# LOW RANK APPROXIMATIONS OF TENSORS AND MATRICES: THEORY, APPLICATIONS, PERSPECTIVES

Eugene Tyrtyshnikov

Institute of Numerical Mathematics of Russian Academy of Sciences
Lomonosov Moscow State Univeristy
Moscow Institute of Science and Technology

eugene.tyrtyshnikov@gmail.com

Atlanta, 30 September 2015

# TENSORS $=d$-DIMENSIONAL MATRICES

Tensor $=$ $d$-linear form $=$
$d$-dimensional matrix (array) of size $n_1 \times \ldots \times n_d$:

$$A = [a(i_1, \ldots, i_d)] = [a_{i_1 i_2 \ldots i_d}]$$

# MAIN PROBLEM

Representing a $d$-tensor $A = [a(i_1, \ldots, i_d)]$ of size $n \times \ldots \times n$ by a list of its entries is intractable:

- if $d = 300$ and $n = 2$ then
  the number of elements is $2^{300} \gg 10^{83}$

  greater than atoms in the universe!

# CANONICAL DECOMPOSITION NOW FOR TENSORS

Nonzero tensor with separated variables

$$u(i)v(j)w(k)$$

is called *rank-one tensor* or *skeleton*.

Canonical decomposition = sum of rank-one tensors (skeletons)

$$a(i,j,k) = \sum_{\alpha=1}^{r} u(i,\alpha)v(j,\alpha)w(k,\alpha)$$

Defined by three matrices: $U = [u(:,1), ..., u(:,r)]$,

$$V = [v(:,1), ..., v(:,r)], \quad W = [w(:,1), ..., w(:,r)].$$

# APPLICATIONS OF CANONICAL DECOMPOSITION

- As a model for data, e.g. in spectrometry: given $n_1$ samples of mixed substances, the data = an array of size $n_1 \times n_2 \times n_3$. $n_2$ and $n_3$ – for frequences of emitters and receivers.

  Canonical decompositions reveal the number of substances and concentrations.

- As a main tool in the complexity theory for the computations of bilinear forms (Strassen, Pan, Bini,...).

- Abundant with difficult problems both in theory and computations!

# MINIMAL CANONICAL DECOMPOSITIONS

$k(A) :=$ minimal natural $k$ s.t. *any* $k$ columns of $A$ are linearly independent.

KRUSKAL THEOREM. Assume that $k(U) + k(V) + k(W) \geqslant 2R - 2$. Then the canonical decomposition is minimal and its rank-one tensors (skeletons) are unique.

Recent results by Domanov & De Lathauwer (2013) – weaker minimality conditions. E.g.:

$$k(U_1) + r(U_2) + r(U_3) \geqslant 2R + 2,$$

$$r(U_1) + k(U_2) + r(U_3) \geqslant 2R + 2,$$

$$r(U_1) + r(U_2) + k(U_3) \geqslant 2R + 2.$$

# A SOURCE OF TROUBLE

tensors of rank $< r$ $\rightarrow$ tensor of rank $r$

## CONJECTURE STILL OPEN

For any tensor $A$ of rank $R <$ maximal possible rank, there exists a rank-one tensor $B$ s.t.

$$\mathrm{rank}\,(A + B) = \mathrm{rank}\,(A) + 1.$$

We can prove that it is true *generically* for rank-$R$ tensors.

# REDUCE TENSORS TO MATRICES !!!

For Canonical and Tucker decompositions see,
e.g. Kolda–Bader survey.
But both are of limited use for our purposes (by different reasons).

New decompositions in numerical anaysis:

- TT (Tensor Train)  –  Moscow, INM (2009)
  Oseledets, Tyr.

- HT (Hierarchical Tucker)  –  Leipzig, MPI (2009)
  Hackbusch, Grasedyck, ...

Both use *low-rank matrices*.

Both use the same *dimensionality reduction tree*.

# ASSUME SEPARATION OF VARIABLES

Tensor converts into a matrix (many ways!):

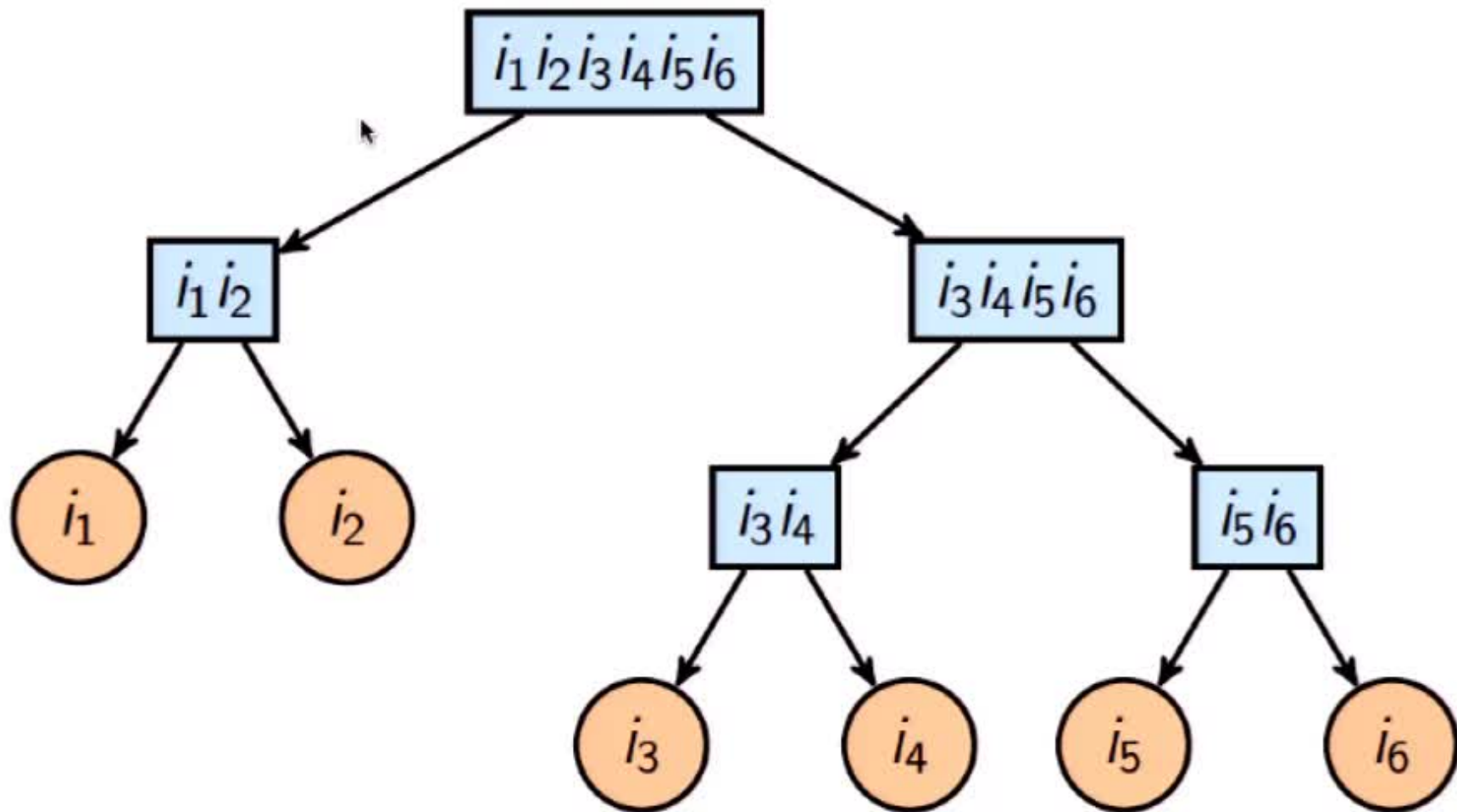$$I = \{1, \ldots, d\} = I_1 \sqcup I_2, \quad b(I_1, I_2) := a(i_1, \ldots, i_d)$$

