

Paralelní a distribuované výpočty (B4B36PDV)

Jakub Mareček, Michal Jakob

`jakub.marecek@fel.cvut.cz`

Artificial Intelligence Center
Department of Computer Science
Faculty of Electrical Engineering
Czech Technical University in Prague

Dnešní přednáška

Techniky paralelizace 3

Chci paralelizovat

učení neuronky



Jak na to?

Dnešní přednáška

Techniky paralelizace 3

Chci paralelizovat maticový algoritmus

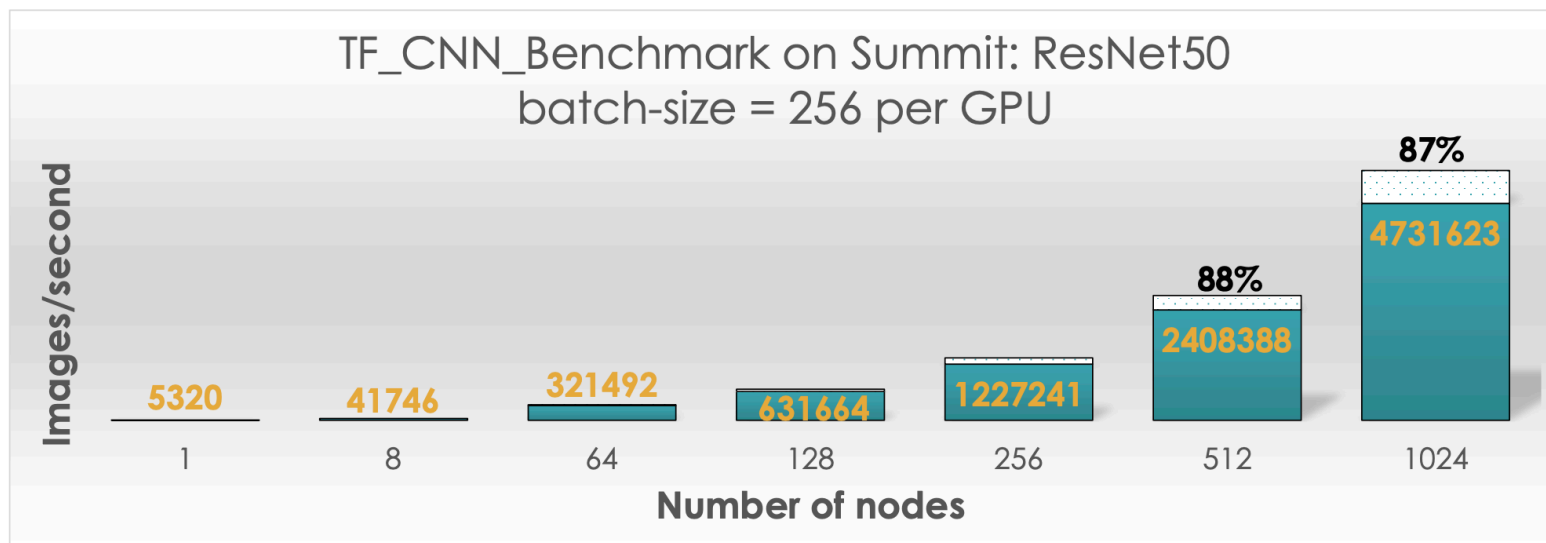


Jak na to?

Demotivace

Prototypování numerických algoritmů

Performance baselines: ResNet50 on ImageNet



Nodes	Mini-batch size	Top-1 Val accuracy	Training time (min)
16	12288	0.750	27
32	12288	0.766	17
64	15360	0.763	12

TF.distribute & Horovod + LARS
[code: TensorFlow distributed example](#)

Motivace

Prototypování numerických algoritmů

```
#include <iostream>
#include <vector>
#include "mkldnn.hpp"

using namespace mkldnn;

...

int main(int argc, char **argv) {
    try {
        simple_net();
        std::cout << "ok\n";
    } catch (error &e) {
        std::cerr << "status: " << e.status << std::endl;
        std::cerr << "message: " << e.message << std::endl;
    }
    return 0;
}
```

https://oneapi-src.github.io/oneDNN/v1.0/cpu_rnn_inference_int8_cpp.html

Motivace

Prototypování numerických algoritmů

```
#include <iostream>
#include <vector>
#include "mkldnn.hpp"

using namespace mkldnn;
using dim_t = mkldnn::memory::dim;

int main(int argc, char **argv) {
    try {
        const dim_t n = 64;
        // column-major order (sloupce souvisle)
        std::vector<float> A(n*n, 1.0f);
        std::vector<float> B(n*n, 1.0f);
        std::vector<float> C(n*n, 1.0f);
        // https://oneapi-src.github.io/oneDNN/v0/group\_\_c\_\_api\_\_blas.html
        mkldnn_status_t status = mkldnn_sgemv('N', 'N', n, n, n,
            1.f, A.data(), n, B.data(), n,
            0.f, C.data(), n);
        std::cerr << "status: " << status << std::endl;
    } catch (error &e) {
        std::cerr << "status: " << e.status << std::endl;
        std::cerr << "message: " << e.message << std::endl;
    }
    return 0;
}
```

Motivace

Prototypování numerických algoritmů

Developer Reference for Intel® oneAPI Math Kernel Library - C

Developer Reference

Version: 2021.2

Last Updated: 03/26/2021

Public Content

[Download as PDF](#)



Developer Reference for Intel® oneAPI Math Kernel Library

[Getting Help and Support](#)

[What's New](#)

[Notational Conventions](#)

> [Overview](#)

> [OpenMP* Offload](#)

> [BLAS and Sparse BLAS Routines](#)

> [BLAS Routines](#)

[Routine Naming Conventions](#)

[C Interface Conventions](#)

cblas_?gemv

Computes a matrix-vector product using a general matrix.

