

Sem vložte zadání Vaší práce.

ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE
FAKULTA INFORMAČNÍCH TECHNOLOGIÍ
KATEDRA TEORETICKÉ INFORMATIKY



Bakalářská práce

Vliv formátu uložení řídské matice na výkonnost násobení řídkých matic

Tomáš Nesrovnal

Vedoucí práce: Ing. Ivan Šimeček, Ph.D

12. května 2014

Poděkování

Děkuji vedoucímu práce Ing. Ivanu Šimečkovi, Ph.D. za cenné rady. Bc. Štefanu Šafárovi za korektury a poznámky.

Prohlášení

Prohlašuji, že jsem předloženou práci vypracoval(a) samostatně a že jsem uvedl(a) veškeré použité informační zdroje v souladu s Metodickým pokynem o etické přípravě vysokoškolských závěrečných prací.

Beru na vědomí, že se na moji práci vztahují práva a povinnosti vyplývající ze zákona č. 121/2000 Sb., autorského zákona, ve znění pozdějších předpisů. V souladu s ust. § 46 odst. 6 tohoto zákona tímto uděluji nevýhradní oprávnění (licenci) k užití této mojí práce, a to včetně všech počítačových programů, jež jsou její součástí či přílohou a veškeré jejich dokumentace (dále souhrnně jen „Dílo“), a to všem osobám, které si přejí Dílo užít. Tyto osoby jsou oprávněny Dílo užít jakýmkoli způsobem, který nesnižuje hodnotu Díla a za jakýmkoli účelem (včetně užití k výdělečným účelům). Toto oprávnění je časově, teritoriálně i množstevně neomezené. Každá osoba, která využije výše uvedenou licenci, se však zavazuje udělit ke každému dílu, které vznikne (byť jen zčásti) na základě Díla, úpravou Díla, spojením Díla s jiným dílem, zařazením Díla do díla souborného či zpracováním Díla (včetně překladu), licenci alespoň ve výše uvedeném rozsahu a zároveň zpřístupnit zdrojový kód takového díla alespoň srovnatelným způsobem a ve srovnatelném rozsahu, jako je zpřístupněn zdrojový kód Díla.

V Praze dne 12. května 2014

.....

České vysoké učení technické v Praze
Fakulta informačních technologií

© 2014 Tomáš Nesrovnal. Všechna práva vyhrazena.

Tato práce vznikla jako školní dílo na Českém vysokém učení technickém v Praze, Fakultě informačních technologií. Práce je chráněna právními předpisy a mezinárodními úmluvami o právu autorském a právech souvisejících s právem autorským. K jejímu užití, s výjimkou bezúplatných zákonných licencí, je nezbytný souhlas autora.

Odkaz na tuto práci

Nesrovnal, Tomáš. *Vliv formátu uložení řádké matice na výkonnost násobení řádkých matic*. Bakalářská práce. Praha: České vysoké učení technické v Praze, Fakulta informačních technologií, 2014.

Abstrakt

Tato práce popisuje formáty pro uložení řídkých matic COO, CSR, BSR, Quadtree a jeho modifikace v podobě snížení výšky stromu. Dále popisuje algoritmy pro násobení matic v těchto formátech. Součástí práce je i implementace těchto formátů a algoritmů v jazyce C a experimentální porovnání výkonnosti s teoretickými předpoklady.

Klíčová slova řídká matice, kvadrantový strom, násobení řídkých matic

Abstract

We describe sparse matrix storage formats COO, CSR, BSR, Quadtree and its modification in form of a lowering its height in this thesis. We also describe algorithms for matrix multiplication in these formats. An implementation of these formats and algorithms in the C language with an experimental comparison of measured results with theoretical assumptions is a part of this thesis too.