This matrix is assumed to be of low rank:

$$b(I_1, I_2) = \sum u(I_1, \alpha) v(\alpha, I_2)$$

Next idea is to repeat same for $u(I_1, \alpha)$ and $v(\alpha, I_2)$.

If straightforwardly, then too many $\alpha$'s arise.

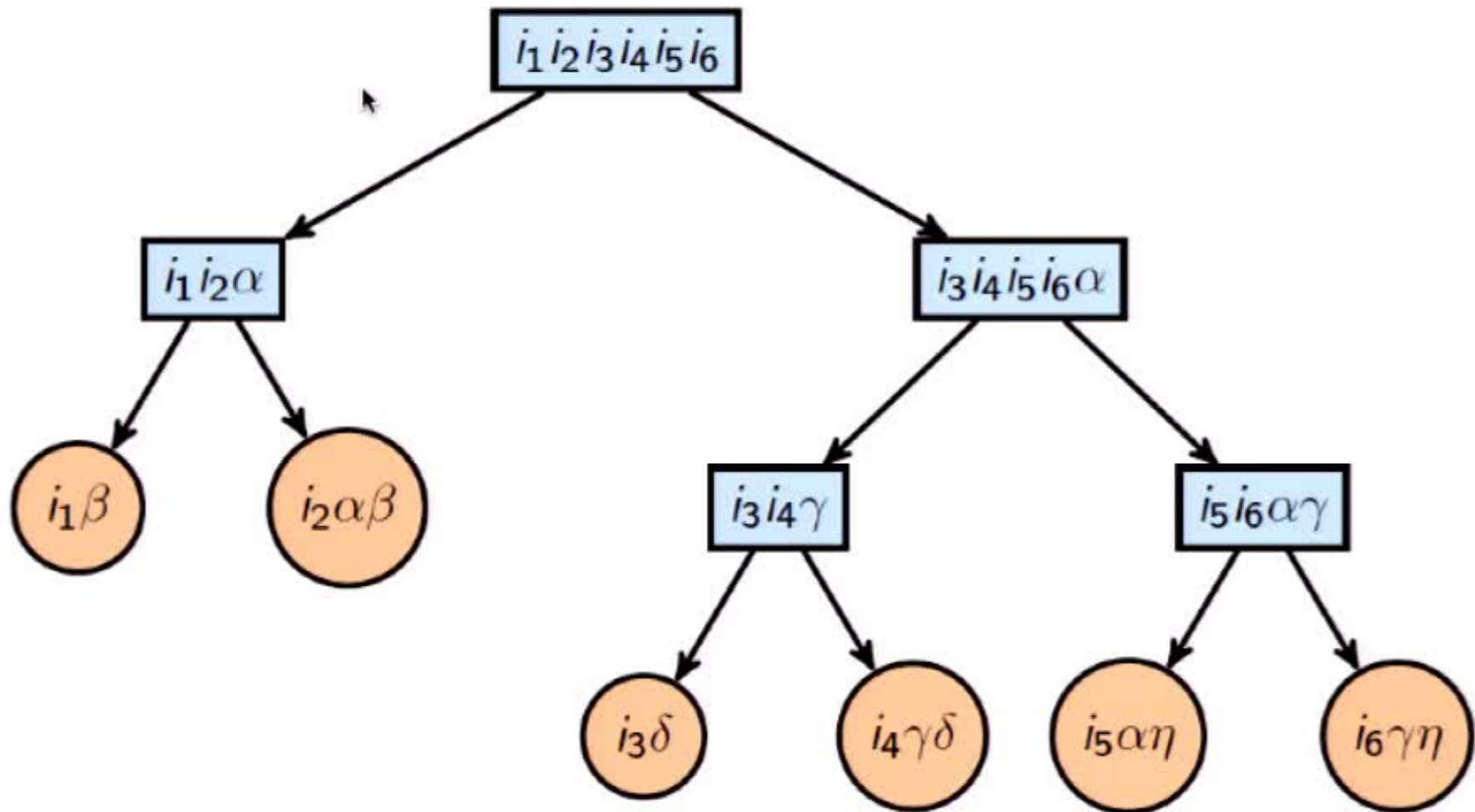# REDUCTION OF DIMENSIONALITY

# THE FIRST STEP IS ESSENTIALLY SAME

$$a(i_1 i_2 ; i_3 i_4 i_5 i_6) = \sum u(i_1 i_2 ; \alpha) v(\alpha ; i_3 i_4 i_5 i_6)$$

Tensor reduces to smaller dimensionality tensors.