Syntax

```
void cblas_sgemv (const CBLAS_LAYOUT Layout, const CBLAS_TRANSPOSE trans, const MKL_INT m, const MKL_INT n, const MKL_INT incx, const float alpha, const float *a, const MKL_INT lda, const float *x, const MKL_INT incx, const float beta, const float *y, const MKL_INT incy);
```

```
void cblas_dgemv (const CBLAS_LAYOUT Layout, const CBLAS_TRANSPOSE trans, const MKL_INT m, const MKL_INT n, const MKL_INT incx, const double alpha, const double *a, const MKL_INT lda, const double *x, const MKL_INT incx, const double beta, const double *y, const MKL_INT incy);
```

```
void cblas_cgemv (const CBLAS_LAYOUT Layout, const CBLAS_TRANSPOSE trans, const MKL_INT m, const MKL_INT n, const MKL_INT incx, const void *alpha, const void *a, const MKL_INT lda, const void *x, const MKL_INT incx, const void *beta, const void *y, const MKL_INT incy);
```

Motivace

Prototypování numerických algoritmů

Level 1 BLAS

	dim	scalar	vector	vector	scalars	5-element array
SUBROUTINE xROTG (A, B, C, S)	
SUBROUTINE xROTMG(D1, D2, A, B,	PARAM)	
SUBROUTINE xROT (N,		X, INCX, Y, INCY,			C, S)	
SUBROUTINE xROTM (N,		X, INCX, Y, INCY,			PARAM)	
SUBROUTINE xSWAP (N,		X, INCX, Y, INCY)				
SUBROUTINE xSCAL (N,	ALPHA,	X, INCX)				
SUBROUTINE xCOPY (N,		X, INCX, Y, INCY)				
SUBROUTINE xAXPY (N,	ALPHA,	X, INCX, Y, INCY)				
FUNCTION xDOT (N,		X, INCX, Y, INCY)				
FUNCTION xDOTU (N,		X, INCX, Y, INCY)				
FUNCTION xDOTC (N,		X, INCX, Y, INCY)				
FUNCTION xxDOT (N,		X, INCX, Y, INCY)				
FUNCTION xNRM2 (N,		X, INCX)				
FUNCTION xASUM (N,		X, INCX)				
FUNCTION IxAMAX(N,		X, INCX)				

Generate plane rotation
Generate modified plane rotation
Apply plane rotation
Apply modified plane rotation
 $x \leftrightarrow y$
 $x \leftarrow \alpha x$
 $y \leftarrow x$
 $y \leftarrow \alpha x + y$
 $dot \leftarrow x^T y$
 $dot \leftarrow x^T y$
 $dot \leftarrow x^H y$
 $dot \leftarrow \alpha + x^T y$
 $nrm2 \leftarrow \|x\|_2$
 $asum \leftarrow \|re(x)\|_1 + \|im(x)\|_1$
 $amax \leftarrow 1^{st} k \ni |re(x_k)| + |im(x_k)|$
 $= \max(|re(x_i)| + |im(x_i)|)$

prefixes
S, D
S, D
S, D
S, D
S, D, C, Z
S, D, C, Z, CS, ZD
S, D, C, Z
S, D, C, Z
S, D, DS
C, Z
C, Z
SDS
S, D, SC, DZ
S, D, SC, DZ
S, D, C, Z

Level 2 BLAS

	options	dim	b-width	scalar	matrix	vector	scalar	vector
xGEMV (TRANS,	M, N,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xGEMV (TRANS,	M, N, KL, KU,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xHEMV (UPLO,		N,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xHEMV (UPLO,		N, K,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xHPMV (UPLO,		N,		ALPHA, AP, X, INCX,	BETA, Y, INCY)			
xSYMV (UPLO,		N,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xSBMV (UPLO,		N, K,		ALPHA, A, LDA, X, INCX,	BETA, Y, INCY)			
xSPMV (UPLO,		N,		ALPHA, AP, X, INCX,	BETA, Y, INCY)			
xTRMV (UPLO, TRANS, DIAG,		N,		A, LDA, X, INCX)				
xTBMV (UPLO, TRANS, DIAG,		N, K,		A, LDA, X, INCX)				
xTPMV (UPLO, TRANS, DIAG,		N,		AP, X, INCX)				
xTRSV (UPLO, TRANS, DIAG,		N,		A, LDA, X, INCX)				
xTBSV (UPLO, TRANS, DIAG,		N, K,		A, LDA, X, INCX)				
xTPSV (UPLO, TRANS, DIAG,		N,		AP, X, INCX)				
	options	dim	scalar	vector	vector	matrix		
xGER (M, N,	ALPHA, X, INCX, Y, INCY,	A, LDA)				
xGERU (M, N,	ALPHA, X, INCX, Y, INCY,	A, LDA)				
xGERC (M, N,	ALPHA, X, INCX, Y, INCY,	A, LDA)				
xHER (UPLO,		N,	ALPHA, X, INCX,	A, LDA)				
xHPR (UPLO,		N,	ALPHA, X, INCX,	AP)				
xHER2 (UPLO,		N,	ALPHA, X, INCX, Y, INCY,	A, LDA)				
xHPR2 (UPLO,		N,	ALPHA, X, INCX, Y, INCY,	AP)				
xSYR (UPLO,		N,	ALPHA, X, INCX,	A, LDA)				
xSPR (UPLO,		N,	ALPHA, X, INCX,	AP)				
xSYR2 (UPLO,		N,	ALPHA, X, INCX, Y, INCY,	A, LDA)				
xSPR2 (UPLO,		N,	ALPHA, X, INCX, Y, INCY,	AP)				