Keywords sparse matrix, quadtree, sparse matrix multiplication

Obsah

Úvod	1
1 Úvod do problematiky	3
1.1 Řídké matice	3
1.2 Použití násobení matic	3
1.3 Matice	3
1.4 Vektor	5
1.5 Definice násobení matice maticí	5
1.6 Složitosti	5
1.7 Řídké matice	6
2 Algoritmy násobení matic	7
2.1 Podle definice	7
2.2 Násobení transponovanou maticí	8
2.3 Násobení po řádcích	8
2.4 Rekurzivní násobení	9
2.5 Strassenův algoritmus	10
2.6 Rychlé algoritmy	14
2.7 Algoritmus podle definice upravený pro řídké matice	14
2.8 Rychlé násobení řídkých matic	15
2.9 Další algoritmy pro řídké matice	16
3 Formáty uložení řídkých matic	19
3.1 COO - Coordinate list	19
3.2 CSR - Compressed sparse row	21
3.3 BSR - Block Sparse Row	23
3.4 Quadtree	26
3.5 Modifikace formátu quadtree	27
4 Analýza a návrh	31

4.1	Software pro práci s řídkými maticemi	31
4.2	Řešení implementace	32
5	Realizace	33
5.1	Prostředí	33
5.2	Design implementace	33
5.3	Implementace KAT	33
5.4	Testování	35
5.5	Měření	36
5.6	Porovnání výkonnosti s předpoklady	41
5.7	Zhodnocení výsledků	42
Závěr		47
	Další práce	47
Literatura		49
A Přílohy		53
A.1	Překlad programu	53
A.2	Překlad textu práce	54
A.3	Práce s programem	54
A.4	Práce s formátem MatrixMarket	54
A.5	Generátor řídkých matic	54
A.6	Spuštění testů	57
B Seznam použitých zkratk		59
Seznam obrázků		59
C Obsah přiloženého CD		63

Seznam obrázků

1.1	3D neorientovaný graf ve tvaru helikoptéry, jeho reprezentace v řídké matici a výsledek simulace	4
2.1	Strassenův algoritmus [12]	11
2.2	Ukázka numerické stability Strassenova algoritmu	14
2.3	Matice orsirr_1 před vynásobením	16
2.4	Matice orsirr_1 po vynásobením sama se sebou	17
3.1	Matice uložená ve formátu CSR	22
3.2	Matice uložená ve formátu BSR	24
3.3	Rozdělení matice	26
3.4	Rozdělení matice	27
3.5	Rozdělení matice	27
5.1	UML diagram tříd programu	34
5.2	Zrychlení výpočtu matice · matice oproti formátu CSR	38
5.3	Zrychlení výpočtu matice · vektor oproti formátu CSR	39
5.4	Velikost datových struktur	40
5.5	Počet uzlů v matici KAT pro uložení matice EX6	40
5.6	Počet uzlů v matici KAT pro uložení matice exdata-1	41
5.7	Počet uzlů v matici KAT pro uložení matice fp	42
5.8	Počet uzlů v matici KAT pro uložení matice gupta3	43
5.9	Počet uzlů v matici KAT pro uložení matice heart1	44
5.10	Počet uzlů v matici KAT pro uložení matice human-gene2	44
5.11	Čas násobení husté matice s hustou maticí v řídkém formátu	45
5.12	Čas násobení husté matice s vektorem v řídkém formátu	45
5.13	Datová velikost hustých matic v řídkých formátech	46
A.1	Matice vygenerovaná generátorem	56

Seznam tabulek

3.1	Matice uložená ve formátu COO	20
5.1	Přehled časových složitostí násobení matic	42
5.2	Přehled paměťových složitostí násobení matic	42

Úvod

Při řešení problémů často hledáme způsob, jak interpretovat data v takovém formátu, se kterým již umíme pracovat. Jedním z takových základních prvků jsou matice. V některých případech matice obsahují nemalý počet nulových prvků. Takové matice obecně nazýváme řídké a při práci s nimi této vlastnosti využíváme.

Binární operace násobení nad množinou matic patří mezi základní operace z lineární algebry. Pomocí násobení matic lze skládat lineární transformace. Tato operace najde uplatnění v mnohých vědeckých disciplínách.

Formát uložení řídké matice má vliv na výkonnost násobení řídkých matic. V této práci popíšeme formáty COO, CSR, BSR, Quadtree a navrhne modifikaci formátu Quadtree snížením výšky stromu. Implementujeme tyto formáty spolu s algoritmy pro násobení matic v těchto formátech a porovnáme teoretické předpoklady s naměřenými hodnotami.