The $\alpha$ index is no longer viewed as a parameter!

# SCHEME FOR TT

Auxiliary indices must go to different descendants.

# WHERE TT AND HT START TO DIFFER

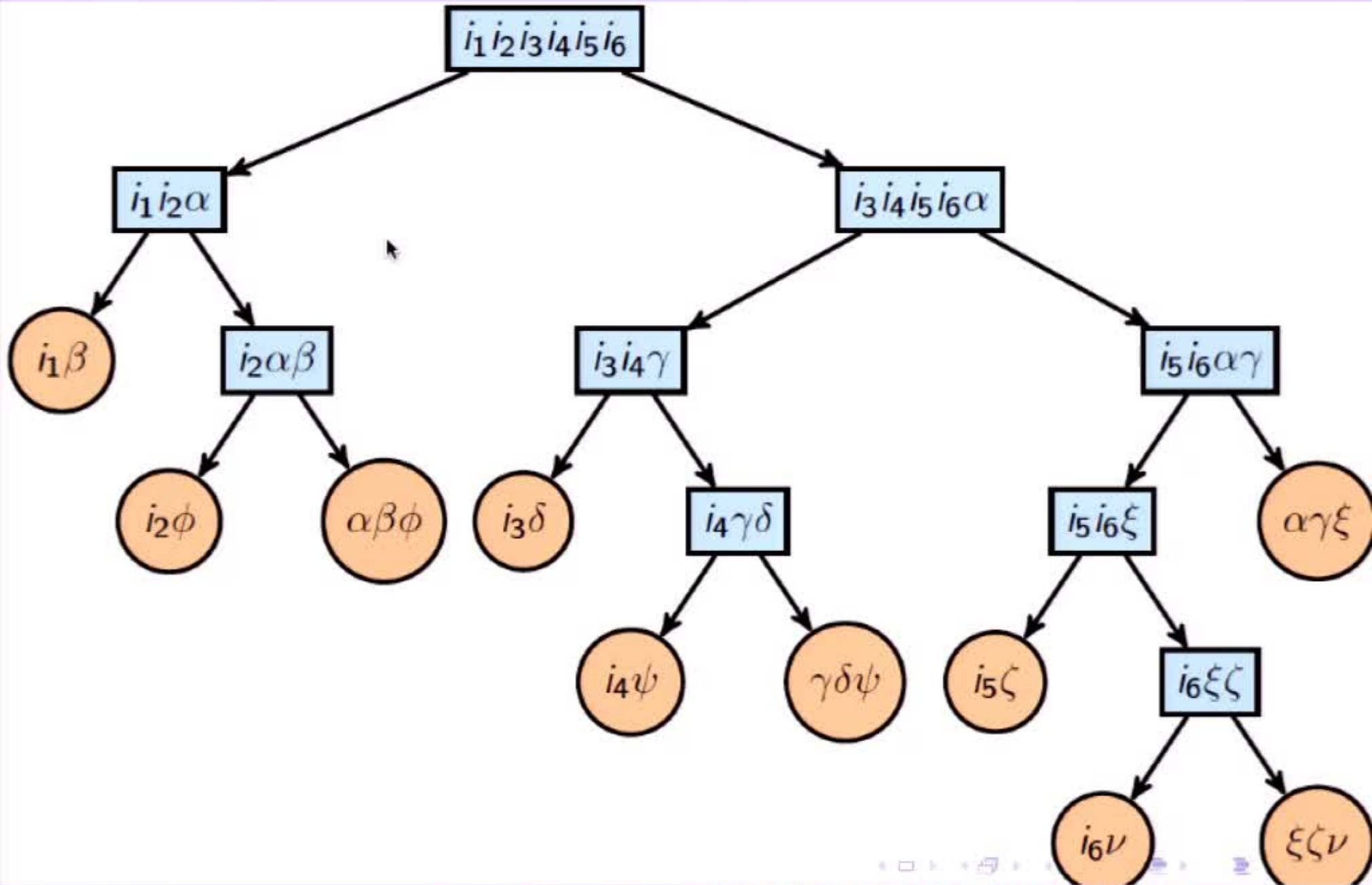In TT, we relegate $\alpha$ and $\gamma$ to different descendants:

$$a(i_5 i_6 \alpha \gamma) = \sum u(i_5 \alpha \, ; \, \eta) v(\eta \, ; \, i_6 \gamma)$$

In HT, we separate $\alpha$ and $\gamma$ from the original indices:

$$a(i_5 i_6 \alpha \gamma) = \sum u(i_5 i_6 \, ; \, \xi) v(\xi \, ; \, \alpha \gamma)$$

The only difference: auxilliary summation indices are treated in different ways!

# SCHEME FOR HT

# TENSOR TRAIN IN $d$ DIMENSIONS

$$a(i_1 \ldots i_d) =$$

$$\sum g_1(i_1\alpha_1)g_2(\alpha_1 i_2\alpha_2)\ldots$$

$$\ldots g_{d-1}(\alpha_{d-2}i_{d-1}\alpha_{d-1})g_d(\alpha_{d-1}i_d)$$

$d$-tensor reduces to 3-tensors $g_k(\alpha_{k-1}i_k\alpha_k)$.

If the maximal size is $r \times n \times r$ then
the number of tensor-train elements does not exceed

$$dnr^2 \ll n^d.$$

# WHAT IS OUR CLASS OF TENSORS?

$$A_k = [a(i_1 \ldots i_k \, ; \, i_{k+1} \ldots i_d)] =$$

$$\left[ \sum u_k(i_1 \ldots i_k \, ; \, \alpha_k) \, v_k(\alpha_k \, ; \, i_{k+1} \ldots i_d) \right] = U_k V_k^\top$$

$$u_k(i_1 \ldots i_k \alpha_k) = \sum g_1(i_1 \alpha_1) \ldots g_k(\alpha_{k-1} i_k \alpha_k)$$

$$v_k(\alpha_k i_{k+1} \ldots i_d) = \sum g_{k+1}(\alpha_k i_{k+1} \alpha_{k+1}) \ldots g_d(\alpha_{k-1} i_d)$$

THE MAIN PROPERTY OF THE CLASS:
all matrices $A_k$ must be (close to) low-rank matrices.