$y \leftarrow \alpha Ax + \beta y, y \leftarrow \alpha A^T x + \beta y, y \leftarrow \alpha A^H x + \beta y, A - m \times n$
 $y \leftarrow \alpha Ax + \beta y, y \leftarrow \alpha A^T x + \beta y, y \leftarrow \alpha A^H x + \beta y, A - m \times n$
 $y \leftarrow \alpha Ax + \beta y$
 $y \leftarrow \alpha Ax + \beta y$
 $y \leftarrow \alpha Ax + \beta y$
 $y \leftarrow \alpha Ax + \beta y$
 $y \leftarrow \alpha Ax + \beta y$
 $y \leftarrow \alpha Ax + \beta y$
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$
 $x \leftarrow Ax, x \leftarrow A^T x, x \leftarrow A^H x$
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$
 $x \leftarrow A^{-1} x, x \leftarrow A^{-T} x, x \leftarrow A^{-H} x$
 $A \leftarrow \alpha xy^T + A, A - m \times n$
 $A \leftarrow \alpha xy^T + A, A - m \times n$
 $A \leftarrow \alpha xy^H + A, A - m \times n$
 $A \leftarrow \alpha xx^H + A$
 $A \leftarrow \alpha xx^H + A$
 $A \leftarrow \alpha xy^H + y(\alpha x)^H + A$
 $A \leftarrow \alpha xy^H + y(\alpha x)^H + A$
 $A \leftarrow \alpha xx^T + A$
 $A \leftarrow \alpha xx^T + A$
 $A \leftarrow \alpha xy^T + \alpha yx^T + A$
 $A \leftarrow \alpha xy^T + \alpha yx^T + A$

S, D, C, Z
S, D, C, Z
C, Z
C, Z
C, Z
S, D
S, D
S, D
S, D, C, Z
S, D, C, Z
S, D, C, Z
S, D, C, Z
S, D, C, Z
S, D, C, Z
S, D
S, D
S, D
S, D
S, D
S, D
S, D
S, D

Dnešní přednáška

Techniky paralelizace 3

Chci paralelizovat maticový algoritmus



Jak na to?

Maticové operace

Násobení matice vektorem

a11	a12	a13	a14	a15		x1	y1
a21	...					x2	y2
...					X	x3	y3
						x4	y4
				a55		x5	y5

$$y_1 = \sum_{\{i=1,\dots,5\}} a_{1i} \cdot x_i$$

Maticové operace

Násobení matice vektorem

- Jak paralelizovat?

a11	a12	a13	a14	a15		x1	y1
a21	...					x2	y2
...					X	x3	y3
						x4	y4
				a55		x5	y5

$$y_1 = \sum_{\{i=1,\dots,5\}} a_{1i} \cdot x_i$$

Maticové operace

Násobení matice vektorem

- Zřejmá paralelizace – výpočet každé složky y je nezávislý a všechny mohou být spočteny paralelně

P1	→	a11	a12	a13	a14	a15		x1	y1
P2	→	a21	...					x2	y2
		...						x3	y3
...								x4	y4
						a55		x5	y5

X =

$$y_1 = \sum_{\{i=1,\dots,5\}} a_{1i} \cdot x_i$$

Maticové operace

Násobení matice vektorem

- Zřejmá paralelizace – výpočet každé složky y je nezávislý a všechny mohou být spočteny paralelně

```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)

#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)
    for (int i=0; i<ROWS; i++) {
        for (int j=0; j<COLS; j++) {
            y[i] += A[i * COLS + j]*x[j];
        }
    }
}
```

Maticové operace

Násobení matice vektorem

- Zřejmá paralelizace – výpočet každé složky y je nezávislý a všechny mohou být spočteny paralelně

```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)

#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)
    for (int i=0; i<ROWS; i++) {
        for (int j=0; j<COLS; j++) {
            y[i] += A[i * COLS + j]*x[j];
        }
    }
}
```

- Co můžeme zlepšit?

Maticové operace

Násobení matice vektorem

- Zřejmá paralelizace – výpočet každé složky y je nezávislý a všechny mohou být spočteny paralelně

```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)

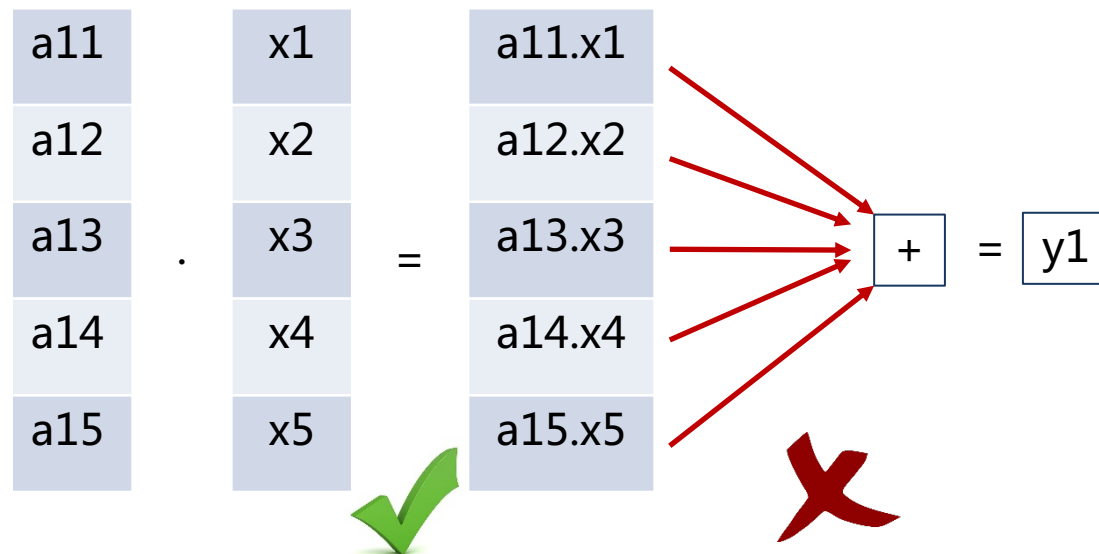
#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)
    for (int i=0; i<ROWS; i++) {
        for (int j=0; j<COLS; j++) {
            y[i] += A[i * COLS + j]*x[j];
        }
    }
}
```

- Co můžeme zlepšit?
- Lze využít vektorizaci?