V první kapitole zmíníme některé příklady použití násobení řídkých matic v praxi a definujeme některé pojmy. Ve druhé kapitole představíme obecné algoritmy pro násobení matic a ve třetí kapitole popíšeme řídké formáty COO, CSR, BSR, Quadtree a jeho modifikaci a algoritmy pro násobení matic v těchto formátech. Ve čtvrté kapitole ukážeme některé současné programy pro práci s řídkými maticemi. V páté kapitole popíšeme naši implementaci a ukážeme naměřené hodnoty.

Úvod do problematiky

1.1 Řídké matice

Pro výstavu v roce 2011 s názvem Art in Enginnering v Harnově Muzeu představila floridská univerzita kolekci ilustrací řídkých matic s názvem The Beauty of Mathematics: As Illustrated by the University of Florida Sparse Matrix Collection [21]. Pro estetické zobrazení řídkých matic byla provedena simulace [20]. Každému uzlu byl přiřazen elektrický náboj a každá hrana představovala pružinu. V simulaci byla celá tato konstrukce položena na tvrdou podložku. Simulace byla zastavena v okamžiku, kdy se konstrukce přestala hýbat. Pro ilustraci přikládáme výsledek simulace 1.1 na modelu helikoptéry, tedy neorientovaném 3D grafu uloženém v řídké matici.

Řídké matice jsou používané ve velké škále oblastí [13], například výpočty diferenciálních rovnic, Google page rank, 3D grafika, statistika, těžení z dat, kryptografie a další.

1.2 Použití násobení matic

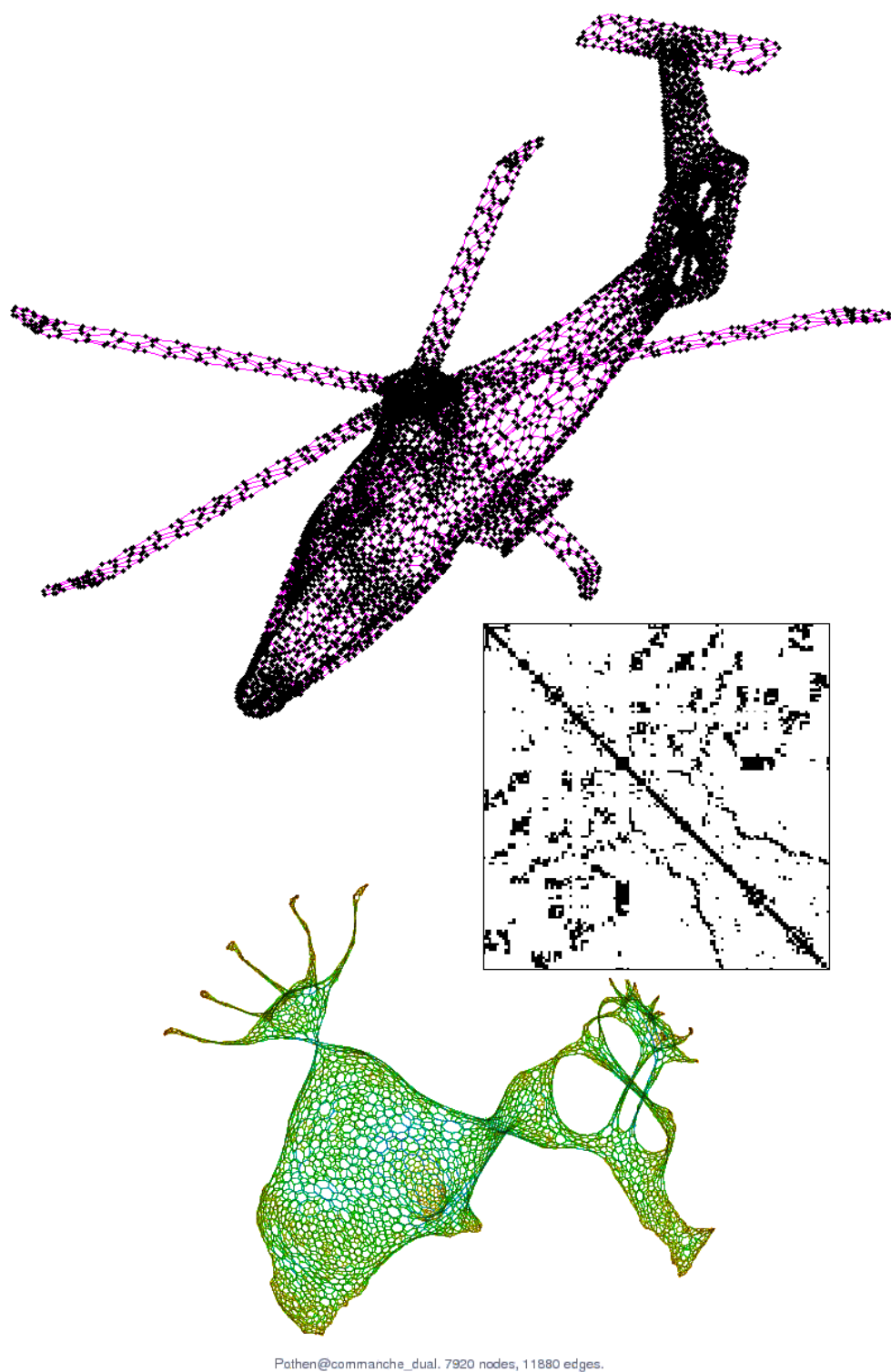
Násobení maticí může znázorňovat nějakou transformaci jednoho stavu do druhého. Největší uplatnění řídkých matic je tedy v simulaci jevů, například fyzikálních, biologických, chemických, ekonomických a dalších.

