of *msec* up to *min*) and simulation times (presently in the region of *psec* up to *nsec*). Detailed examination of the problem reveals that the computation of molecular dynamics has a hidden mathematical limitation: the arising trajectories are Hamiltonian and as such chaotic. This implies that long term trajectory simulations – as typically performed in classical molecular dynamics (MD) – can, at best, only yield information about time averages. On the basis of this insight, an investigation of the dynamics of molecular systems over the time scales of interest will require a rather different mathematical approach. Such an approach, now called *conformation dynamics*, has been derived and worked out in a series of papers [15, 14, 6] by DEUFLHARD and SCHÜTTE and their coworkers since 1997 and is still under investigation.

Geometrical conformations. The term 'conformation' used in conformation dynamics is quite different from the earlier classical term 'conformation' that has been used within chemistry for quite a while. In fact, the chemical term condenses the scientific experience that certain geometrical forms of molecules – often with given extra names like *cis*, *trans*, or *gauche* – play a dominant role in certain reaction mechanisms. This term does not characterize a single molecular 'configuration', which is just a *point* in 3D space, but a whole *set* of 'similar' configurations – whatever similar means. Therefore, the name *geometrical conformation* has been coined for this meaning.

Metastable conformations. In order to elucidate the meaning of the term 'conformation' used in conformation dynamics, we have to outreach somehow.

Hamiltonian differential equations. Let N atoms of a molecular system be spec-

ified in terms of their spatial coordinates $q = (q_1, \dots, q_N)$ and their generalized momenta $p = (p_1, \dots, p_N)$. Then, usually, the Hamilton function H has the separated form

$$H(q, p) = \frac{1}{2} p^T M^{-1} p + V(q) ,$$

where the first term is the kinetic energy $T(p)$, the second term the potential energy. From given H , the Hamiltonian differential equations are defined as

$$q'_i = \frac{\partial H}{\partial p_i}, \quad p'_i = -\frac{\partial H}{\partial q_i}, \quad i = 1, \dots, N.$$

Given initial values $x_0 = (q_0, p_0)$, we may assume that the above initial value problem has a *unique* solution, formally written in terms of the flow Φ as

$$x(t) = (q(t), p(t)) = \Phi^t x_0 .$$

In addition, we have to study the *condition number* κ , which characterizes the sensitivity of the unique solution under perturbation of the initial values. As shown in Section 3.1.2 of the textbook [3] by DEUFLHARD and BORNEMANN (2002), such a quantity can be defined (in first order perturbation theory) as

$$\|\delta x(t)\| \leq \kappa(t) \|\delta x_0\| , \quad \kappa(t) = \|\partial \Phi^t / \partial x_0\| .$$

As already discovered by POINCARÉ (1881–1885), Hamiltonian systems are *chaotic*, which implies that $\kappa(\infty) = \infty$. In the context of numerical analysis, however, the interesting question is, after which characteristic critical time the condition number exceeds the inverse initial accuracy or, colloquially speaking, after which time the molecule 'forgets its history'. For so-called integrable Hamiltonian systems (such as the popular Kepler problem) the condition number is known to grow linearly. In real life molecular dynamics problems, however, the growth is exponential

$$\kappa(t) \sim \exp(t/t_{\text{crit}}) . \tag{3.1}$$

The critical times t_{crit} turned out to be typically no longer than a few ps – which had been a surprising phenomenon even to the experts!

Example: Trinucleotide ACC. We illustrate the effect for a short RNA segment consisting of 94 atoms and containing three genetic letters. Fig. 3 shows simulated configurations at times $t = 0.0$ ps, and $t = 20$ ps. At the start, the two molecular configurations are nearly identical. After only 20 ps they differ completely.

The resulting configurations – a spherical shape on the left, a stretched one on the right – remain essentially the same over quite long time spans. In other words: these forms are metastable, which motivates the name *metastable conformations*.