Maticové operace

Násobení matice vektorem

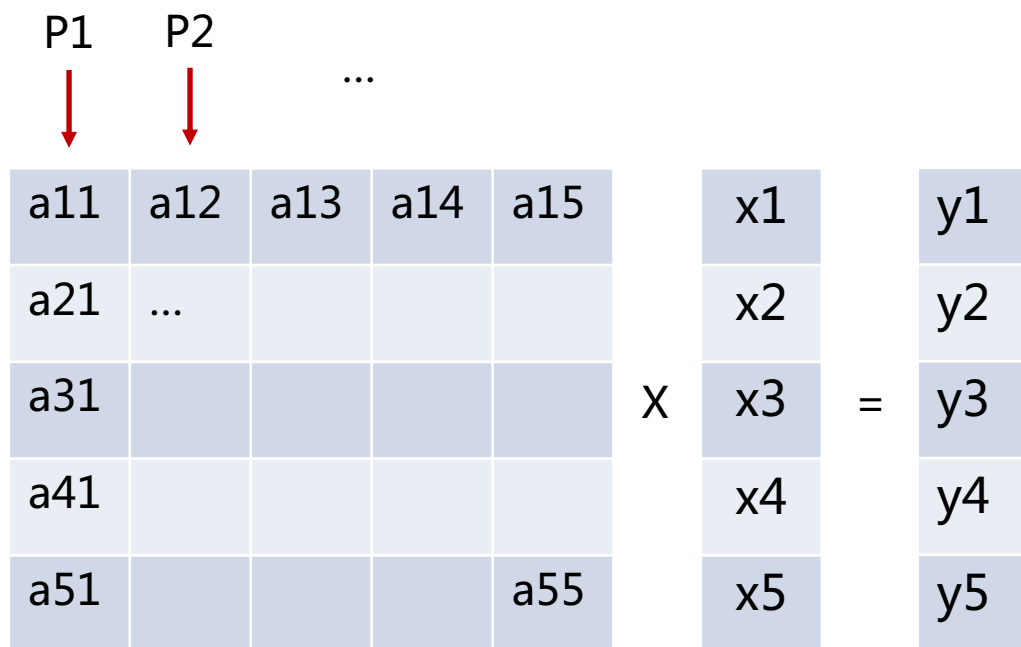
```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {  
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \  
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>()) \  
    initializer(omp_priv = omp_orig)  
  
#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)  
    for (int i=0; i<ROWS; i++) {  
        for (int j=0; j<COLS; j++) {  
            y[i] += A[i * COLS + j]*x[j];  
        }  
    }  
}
```



Maticové operace

Násobení matice vektorem

- Co když budeme násobit po sloupcích?





$$z_{ij} = a_{ij} \cdot x_j$$


$$y_i = \sum_{\{j=1, \dots, 5\}} z_{ij}$$

Maticové operace

Násobení matice vektorem

$$\begin{array}{c} a_{11} \\ a_{21} \\ a_{31} \\ a_{41} \\ a_{51} \end{array} \cdot \begin{array}{c} x_1 \\ x_1 \\ x_1 \\ x_1 \\ x_1 \end{array} = \begin{array}{c} a_{11} \cdot x_1 \\ a_{21} \cdot x_1 \\ a_{31} \cdot x_1 \\ a_{41} \cdot x_1 \\ a_{51} \cdot x_1 \end{array} \quad z_1$$


$$\begin{array}{c} a_{12} \\ a_{22} \\ a_{32} \\ a_{42} \\ a_{52} \end{array} \cdot \begin{array}{c} x_2 \\ x_2 \\ x_2 \\ x_2 \\ x_2 \end{array} = \begin{array}{c} a_{12} \cdot x_2 \\ a_{22} \cdot x_2 \\ a_{32} \cdot x_2 \\ a_{42} \cdot x_2 \\ a_{52} \cdot x_2 \end{array} \quad z_2$$


$$\begin{array}{c} a_{11} \cdot x_1 \\ a_{21} \cdot x_1 \\ a_{31} \cdot x_1 \\ a_{41} \cdot x_1 \\ a_{51} \cdot x_1 \end{array} + \begin{array}{c} a_{12} \cdot x_2 \\ a_{22} \cdot x_2 \\ a_{32} \cdot x_2 \\ a_{42} \cdot x_2 \\ a_{52} \cdot x_2 \end{array} + \dots + = \begin{array}{c} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{array}$$


Maticové operace

Násobení matice vektorem

- Bude to fungovat?