Jako příklad simulace můžeme uvést metodu konečných prvků [23][8].

Pokud potřebujeme vynásobit jednu matici s více vektory, můžeme vektory složit do matice a provést násobení matice s maticí. To se hodí například v real-time aplikacích.

1.3 Matice

Matice \mathbf{A} typu (m, n) je mn uspořádaných prvků z množiny \mathbf{R} . O prvku $a_{r,s} \in \mathbf{R}, r \in \{1, 2, \dots, m\}, s \in \{1, 2, \dots, n\}$ říkáme, že je na r -tém řádku



Obrázek 1.1: 3D neorientovaný graf ve tvaru helikoptéry, jeho reprezentace v řídké matici a výsledek simulace

a s-tém sloupci matice \mathbf{A} . Matici \mathbf{A} zapisujeme do řádků a sloupců takto:

$$\mathbf{A} = \begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & a_{m,2} & \cdots & a_{m,n} \end{pmatrix} \quad (1.1)$$

Matici \mathbf{M} typu (m, n) , kde všechny její prvky jsou rovny nule, nazýváme *nulovou maticí*.

O matici typu (m, n) budeme říkat, že je m široká a n vysoká. Pokud o matici řekneme, že má velikost n , myslíme tím, že je typu (n, n) .

1.4 Vektor

Matici \mathbf{V} typu $(1, n)$ nazveme vektorem.

1.5 Definice násobení matice maticí

Buď \mathbf{A} matice typu (m, n) s prvky $a_{i,j}$ a \mathbf{B} matice typu (n, p) s prvky $b_{j,k}$. Definujeme součin matic $\mathbf{A} \cdot \mathbf{B}$ jako matici \mathbf{C} typu (m, p) s prvky $c_{i,k}$, které vypočteme jako:

$$c_{row,col} = \sum_{k=1}^N a_{row,k} \cdot b_{k,col} \quad (1.2)$$

Výsledek součinu matic se nezmění, pokud matice doplníme o libovolný počet nulových řádků a nebo sloupců. Této vlastnosti můžeme využít pro získání potřebných rozměrů:

1. Při násobení matice \mathbf{A} typu (m, n) s maticí \mathbf{B} typu (o, p) , kde $n \neq o$.
2. Pokud potřebujeme matice stejné velikosti.
3. Pokud potřebujeme matice určité velikosti, například $2^{\mathbb{N}}$.

Násobení matice vektorem je pouze případem násobení matice typu (m, n) s maticí typu $(1, n)$.