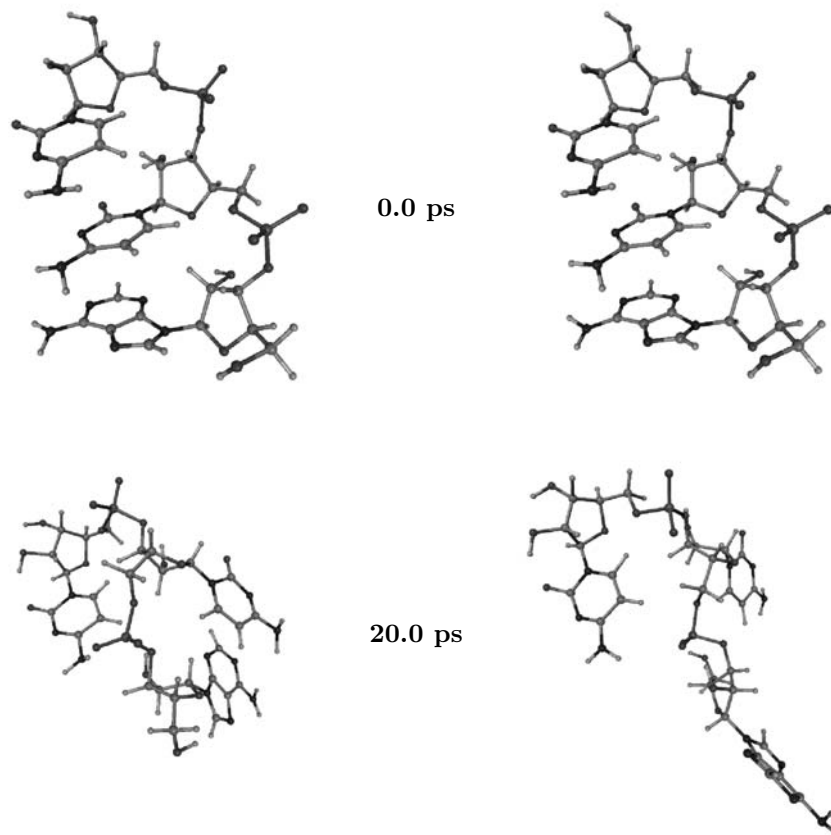


Figure 3: Trinucleotide ACC: Development of distinct conformations from nearly identical initial configurations

This kind of observation stimulated severe changes in the mathematical modelling of molecular dynamics. Instead of the *point concept* of classical mechanics a *set concept* turned out to be actually needed.

Stochastic transition operator. A first attempt to find such a different mathematical model has been suggested in [4] by DEUFLHARD ET AL. (1999). This approach has been based on the so-called PERRON–FROBENIUS operator, a special stochastic operator defined over the space of both position variables q and generalized momenta p . In [15, 14], SCHÜTTE (1999) constructed a more appropriate stochastic operator T , defined only over the space of position variables q .

Starting point of his construction was the fact that in a chemical lab, with constant temperature and constant volume, the deterministic model should be embedded into a canonical distribution f_0 . With β essentially the inverse tem-

perature we may factorize this distribution according to

$$f_0 = \mathcal{P}\mathcal{Q}, \quad Z = Z_p Z_q, \quad \int \mathcal{P}(p) dp = \int \mathcal{Q}(q) dq = 1 \quad (3.2)$$

Observe that \mathcal{P} represents the (Gaussian) distribution of the momenta p , while \mathcal{Q} represents the Boltzmann partial distribution of the position variables q – containing all information about the potentials $V(q)$. In this setting the probability for the dynamical system to *be* within some subset A of position space can be written as

$$\pi(A) = \int_{\Gamma(A)} f_0(p, q) dq dp = \int_A \mathcal{Q}(q) dq = \int_{\Omega} \chi_A^2 \mathcal{Q}(q) dq =: \langle \chi_A, \chi_A \rangle_{\mathcal{Q}}, \quad (3.3)$$

where we introduced some inner product with weighting \mathcal{Q} .

After this preparation, the operator T is constructed as the restriction of the Perron-Frobenius operator to position space via averaging over the momentum part of the canonical distribution. Once this operator has been defined, the conditional probability for the system to *move* from some subset A to some subset B in position space during time τ can be written as

$$w(A, B, \tau) = \frac{\langle \chi_A, T\chi_B \rangle_{\mathcal{Q}}}{\langle \chi_A, \chi_A \rangle_{\mathcal{Q}}}. \quad (3.4)$$

In the same manner, the probability for the system to *stay* in A during time τ comes out as

$$w(A, A, \tau) = \frac{\langle \chi_A, T\chi_A \rangle_{\mathcal{Q}}}{\langle \chi_A, \chi_A \rangle_{\mathcal{Q}}}. \quad (3.5)$$

The elements of the operator are computed via some hybrid Monte Carlo method.