# WHAT IS OUR CLASS OF TENSORS?

THEOREM (Oseledets-Tyr.'2009)

Given a tensor $A$, assume that $\operatorname{rank}(A_k + E_k) = r_k$. Then a tensor train $T$ exists with ranks $r_1, \ldots, r_{d-1}$ s.t.

$$\|A - T\|_F \leqslant \sqrt{\sum_{i=1}^{d-1} \|E_k\|_2^2}$$

L. Graesedyck: a similar result for HT.

# EVERYTHING REDUCES TO MATRICES

Tensor train can be viewed as a *rank-structured representation* for matrices $A_1, \ldots, A_{d-1}$. Structured SVD can be computed for them simultaneously just in $O(dnr^3)$ operations!

Tensor train can be constructed if we know low-rank decompositons for matrices $A_1, \ldots, A_{d-1}$.

**Moreover, it can be construced from cleverly chosen *crosses* in some small submatrices of those matrices.**

# INTERPOLATION ERROR

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} - \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} A_{11}^{-1} \begin{bmatrix} A_{11} & A_{12} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 \\ 0 & A_{22} - A_{21} A_{11}^{-1} A_{12} \end{bmatrix}$$

# MAXIMAL VOLUME PRINCIPLE

It is still not illegal to teach determinants

**THEOREM** (Goreinov, Tyr.)  *Let*

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad A_{11} \text{ is } r \times r$$

*with maximal volume (determinant in modulus) among all $r \times r$ blocks in A, and set*

$$A_r = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} A_{11}^{-1} \begin{bmatrix} A_{11} & A_{12} \end{bmatrix}.$$

*Then*

$$\|A - A_r\|_C \leqslant (r+1)^2 \min_{\operatorname{rank} B \leqslant r} \|A - B\|_C.$$

# MAXIMAL VOLUME PRINCIPLE

**PREVIOUS RESULT** (Goreinov, Tyr.'2000)

$$\|A - A_r\|_C \leqslant (r + 1) \min_{\mathrm{rank}\, B \leqslant r} \|A - B\|_2 = \sigma_{r+1}(A).$$

Coming soon: generalizations for using larger or even rectangular cross-intersection blocks.

# PROOF

$$Q = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \\ q_{r+1,1} & \cdots & q_{r+1,r} \\ \cdots & \cdots & \cdots \\ q_{n1} & \cdots & q_{nr} \end{bmatrix}$$

Necessary for the maximal volume:

$$|q_{ij}| \leqslant 1, \quad r+1 \leqslant i \leqslant n, \quad 1 \leqslant j \leqslant r.$$

Otherwise, swapping the rows increases the volume!

# CROSS INTERPOLATION HISTORY

1985 Knuth:  Semi-optimal bases for linear dependencies

1995 Tyr., Goreinov, Zamarashkin:  $A = CGR$ pseudoskeleton

2000 Tyr.:  incomplete cross approximation with ALS maxvol

2000 Bebendorf:  ACA = Gaussian elimination

2001 Tyr., Goreinov:  maximum volume principle,
quasioptimality $\| \text{cross} \|_C \le (r+1) \| \text{best} \|_2$

2006 Mahoney et al:  randomized $CUR$ algorithm

2008 Oseledets, Savostyanov, Tyr.:  Cross3D

2009 Oseledets, Tyr.:  TT-Cross

2010 J.Schneider:  function-related quasioptimality
$\| \text{cross} \|_C \le (r+1)^2 \| \text{best} \|_C$

2011 Tyr., Goreinov:  quasioptimality
$\| \text{cross} \|_C \le (r+1)^2 \| \text{best} \|_C$

2013 Ballani, Grasedyck, Kluge:  HT-Cross

2013 Townsend, Trefethen -- Chebfun2

# *What do YOU possess that you have not RECEIVED?*

TT, HT and some algorithms (most famous is DMRG by White) can be found in theoretical physics.

A useful outcome with the new name are some
NEW ALGORITHMS:

- ▸ TT-CROSS (Oseledets, Tyr., 2009)

- ▸ Wavelet-TT (Oseledets, Tyr., 2011)

- ▸ AMEn (Dolgov & Savostyanov, 2013)

# TT-CROSS

Seek crosses in the unfolding matrices. Let $a_1 = a(i_1, i_2, i_3, i_4)$.
On input: $r$ initial columns in each. Select *good* rows.

$$A_1 = [a(i_1\,;\ i_2, i_3, i_4)], \quad J_1 = \{i_2^{(\beta_1)} i_3^{(\beta_1)} i_4^{(\beta_1)}\}$$

$$A_2 = [a(i_1, i_2\,;\ i_3, i_4)], \quad J_2 = \{i_3^{(\beta_2)} i_4^{(\beta_2)}\}$$

$$A_3 = [a(i_1, i_2, i_3\,;\ i_4)], \quad J_3 = \{i_4^{(\beta_3)}\}$$

| rows | matrix | skeleton decomposition |
|---|---|---|
| $l_1 = \{i_1^{(\alpha_1)}\}$ | $a_1(i_1\,;\ i_2, i_3, i_4)$ | $a_1 = \sum_{\alpha_1} g_1(i_1; \alpha_1)\, a_2(\alpha_1;\ i_2, i_3, i_4)$ |
| $l_2 = \{i_1^{(\alpha_2)} i_2^{(\alpha_2)}\}$ | $a_2(\alpha_1, i_2\,;\ i_3, i_4)$ | $a_2 = \sum_{\alpha_2} g_2(\alpha_1, i_2;\ \alpha_2)\, a_3(\alpha_2, i_3;\ i_4)$ |
| $l_3 = \{i_1^{(\alpha_3)} i_2^{(\alpha_3)} i_3^{(\alpha_3)}\}$ | $a_3(\alpha_2, i_3\,;\ i_4)$ | $a_3 = \sum_{\alpha_3} g_3(\alpha_2, i_3;\ \alpha_3)\, g_4(\alpha_3;\ i_4)$ |

$$a = \sum_{\alpha_1, \alpha_2, \alpha_3, \alpha_4} g_1(i_1, \alpha_1)\, g_2(\alpha_1, i_2, \alpha_2)\, g_3(\alpha_2, i_3, \alpha_3)\, g_4(\alpha_3, i_4)$$