1.6 Složitosti

K označení složitostí, ať již časových nebo prostorových (paměťových), používáme \mathcal{O} -notaci, která označuje množinu funkcí, asymptoticky rostoucí řádově stejně rychle nebo rychleji.

$$\mathcal{O}(g(n)) = \left\{ f(n) : \exists c \in \mathbb{R}^+ \exists n_0 \in \mathbb{N} \forall n \geq n_0 : 0 \leq f(n) \leq cg(n) \right\} \quad (1.3)$$

\mathcal{O} notací budeme popisovat nejhorší možný případ. Obdobně Θ značí funkce rostoucí stejně rychle a Ω stejně rychle nebo pomaleji.

Pro výpočet složitostí rekurzivních algoritmů používáme mistrovskou metodu [11]. Pokud $a \geq 1, b > 0$ jsou konstanty a $f(n)$ je funkce o jedné proměnné, tak rekurence

$$t(n) = at(n/b) + f(n) \tag{1.4}$$

má asymptotickou složitost:

1. Pro $f(n) = \mathcal{O}(n^{\log_b a - \epsilon})$, kde $\epsilon > 0$ je konstanta platí, že $t(n) = \Theta(n^{\log_b a})$.
2. Pro $f(n) = \Theta(n^{\log_b a})$ platí, že $t(n) = \Theta(n^{\log_b a} \log n)$.
3. Pro $f(n) = \Omega(n^{\log_b a + \epsilon})$, kde $\epsilon > 0$ je konstanta a $af(n/b) \leq cf(n)$ pro konstantu $c < 1$ a $\forall n \geq n_0$, tak platí, že $t(n) = \Theta(f(n))$.

1.7 Řídké matice

Matice, které obsahují velké množství nulových prvků, nazýváme řídké. Nebudeme přesně uvádět, kolik procent z celkového počtu prvků musí být nulových, abychom matici nazývali řídkou. Stejně jako řídkou matici můžeme uložit do formátu pro husté matice, můžeme hustou matici uložit do formátu pro řídké matice.

Řídkost matice budeme vyjadřovat pomocí *nnz* (Number of NonZero elements), tedy počtem nenulových prvků z celkových mn , pro matici A typu (m, n) .

Algoritmy násobení matic

V této kapitole představíme některé základní a pokročilejší algoritmy pro násobení matic. Základní algoritmy mají stejnou asymptotickou složitost $\mathcal{O}(n^3)$ a liší se pouze přístupem k prvkům, což je důležité pro řídké formáty. Pokročilejší algoritmy jsou sice asymptoticky rychlejší, ale přinášejí nevýhodu ve formě numerické stability a velké skryté konstanty.

2.0.1 Pseudokódy

Pseudokódy v této práci jsou ve stylu syntaxe jazyka Fortran, ale budou představovat zápis jazyka C99. Pro pole platí, že jsou indexovaná od nuly a for cyklus `for i ← 0 to 10` bude iterovat od nultého do devátého prvku včetně.

2.1 Podle definice

Základním algoritmem násobení dvou matic je podle definice 1.5. Ve třech for cyklech postupně vybíráme řádky matice A, sloupce matice B a v N krocích násobíme. N je jak šířka matice A, tak i výška matice B.

Algorithm 1 Násobení matic podle definice

```
1: procedure MMM-DEFINITION( $A, B, C$ )                                ▷ A,B,C jsou matice
2:   for  $row \leftarrow 0$  to  $A.height$  do                                ▷ řádky
3:     for  $col \leftarrow 0$  to  $B.width$  do                                ▷ sloupce
4:        $sum \leftarrow 0$ ;
5:       for  $i \leftarrow 0$  to  $A.height$  do
6:          $sum \mathrel{+}= A[row][i] * B[i][col]$ ;
7:       end for
8:        $C[row][col] \leftarrow sum$ ;
9:     end for
10:  end for
11: end procedure
```

Z pseudokódu je vidět, že ve dvou for cyklech provádíme N násobení a N sčítání. Asymptotická složitost je tedy $\mathcal{O}(n^2(n+n)) = \mathcal{O}(n^3)$. V ukázkových výpočtech je násobení pouze $N - 1$ krát, to proto, že neuvádíme přičítání k nule (řádek 6).