Perron cluster analysis. Suppose we have k sets of configurations, say 'metastable conformations' $\mathcal{S}_1, \dots, \mathcal{S}_k$, each of which captures the molecular system 'for a long time', once it is in there. For the transition probabilities (3.4) and (3.5) this means that

$$w(\mathcal{S}_i, \mathcal{S}_i, \tau) = 1 - O(\epsilon), \quad w(\mathcal{S}_i, \mathcal{S}_j, \tau) = O(\epsilon), \quad i \neq j \quad (3.6)$$

in terms of some 'small' perturbation parameter not specified here. Assuming a 'reasonable' discretization of the above operator T , the obtained transition matrix T contains the desired information about the metastable sets in the *Perron cluster* eigenproblem associated with the eigenvalues

$$\lambda_1 = 1, \quad \lambda_2 = 1 - O(\epsilon), \dots, \lambda_k = 1 - O(\epsilon).$$

Details are omitted here – see the quoted literature or the recent survey paper [2] by DEUFLHARD (2002). The name 'Perron cluster analysis' characterizes a cluster analysis technique based on some analysis of the arising Perron cluster

of eigenvalues of the transition matrix of a Markov chain. For this reason it is more correctly named as **Perron Cluster Cluster Analysis**, abbreviated as PCCA.

In summary, the PCCA algorithm supplies the following information:

- the probabilities $\pi(\mathcal{S}_i)$ for the system to *be* within each of the subsets \mathcal{S}_i ,
- the probabilities $w_{ii} = w(\mathcal{S}_i, \mathcal{S}_i, \tau)$ for the system to *stay* within subset \mathcal{S}_i during time τ , once it is in there, and
- the probabilities $w(\mathcal{S}_i, \mathcal{S}_j, \tau)$, $i \neq j$, for the system to *move* from subset \mathcal{S}_i to subset \mathcal{S}_j .

In other words: The Perron cluster analysis supplies the number, the life times, and the decay pattern of the metastable chemical conformations. The characteristic life times for each of the \mathcal{S}_i are roughly found to be

$$\tau_{\mathcal{S}_i} \approx \frac{\tau}{1 - w_{ii}} .$$

This formula nicely shows that the blow-up from the deterministic time scale τ to the time scales $\tau_{\mathcal{S}_i}$ of the metastable conformations may be significant.

Example: HIV protease inhibitor VX-478. This molecule is the basis for the anti-AIDS drug Agenerase distributed by Glaxo Wellcome. As is well-known among chemists, HIV is a retrovirus and therefore hard to attack directly by drugs. The HIV protease is an enzyme regulating the passage of HIV through the cell membrane. The here selected molecule has been exactly designed (by Vertex) to inhibit this passage. The molecular data were taken from the public domain Protein Data Bank (PDB).

The conformation analysis yielded $k = 3$ metastable conformations at a virtual temperature of 1400 K in the Boltzmann distribution part \mathcal{P} . At a lower temperature level (1000 K), more substructures came into sight, two of which are shown in Fig. 4. The representation there is via some volume rendering of the corresponding molecular probability density for the system dynamics.

Summarizing, *the mathematical concept of metastable conformations supplies a much deeper understanding of the original intuitive chemical concept of geometrical conformations, especially in view of their dynamical properties.*

Conclusion

The selected comparison of three related concepts from chemistry and mathematics gives a rather homogeneous common picture. As a result of close interdisciplinary interaction over decades, certain intuitive chemical concepts, which were mathematically deficient, have turned into less intuitive, but more precise

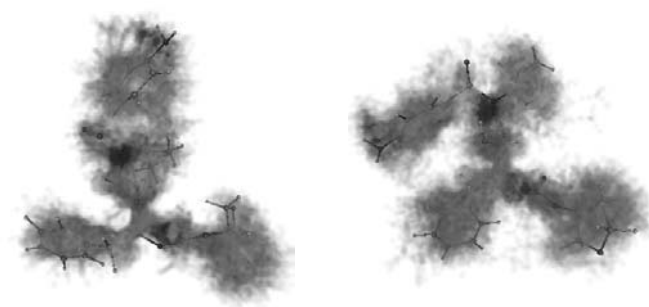


Figure 4: HIV protease inhibitor: T-bone and double T conformations

mathematical concepts. The firmer mathematical basis, in turn, gave rise to further fruitful developments, which now play an important role in chemical research and engineering. Along this line, further progress can be expected in the future from the dialogue between the two scientific languages.

Acknowledgement. This work was supported by the DFG Research Center 'Mathematics for Key Technologies' in Berlin.