# TENSOR TRAIN COMES

## FROM SMALL CROSSES IN THE UNFOLDING MATRICES

$$A(i_1 \ldots i_d) = \prod_{k=1}^{d} A(J_{\leqslant k-1}, \, i_k, \, J_{>k}) \left[ A(J_{\leqslant k}, \, J_{>k}) \right]^{-1}$$

# PSEUDO-QUASI-OPTIMALITY RESULT

THEOREM (Savostyanov'2013)

Assume that a d-tensor $A$ is approximated by $\widetilde{A}$ on the maximal volume crosses in the unfolding matrices, and let the error is upper bounded by $\varepsilon \|A\|_C$ in each matrix. Then for sufficiently small $\varepsilon$ we have

$$\|A - \widetilde{A}\|_C \leqslant 2dr\varepsilon\|A\|_C.$$

# WHAT HAMLET WOULD SAY

Tensor Train = Matrix Product State = Linear Tensor Network

$$a(i_1, i_2, i_3, i_4, i_5) =$$

$$\sum g_1(i_1, \alpha_1) g_2(\alpha_1, i_2, \alpha_2) g_3(\alpha_2, i_3, \alpha_3) g_4(\alpha_3, i_4, \alpha_4) g_5(\alpha_4, i_3)$$

$$= \underbrace{A_1^{(i_1)}}_{1 \times r_1} \underbrace{A_2^{(i_2)}}_{r_1 \times r_2} \underbrace{A_3^{(i_3)}}_{r_2 \times r_3} \underbrace{A_4^{(i_4)}}_{r_3 \times r_4} \underbrace{A_5^{(i_5)}}_{r_4 \times 1}$$
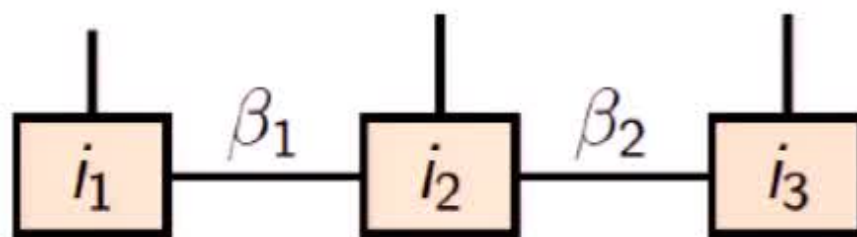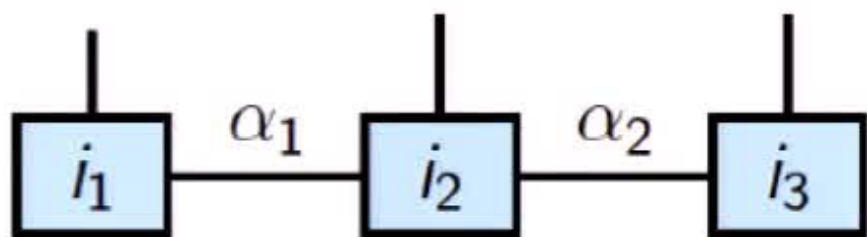
# EASY OPERATIONS ON TENSORS

e.g. summation

$$a(i_1, i_2, i_3) = \underbrace{A_1^{(i_1)}}_{1 \times r_1} \underbrace{A_2^{(i_2)}}_{r_1 \times r_2} \underbrace{A_3^{(i_3)}}_{r_2 \times 1}, \qquad b(i_1, i_2, i_3) = \underbrace{B_1^{(i_1)}}_{1 \times s_1} \underbrace{B_2^{(i_2)}}_{s_1 \times s_2} \underbrace{B_3^{(i_3)}}_{s_2 \times 1}$$



$$(a + b)(i_1, i_2, i_3) = \begin{bmatrix} A_1^{(i_1)} & B_1^{(i_1)} \end{bmatrix} \begin{bmatrix} A_2^{(i_2)} & \\ & B_2^{(i_2)} \end{bmatrix} \begin{bmatrix} A_3^{(i_3)} \\ B_3^{(i_3)} \end{bmatrix}$$

# A NEW PARADIGM OF COMPUTATIONS
only through low-parametric formats

- $A = A(p)$, $B = B(q)$, $C = C(s)$
- To implement $C = A * B$ we should devise *fast algorithms* for getting $s$ from $p$ and $q$.

- *General algebraic method* for a wide class of applications!

- We can use classical methods of numerical analysis together with TT-approximation.

# AND EVEN NEWER APPROACH
## MINIMIZE THE ERROR FUNCTIONAL ON THE TT MANIFOLD

AMEn does it for quadratic functionals.

$$A = A_1 A_2 \ldots A_d \quad \text{(approximation)}$$

$$B = B_1 B_2 \ldots B_d \quad \text{(gradient)}$$

$$A := A_1 \ldots A_{i-2} \, [A_{i-1}, B_{i-1}] \begin{bmatrix} A_i \\ 0 \end{bmatrix} A_{i+1} \ldots A_d$$

# TENSORIZATION OF VECTORS

Any vector of size $mn$ can be viewed as a matrix of size $m \times n$.

Any vector of size $N = n_1 \ldots n_d$ can be viewed as a $d$-dimensional matrix (tensor) of size $n_1 \times \ldots \times n_d$.

# TENSORIZATION OF MATRICES

Let $N = n_1 \ldots n_d$. Any matrix of size $N \times N$

$$a(i, j) = a(i_1 \ldots i_d, \; j_1 \ldots j_d)$$

can be viewed as a $2d$-tensor, and as a $d$-tensor, e.g.

$$a(i_1 j_1, \ldots, i_d j_d)$$

of size $n_1^2 \times \ldots \times n_d^2$.

## Tensorization with TT may crucially decrease the number of representation parameters!

# EXAMPLES OF TENSORIZATION

$f(x)$ is a function on $[0, 1]$

$$a(i_1, \ldots, i_d) = f(ih), \quad i = \frac{i_1}{2} + \frac{i_2}{2^2} + \cdots + \frac{i_d}{2^d}$$

The array of values of $f$ is viewed as a tensor of size $2 \times \cdots \times 2$.