Maticové operace

Násobení matice vektorem

- Bude to fungovat?
 - Operace sice mohou být vektorizovány, nicméně přístup není vhodný pro cache (mnoho dotazů)
 - Můžeme data v matici uspořádat po sloupcích

a11	a12	a13	a14	a15
a21	...			
a31				
a41				
a51				a55



a11	a21	a31	a41	a51	...		
-----	-----	-----	-----	-----	-----	--	--

Maticové operace

Násobení matice vektorem

- Bude to teď fungovat?

```
...
// data has to be ordered by columns in memory
for (int i = 0; i < COLS; i++) {
    x[i] = rand() % 1000;
    for (int j = 0; j < ROWS; j++) {
        A[i * ROWS + j] = rand() % 1000;
    }
}
...

void multiply_column(std::vector<int> &A, std::vector<int> &x, std::vector<int> &y) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)

#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)
    for (int i = 0; i < COLS; i++) {
        for (int j = 0; j < ROWS; j++) {
            y[j] += A[i * ROWS + j]*x[i];
        }
    }
}
```

Maticové operace

Násobení matice vektorem

- Lze dále zefektivnit původní přístup?
 - Vzpomeňte si na false sharing ...

```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)

#pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)
    for (int i=0; i<ROWS; i++) {
        for (int j=0; j<COLS; j++) {
            y[i] += A[i * COLS + j]*x[j];
        }
    }
}
```

Maticové operace

Násobení matice vektorem

- Lze dále zefektivnit původní přístup?
 - Vzpomeňte si na false sharing ...

```
void multiply(std::vector<int>& A, std::vector<int>& x, std::vector<int>& y) {  
    #pragma omp declare reduction(vec_int_plus : std::vector<int> : \  
        std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \  
        initializer(omp_priv = omp_orig)  
  
    #pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)  
        for (int i=0; i<ROWS; i++) {  
            for (int j=0; j<COLS; j++) {  
                y[i] += A[i * COLS + j]*x[j];  
            }  
        }  
}
```

- Nahradíme pole lokální proměnnou

Maticové operace

Násobení matice vektorem

- Lokální proměnná

```
void multiply(std::vector<int> &A, std::vector<int> &x, std::vector<int> &y) {  
    #pragma omp declare reduction(vec_int_plus : std::vector<int> : \  
        std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \  
        initializer(omp_priv = omp_orig)  
  
    int tmp;  
    #pragma omp parallel for num_threads(thread_count) reduction(vec_int_plus : y)  
    for (int i=0; i<ROWS; i++) {  
        tmp = 0;  
        for (int j=0; j<COLS; j++) {  
            tmp += A[i * COLS + j]*x[j];  
        }  
        y[i] += tmp;  
    }  
}
```


Maticové operace

Násobení matice vektorem

Computer Languages, Systems & Structures 51 (2018) 158–175



ELSEVIER

Contents lists available at [ScienceDirect](#)

Computer Languages, Systems & Structures

journal homepage: www.elsevier.com/locate/cl

Effective Implementation of Matrix–Vector Multiplication on Intel’s AVX multicore Processor

Somaia A. Hassan^{a,*}, Mountasser M.M. Mahmoud^a, A.M. Hemeida^b,
Mahmoud A. Saber^a

Maticové operace

Násobení dvou matic

a11	a12	a13		b11	b12	b13		c11	c12	c13
a21	...		x	b21	...		=	c21	...	
...					

$$c_{ij} = \sum_{\{k=1, \dots, n\}} a_{ik} \cdot b_{kj}$$

Maticové operace

Násobení dvou matic

a11	a12	a13		b11	b12	b13		c11	c12	c13
a21	...		x	b21	...		=	c21	...	
...					

Výpočet prvků c je opět nezávislý a lze paralelizovat

$$c_{ij} = \sum_{\{k=1, \dots, n\}} a_{ik} \cdot b_{kj}$$

Nevýhody?

Velké množství úloh, malé úlohy

Maticové operace

Násobení dvou matic

a11	a12	a13		b11	b12	b13		c11	c12	c13
a21	...		×	b21	...		=	c21	...	
...					

Můžeme zvětšit úkoly spojením několika řádků

```
void multiply(std::vector<int>& A, std::vector<int>& B, std::vector<int>& C) {
#pragma omp declare reduction(vec_int_plus : std::vector<int> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \
    initializer(omp_priv = omp_orig)
    int tmp;
#pragma omp parallel for collapse(2) num_threads(thread_count) reduction(vec_int_plus : C)
    for (int i=0; i<ROWS; i++) {
        for (int j=0; j<COLS; j++) {
            tmp = 0;
            for (int k=0; k<ROWS; k++) {
                tmp += A[i * COLS + k] * B[k * COLS + j];
            }
            C[i * COLS + j] += tmp;
        }
    }
}
```

Maticové operace

Násobení dvou matic

- Rozdělení na bloky
- 1 úkol odpovídá
částečnému výsledku
submatice

Maticové operace

Násobení dvou matic

- Rozdělení na bloky
- 1 úkol odpovídá částečnému výsledku submatice (např. $c_{11}, c_{12}, c_{21}, c_{22}$)

b11	b12	b13	b14
b21	...		
...			

a11	a12	a13	a14
a21	...		
...			

c11	c12	c13	c14
c21	...		
...			

$$c_{11} += a_{11} \cdot b_{11} + a_{12} \cdot b_{21}$$

$$c_{12} += a_{11} \cdot b_{12} + a_{12} \cdot b_{22}$$

...