2.2 Násobení transponovanou maticí

Pokud nám formát uložení matice nedovolí procházet prvky po sloupcích, je řešením druhou matici transponovat. Poté můžeme násobit řádky matice A s řádky transponované matice B.

Algorithm 2 Násobení transponovanou maticí

```
1: procedure MMM-TRANPOSE( $A, B, C$ ) ▷ A,B,C jsou matice
2:    $B \leftarrow \text{transpose}(B)$ 
3:   for  $\text{rowA} \leftarrow 0$  to  $A.\text{height}$  do ▷ řádky
4:     for  $\text{rowB} \leftarrow 0$  to  $B.\text{height}$  do ▷ sloupce
5:        $\text{sum} \leftarrow 0$ ;
6:       for  $i \leftarrow 0$  to  $A.\text{height}$  do
7:          $\text{sum} += A[\text{rowA}][i] * B[\text{rowB}][i]$ ;
8:       end for
9:        $C[\text{rowA}][\text{rowB}] \leftarrow \text{sum}$ ;
10:    end for
11:  end for
12: end procedure
```

Podobný algoritmus můžeme použít i pokud nám formát nedovolí procházet prvky po řádcích, ale pouze po sloupcích. Například v této práci neuvedený Compressed Sparse Columns.

2.3 Násobení po řádcích

Další možností, jak násobit dvě matice, kde nám formát uložení nedovolí procházet po sloupcích, je procházet současně řádky matice A i B a přičítat jednotlivé součiny na správné místo ve výsledné matici C.

Nevýhodou tohoto řešení je velký počet náhodných přístupů do pole C. Protože k prvkům přičítáme, tedy načítáme a sčítáme, je potřeba před samotným násobením nastavit všechny prvky matice C na hodnotu nula.

Algorithm 3 Násobení po řádcích

```

1: procedure MMM-BY-ROWS(A, B, C)                                ▷ A, B, C jsou matice
2:   for r ← 0 to A.height do                                     ▷ řádky matice A i B
3:     for cA ← 0 to A.width do                                   ▷ sloupce matice A
4:       for c ← 0 to B.width do                                   ▷ sloupce matice B
5:         C[r][cA] += A[r][cA] * B[r][cB];
6:       end for
7:     end for
8:   end for
9: end procedure

```

2.4 Rekurzivní násobení

Pro matice *A* i *B* o stejné velikosti 2^N můžeme použít rekurzivní přístup, tedy programovací techniku rozděl a panuj, kdy rozdělíme větší problémy na menší podproblémy.

Každou z matic rozdělíme na čtvrtiny a jednotlivé podmatice násobíme algoritmem podle definice, tedy jako matice o velikosti dva.

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \cdot \begin{pmatrix} e & f \\ g & h \end{pmatrix} = \begin{pmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{pmatrix} \quad (2.1)$$

Tento postup opakujeme, dokud velikostí podmatic nenarazíme na práh, tedy hodnotu, při které opustíme rekurzivní algoritmus a použijeme algoritmus lineární. V ukázkovém pseudokódu dělíme podmatice až na velikost prahu jedna, podmatice tedy obsahují pouze jeden prvek.

Algorithm 4 Rekurzivní násobení

```

1: procedure MMM-RECURSIVE(A, B, C, ay, ax, by, bx, cy, cx, n)
2:   if n = 1 then
3:     C[cy][cx] ← C[cy][cx] + A[ay][ax] · B[by][bx];
4:     return;
5:   end if
6:   for all r ∈ {0, n/2} do
7:     for all c ∈ {0, n/2} do
8:       for all i ∈ {0, n/2} do
9:         MMM-recursive(A, B, C, ay + i, ax + r, by + c, bx + i, cy +
10:          c, cx + r, n/2);
11:       end for
12:     end for
13:   end for
14: end procedure

```

Pro ilustraci jako příklad uvádíme výpočet horního levého prvku v násobení dvou matic o velikosti 2^2 . Pro větší přehlednost značíme prvky malým písmem z názvu matice a indexy o jejich pozicích.