EXAMPLE 1. $f(x) = e^x + e^{2x} + e^{3x}$
ttrank= 2.7     ERROR=1.5e-14

EXAMPLE 2. $f(x) = 1 + x + x^2 + x^3$
ttrank= 3.4     ERROR=2.4e-14

EXAMPLE 3. $f(x) = 1/(x - 0.1)$

ttrank= 10.1     ERROR=5.4e-14

# EXAMPLE OF TT INTERPOLATION

$f(x) = \frac{\sin(1000x)}{\sqrt{x}}$ on $[0, 1000]$.

$2^{63}$ nodes.

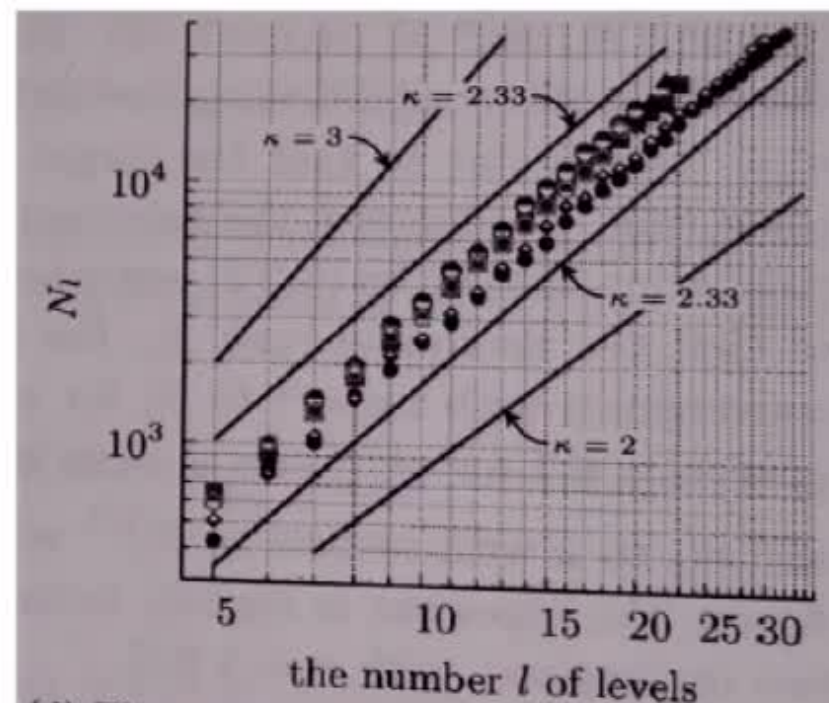| tol | $\|E\|_F$ | $\|E\|_C$ | Number of function calls | TT-rank |
|---|---|---|---|---|
| $10^{-2}$ | $6.46 * 10^{-2}$ | $1.39 * 10^{-3}$ | 73668 | 3.00 |
| $10^{-4}$ | $1.71 * 10^{-4}$ | $3.11 * 10^{-5}$ | 164729 | 4.67 |
| $10^{-6}$ | $2.12 * 10^{-6}$ | $1.82 * 10^{-7}$ | 288449 | 6.18 |
| $10^{-7}$ | $1.54 * 10^{-7}$ | $3.54 * 10^{-8}$ | 348201 | 6.88 |

# TT-FE APPROXIMATION

$$u_\Gamma(x) = r^\alpha \sin \alpha \phi(x), \quad x \in \Omega = (0,1)^2$$

$$\varepsilon_l \leqslant \exp\{-cN_l^{\frac{1}{\kappa}}\}, \quad N_l - \text{the number of TT-elements}$$

THEOREM (V.Kazeev & C.Schwab). $\kappa \leqslant 5$.



(a) Convergence with respect to $l$. The reference lines correspond to the exponential convergence $\varepsilon_l = C_\alpha 2^{-\tilde{\alpha} l}$ with $C_\alpha$ independent of $l$ and with $\bar{\alpha} = \min\{\alpha, 1\}$.
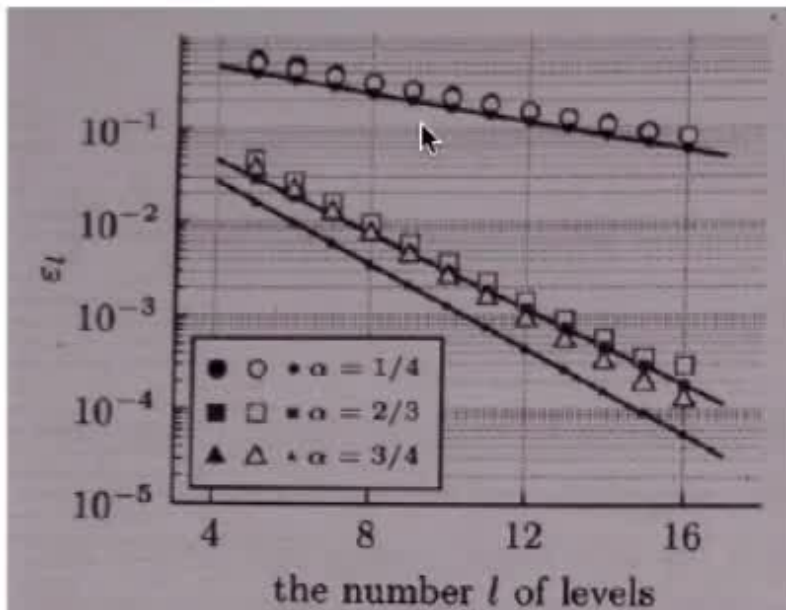
(d) The number $N_l$ (3.3.2) of QTT parameters vs. $l$. The reference lines correspond to the algebraic growth $N_l = C_\alpha l^\kappa$ with $\kappa$ and $C_\alpha$ independent of $l$.
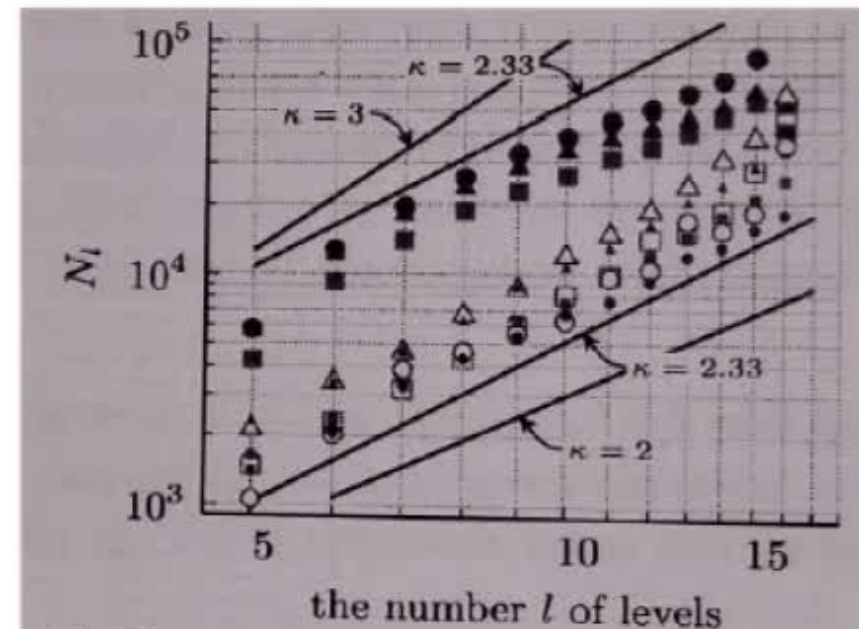
$$\Omega = (0,1)^2 \setminus [0,1) \times (-1,0] \qquad \Omega = (0,1)^2 \setminus [0,1) \times \{0\}$$