Maticové operace

Násobení dvou matic

```
void multiply_blocks(std::vector<int>& A, std::vector<int>& B, std::vector<int>& C) {  
    #pragma omp declare reduction(vec_int_plus : std::vector<int> : \  
        std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \  
        initializer(omp_priv = omp_orig)  
  
    const int ROWS_IN_BLOCK = 10;  
    const int BLOCKS_IN_ROW = ROWS/ROWS_IN_BLOCK;  
    int tmp;  
  
    #pragma omp parallel for collapse(2) num_threads(thread_count) reduction(vec_int_plus : C) private(tmp)  
    for (int br1=0; br1<BLOCKS_IN_ROW; br1++) {  
        for (int bb=0; bb<BLOCKS_IN_ROW; bb++) {  
            for (int bc2 = 0; bc2 < BLOCKS_IN_ROW; bc2++) {  
                for (int r = br1 * ROWS_IN_BLOCK; r < (br1 + 1) * ROWS_IN_BLOCK; r++) {  
                    for (int c = bc2 * ROWS_IN_BLOCK; c < (bc2 + 1) * ROWS_IN_BLOCK; c++) {  
                        tmp = 0;  
                        for (int k = 0; k < ROWS_IN_BLOCK; k++) {  
                            tmp += A[r * COLS + (k + bb*ROWS_IN_BLOCK)] * B[(bb*ROWS_IN_BLOCK + k) * COLS + c];  
                        }  
                        C[r * COLS + c] += tmp;  
                    }  
                }  
            }  
        }  
    }  
}
```

Maticové operace

Násobení dvou matic

$$\begin{array}{|c|c|c|} \hline a_{11} & a_{12} & a_{13} \\ \hline a_{21} & \dots & \\ \hline \dots & & \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline b_{11} & b_{12} & b_{13} \\ \hline b_{21} & \dots & \\ \hline \dots & & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline c_{11} & c_{12} & c_{13} \\ \hline c_{21} & \dots & \\ \hline \dots & & \\ \hline \end{array}$$

- A co dál? Lze využít vektorizaci?
- Můžeme najednou spočítat vektor částečných hodnot?

Maticové operace

Násobení dvou matic

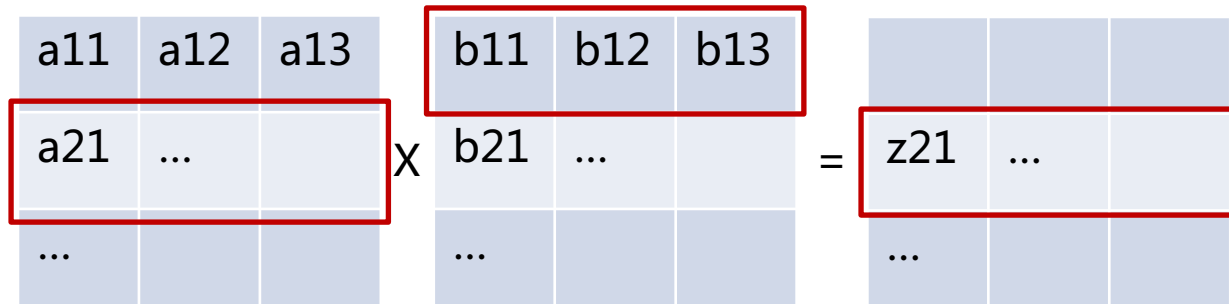
a11	a12	a13		b11	b12	b13		c11	c12	c13
a21	...		×	b21	...		=	c21	...	
...					

- A co dál? Lze využít vektorizaci?
- Můžeme najednou spočítat vektor částečných hodnot?
- Co když budeme násobit řádek i matice A a řádek j matice B?

a11	a12	a13		b11	b12	b13				
a21	...		×	b21	...		=	z21	...	
...					

Maticové operace

Násobení dvou matic



- Opět máme vektor dílčích výsledků z
- Násobení je vektorové, součet různých vektorů z lze také vektorizovat


```
void multiply(std::vector<int>& A, std::vector<int>& B, std::vector<int>& C) {  
    #pragma omp declare reduction(vec_int_plus : std::vector<int> : \  
        std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<int>())) \  
        initializer(omp_priv = omp_orig)  
  
    #pragma omp parallel for collapse(2) num_threads(thread_count) reduction(vec_int_plus : C)  
    for (int r1=0; r1<ROWS; r1++) {  
        for (int r2=0; r2<ROWS; r2++) {  
            for (int k=0; k<ROWS; k++) {  
                C[r1 * COLS + k] += A[r1 * COLS + r2] * B[r2 * COLS + k];  
            }  
        }  
    }  
}
```

Maticové operace

Gaussova eliminace

- Dalším typickým úkolem je řešení soustavy lineárních rovnic
- Lze využít Gaussovu eliminaci

a11	a12	a13	b1
a21	...		b2
...			b3



a11	a12	a13	b1
0	a'22	a'23	b'2
0	0	a'33	b'3

- Jak můžeme paralelizovat?

Maticové operace

Gaussova eliminace

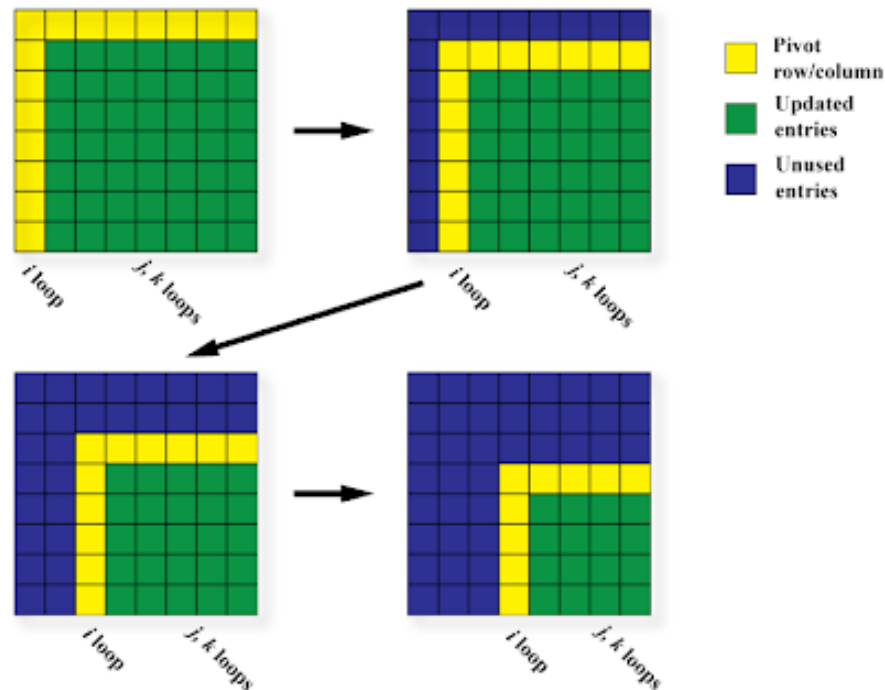
- Jaké jsou závislosti mezi hodnotami?
 - Po výběru pivotu se změny všechny hodnoty pro řádky a sloupce větší než pozice pivotu