References

- [1] DEUFLHARD, P. Differential Equations in Technology and Medicine: Computational Concepts, Adaptive Algorithms, and Virtual Labs. In: R. Burkard, P. Deuffhard, A. Jameson, J.-L. Lions, G. Strang, V. Capasso, H. Engl, J. Periaux (eds.): *Computational Mathematics Driven by Industrial Problems*. Springer Lecture Notes in Mathematics, vol. 1739, 69–126 (2000).
- [2] DEUFLHARD, P. From Molecular Dynamics to Coformation Dynamics in Drug Design. In: M. Kirkilionis, S. Kroemker, R. Rannacher, F. Tomi (eds.): *Trends in Nonlinear Analysis*. Springer Berlin, Heidelberg, New York, 267–286 (2002).
- [3] DEUFLHARD, P., AND BORNEMANN, F. A. *Scientific Computing with Ordinary Differential Equations*. Springer, New York, Texts in Applied mathematics vol. 42 (2002).
- [4] DEUFLHARD, P., DELLNITZ, M., JUNGE, O., AND SCHÜTTE, CH. *Computation of essential molecular dynamics by subdivision techniques*. In: Deuffhard, P., Hermans, J., Leimkuhler, B., Mark, A. E., Reich, S., and Skeel, R. D. (eds). Computational Molecular Dynamics: Challenges, Methods, Ideas. Lecture Notes in Computational Science and Engineering. Springer, Berlin, Heidelberg, New York, vol. 4, pp. 98–115 (1999).
- [5] DEUFLHARD, P., AND HEROTH, J. Dynamic dimension reduction in ODE models. In *Scientific Computing in Chemical Engineering*, F. Keil, W. Mackens, H. Voß, and J. Werther, eds. Springer-Verlag, Berlin, Heidelberg, New York, 29–43 (1996).
- [6] DEUFLHARD, P. HUISINGA, W., FISCHER, A., AND SCHÜTTE, CH. *Identification of almost invariant aggregates in reversible nearly uncoupled Markov chains*. Lin. Alg. Appl. 315, pp. 39–59 (2000).
- [7] DEUFLHARD, P., AND NOWAK, U. Efficient Numerical Simulation and Identification of Large Chemical Reaction Systems. *Ber. Bunsenges.* 90, 940–946 (1986).
- [8] DEUFLHARD, P., AND WULKOW, M. Computational treatment of polyreaction kinetics. *IMPACT Comput. Sci. Engrg.* 1, 269–301 (1989).
- [9] HESSE, H. *The Glass Bead Game*. Penguin Book, 1972. Translated from the German edition *Das Glasperlenspiel*. Fretz and Wasmuth, Zürich (1943).
- [10] HOPPENSTEADT, F. C., ALEFELD, P., AND AIKEN, R. Numerical treatment of rapid chemical kinetics by perturbation and projection methods. In *Modelling of Chemical Reaction Systems*, K. H. Ebert, P. Deuffhard, and W. Jäger, eds. Springer-Verlag, Berlin, Heidelberg, New York, 31–37 (1981).

- [11] KUO, J.C.W., AND WEI, J. A lumping analysis in monomolecular reaction systems: Analysis of approximately lumpable system: Analysis of the exactly lumpable system. *I&EC Fundamentals* 8, 114–123 (1969).
- [12] MAAS, U. Efficient calculation of intrinsic low-dimensional manifolds for the simplification of chemical kinetics. *Computing and Visualization in Science* 1, 69–82 (1998).
- [13] MAAS, U., AND POPE, S. B. Simplifying chemical reaction kinetics: Intrinsic low-dimensional manifolds in composition space. *Combustion and Flame* 88, 239–264 (1992).
- [14] SCHÜTTE, CH. *Conformational Dynamics: Modelling, Theory, Algorithm, and Application to Biomolecules*. Habilitation thesis, Department of Mathematics and Computer Science, Free University of Berlin, 1998. Available as ZIB-Report SC-99-18 via <http://www.zib.de/bib/pub/pw/>.
- [15] SCHÜTTE, C., FISCHER, A., HUISINGA, W., AND DEUFLHARD, P. A direct approach to conformational dynamics based on hybrid Monte Carlo. *J. Comp. Phys.* 151, 146–168 (1999).
- [16] WULKOW, M. The simulation of molecular weight distributions in polyreaction kinetics by discrete Galerkin methods. *Macromol. Theory Simul.* 5, 393–416 (1996).