bright labels $\qquad\qquad\qquad\qquad$ black labels



(a) Convergence with respect to $l$. The reference lines correspond to the exponential convergence $\varepsilon_l = C_\alpha 2^{-\alpha l}$ with $C_\alpha$ independent of $l$. The markers for $u_{sol}^l$ and $u_{tr}^l$ mostly coincide.

(d) The number $N_l$ (3.3.2) of QTT parameters vs. $l$. The reference lines correspond to the algebraic growth $N_l = C_\alpha l^\kappa$ with $\kappa$ and $C_\alpha$ independent of $l$.

(V.Kazeev)

# TENSOR TRAINS
# FOR PARAMETRIC EQUATIONS

Diffusion domain $= [0, 1]^2$ consists of $p \times p$ square subdomains with constant diffusion coefficient, $p^2$ parameters varying from 0.1 to 1.

256 knots in each parameter.    Space grid of size $256 \times 256$.

Solutions *for all values of parameters* are approximated by a tensor train with relative accuracy $10^{-5}$:

| Number of parameters | Memory |
|:---:|:---:|
| 4 | 8 Mb |
| 16 | 24 Mb |
| 64 | 78 Mb |

(I.Oseledets)

# LOW RANK AND TENSOR TRAINS IN THE SMOLUCHOWSKI EQUATIONS

- $\overline{v} = (v_1, \ldots, v_d)$ – volumes of different substances of a particle
- $t$ – time
- $n(\overline{v}, t)$ – concentration function for volume components of a particle

$$\frac{\partial n(\overline{v}, t)}{\partial t} = \frac{1}{2} \int_0^{v_1} du_1 \ldots \int_0^{v_d} K(\overline{v} - \overline{u}; \overline{u}) n(\overline{u}, t) n(\overline{v} - \overline{u}, t) du_d -$$

$$-n(\overline{v}, t) \int_0^{\infty} du_1 \ldots \int_0^{\infty} K(\overline{v}; \overline{u}) n(\overline{u}, t) du_d,$$

$$n(\overline{v}, 0) = n_0(\overline{v}).$$

Joint work with S.Matveev and A.Smirnov

# GLOBAL SEARCH CAN BADLY GAIN WHEN IT USES TENSOR TRAIN

THEOREM. *If $A_\blacksquare$ is of maximal volume among $r \times r$ blocks in $A$, then*

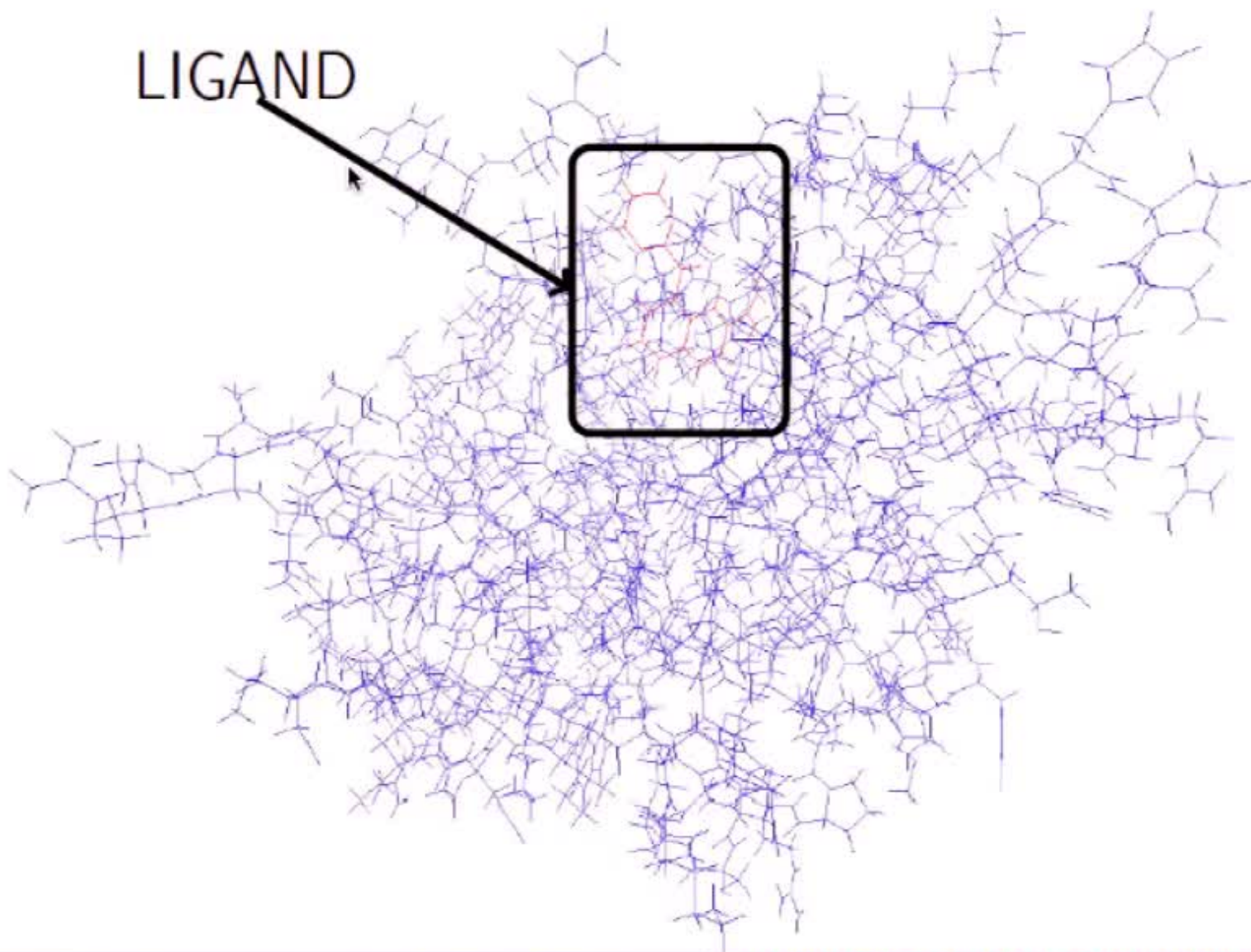$$\|A_\blacksquare\|_C \geq \|A\|_C/(2r^2 + r).$$