Maticové operace

Gaussova eliminace

- Jaké jsou závislosti mezi hodnotami?
 - Po výběru pivotu se změny všech hodnot pro řádky a sloupce větší než pozice pivotu
- Kterou část můžeme paralelizovat?

```
void gauss(std::vector<double>& A) {  
    for (int i=0; i<ROWS; i++) {  
        // Make all rows below this one 0 in current column  
        for (int k=i+1; k<ROWS; k++) {  
            double c = -A[k * COLS + i]/A[i*COLS + i];  
            for (int j=i; j<ROWS; j++) {  
                if (i==j) {  
                    A[k * COLS + j] = 0;  
                } else {  
                    A[k * COLS + j] += c * A[i * COLS + j];  
                }  
            }  
        }  
    }  
}
```



Maticové operace

Gaussova eliminace

- Jaké jsou závislosti mezi hodnotami?
 - Po výběru pivotu se změny všech hodnot pro řádky a sloupce větší než pozice pivotu

```
void gauss_par(std::vector<double>& A) {
#pragma omp declare reduction(vec_int_plus : std::vector<double> : \
    std::transform(omp_out.begin(), omp_out.end(), omp_in.begin(), omp_out.begin(), std::plus<double>())) \
    initializer(omp_priv = omp_orig)

    for (int i=0; i<ROWS; i++) {
        // Make all rows below this one 0 in current column
#pragma omp parallel for num_threads(thread_count)
        for (int k=i+1; k<ROWS; k++) {
            double c = -A[k * COLS + i]/A[i*COLS + i];
            for (int j=i; j<ROWS; j++) {
                if (i==j) {
                    A[k * COLS + j] = 0;
                } else {
                    A[k * COLS + j] += c * A[i * COLS + j];
                }
            }
        }
    }
}
```

Maticové operace

Gaussova eliminace

- Co lze dále zefektivnit?

Maticové operace

Gaussova eliminace

- Co lze dále zefektivnit?

(0,0)	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(a) Iteration $k = 0$ starts

1	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(b)

1	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(c)

1	(0,1)	(0,2)	(0,3)	(0,4)
(1,0)	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(d)

1	(0,1)	(0,2)	(0,3)	(0,4)
0	(1,1)	(1,2)	(1,3)	(1,4)
(2,0)	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(e) Iteration $k = 1$ starts

1	(0,1)	(0,2)	(0,3)	(0,4)
0	1	(1,2)	(1,3)	(1,4)
0	(2,1)	(2,2)	(2,3)	(2,4)
(3,0)	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(f)

1	(0,1)	(0,2)	(0,3)	(0,4)
0	(1,1)	(1,2)	(1,3)	(1,4)
0	(2,1)	(2,2)	(2,3)	(2,4)
0	(3,1)	(3,2)	(3,3)	(3,4)
(4,0)	(4,1)	(4,2)	(4,3)	(4,4)

(g) Iteration $k = 0$ ends

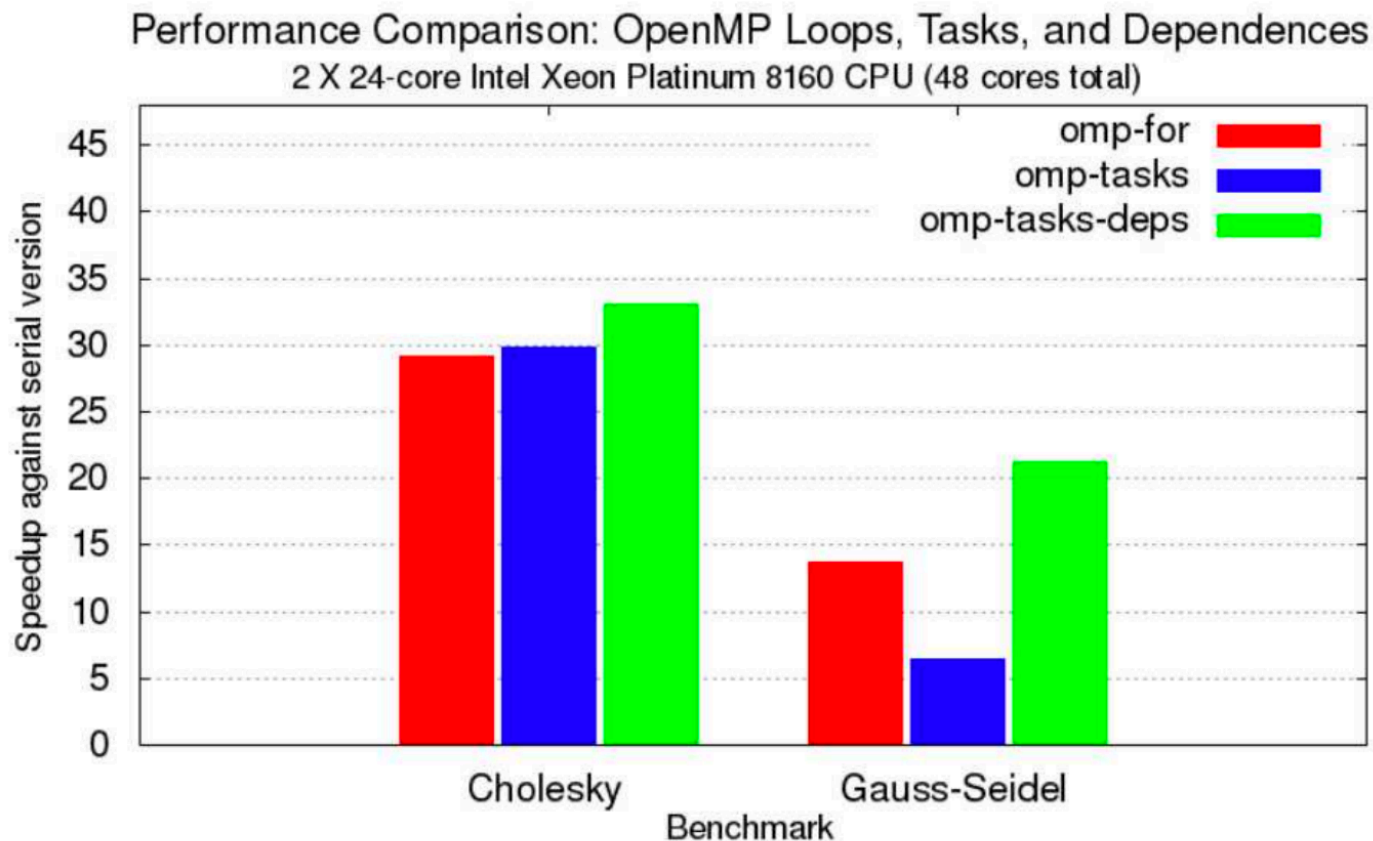
1	(0,1)	(0,2)	(0,3)	(0,4)
0	1	(1,2)	(1,3)	(1,4)
0	(2,1)	(2,2)	(2,3)	(2,4)
0	(3,1)	(3,2)	(3,3)	(3,4)
0	(4,1)	(4,2)	(4,3)	(4,4)