S. Goreinov, I. Oseledets, D. Savostyanov, E. Tyrtyshnikov, N. Zamarashkin, How to find a good submatrix, *Matrix Methods: Theory, Algorithms and Applications. Devoted to the Memory of Gene Golub* (eds. V.Olshevsky and E.Tyrtyshnikov), World Scientific Publishers, Singapore, 2010, pp. 247–256.
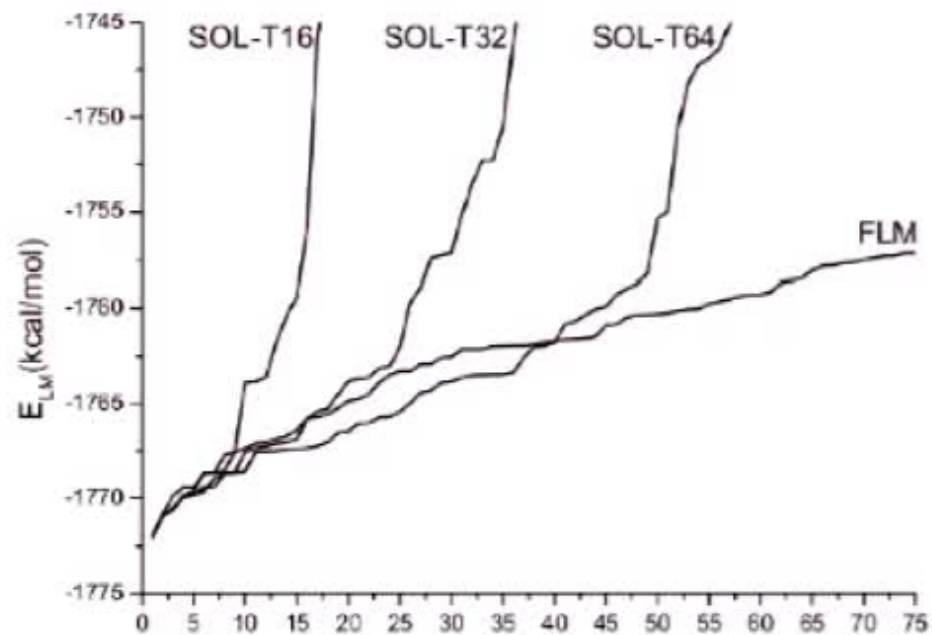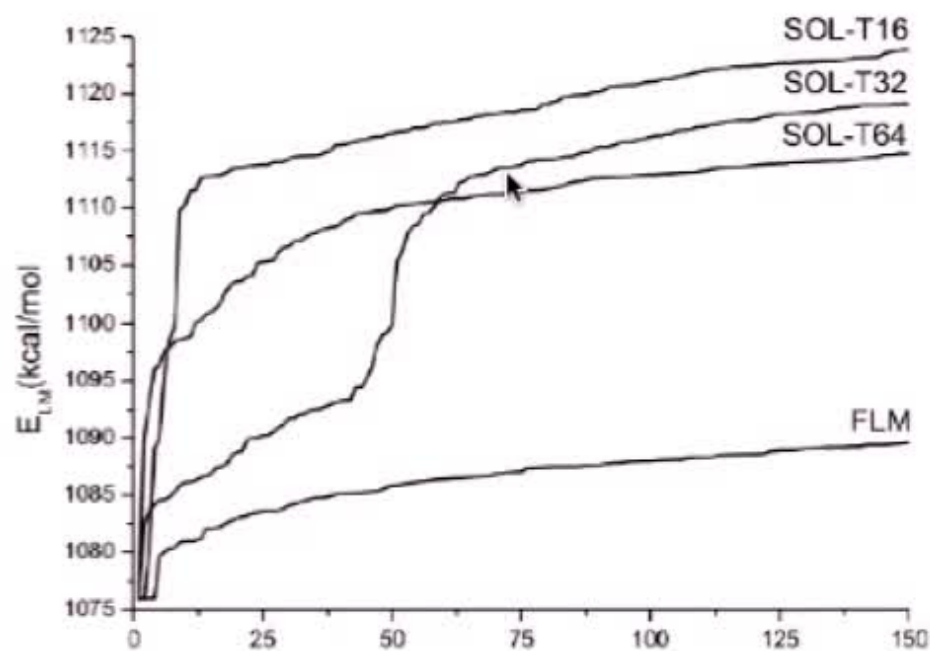
# DIRECT DOCKING IN THE DRUG DESIGN

ACCOMMODATION OF LIGAND INTO PROTEIN

LIGAND

# DIRECT DOCKING IN THE DRUG DESIGN



Joint work with D.Zheltkov and V.Sulimov