(h)

OpenMP: Novinky za posledních 5 let

Ze třetí přednášky

- OpenMP je specifikace, která se vyvíjí (podobně jako C++):



"The Ongoing Evolution of OpenMP"

<https://www.osti.gov/pages/servlets/purl/1465188>

OpenMP: Novinky za posledních 5 let

Offloading na GPGPU:

```
#pragma omp target data map (to: pA[0:N*N],pB[0:N*N]) map
(tofrom: pC[0:N*N])
#pragma omp target
#pragma omp teams distribute parallel for collapse(2)
private(i,j,k)
for (i=0;i<N;i++) {
    for (j=0;j<N;j++) {
        for (k=0;k<N;k++) {
            pC(i,j)+=pA(i,k)*pB(k,j);
        }
    }
}
```

Shrnutí paralelní části

- Základní nástroje pro psaní paralelního programu
 - vlákna a práce s nimi
 - Atomické proměnné
 - OpenMP
 - (Vektorizace, SIMD paralelizace)

Shrnutí paralelní části

- Základní nástroje pro psaní paralelního programu
 - vlákna a práce s nimi
 - atomické proměnné
 - OpenMP
 - Vektorizace (SIMD paralelizace)
- Základní techniky paralelizace
 - Rozděluj a panuj
 - Threadpool
 - Dekompozice, nalezení co možná nejvíc paralelně bezkonfliktně vykonatelných operací

Shrnutí paralelní části

- Základní nástroje pro psaní paralelního programu
 - vlákna a práce s nimi
 - Atomické proměnné
 - OpenMP
 - (Vektorizace, SIMD paralelizace)
- Základní techniky paralelizace
 - Rozděluj a panuj
 - Threadpool
 - Dekompozice, nalezení co možná nejvíc paralelně bezkonfliktně vykonatelných operací
- Základní algoritmy
 - Řazení
 - Maticové operace

Shrnutí paralelní části

- Paralelizujete s cílem zefektivnit běh programu/algoritmu
- Musíte (alespoň částečně) rozumět vykonávání programu na HW
 - False sharing
 - Cache optimization

Shrnutí paralelní části

- Paralelizujete s cílem zefektivnit běh programu/algoritmu
- Musíte (alespoň částečně) rozumět vykonávání programu na HW
 - False sharing
 - Cache optimization
- Vynechali jsme spoustu věcí
 - Úvodní kurz, aby jste získali základní znalosti a zkušenosti
- Pokud Vás paralelní programování zaujalo
 - Paralelní algoritmy (B4M35PAG)
 - Obecné výpočty na grafických procesorech (B4M39GPU)

Shrnutí paralelní části

- Pro implementační zkoušku – programujte, programujte, programujte!

Shrnutí paralelní části

- Pro implementační zkoušku – programujte, programujte, programujte!
 - Dostanete problém + sériový algoritmus
 - Cílem bude **zrychlit** algoritmus **díky paralelizaci**
 - **Paralelizace musí být korektní!** (musíte přemýšlet, ne všechny chyby se projeví, paralelní programování je nedeterministické)

Shrnutí paralelní části

- Pro implementační zkoušku – programujte, programujte, programujte!
 - Dostanete problém + sériový algoritmus
 - Cílem bude **zrychlit** algoritmus **díky paralelizaci**
 - **Paralelizace musí být korektní!** (musíte přemýšlet, ne všechny chyby se projeví, paralelní programování je nedeterministické)
- V dalším studiu/práci – paralelizujte, pokud je to potřeba!
 - Pracujte iterativně – nejdřív je potřeba mít korektní sériovou variantu
 - Pokud je pomalá – zrychlujeme, paralelizujeme, atd.

Shrnutí paralelní části