$$\begin{aligned}
& \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ a_{3,1} & a_{3,2} & a_{3,3} & a_{3,4} \\ a_{4,1} & a_{4,2} & a_{4,3} & a_{4,4} \end{pmatrix} \cdot \begin{pmatrix} b_{1,1} & b_{1,2} & b_{1,3} & b_{1,4} \\ b_{2,1} & b_{2,2} & b_{2,3} & b_{2,4} \\ b_{3,1} & b_{3,2} & b_{3,3} & b_{3,4} \\ b_{4,1} & b_{4,2} & b_{4,3} & b_{4,4} \end{pmatrix} = \\
& \begin{pmatrix} \begin{pmatrix} a_{1,1} & a_{1,2} \end{pmatrix} \cdot \begin{pmatrix} b_{1,1} & b_{1,2} \end{pmatrix} + \begin{pmatrix} a_{1,3} & a_{1,4} \end{pmatrix} \cdot \begin{pmatrix} b_{3,1} & b_{3,2} \end{pmatrix} & \dots \\ a_{2,1} & a_{2,2} & \dots & \dots \end{pmatrix} = \\
& \begin{pmatrix} \begin{pmatrix} a_{1,1}b_{1,1} + a_{1,2}b_{2,1} & \dots \\ \dots & \dots \end{pmatrix} + \begin{pmatrix} a_{1,3}b_{3,1} + a_{1,4}b_{4,1} & \dots \\ \dots & \dots \end{pmatrix} & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix} = \\
& \begin{pmatrix} \begin{pmatrix} a_{1,1}b_{1,1} + a_{1,2}b_{2,1} + a_{1,3}b_{3,1} + a_{1,4}b_{4,1} & \dots \\ \dots & \dots \end{pmatrix} & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix}
\end{aligned}$$

Asymptotická složitost je samozřejmě stejná jako u algoritmu podle definice. Asymptotickou složitost rekursivního algoritmu můžeme spočítat pomocí mistrovské metody.

$$T(n) = \begin{cases} \Theta(1) & \text{if } n = 1 \\ 8T(n/2) + \Theta(1) & \text{if } n > 1 \end{cases}$$

Protože platí, že $a = 8, b = 2, r = \log_2 8, n^r = n^{\log_2 8} = n^3 = \Omega(1)$, tak asymptotická složitost podle mistrovské metody je $\text{MMM-recursive}(n) = \mathcal{O}(n^3)$.

Kvůli režii rekursivního dělení v praxi nezmenšujeme podmatice až na velikost jedna. Vhodný práh velikosti podmatice je například takový, co se vejde do L1 cache.

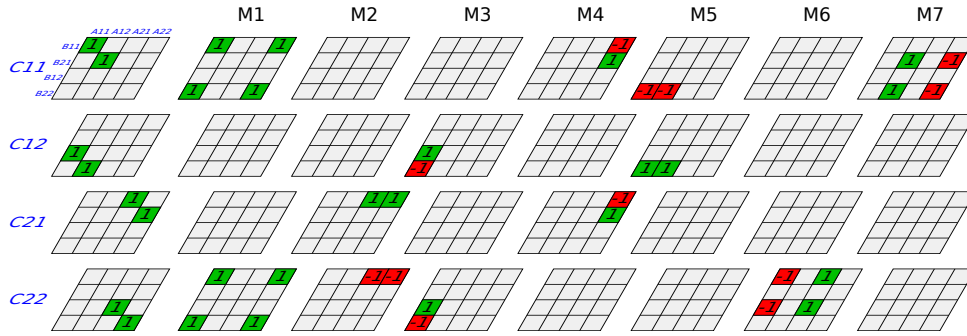
2.5 Strassenův algoritmus

V roce 1969 Volker Strassen v časopise *Numerische Mathematik* publikoval článek [35], ve kterém jako první představil algoritmus násobení dvou matic s menší asymptotickou složitostí než algoritmus podle definice, tedy $\mathcal{O}(n^3)$.

Algoritmus je založen na myšlence, že sčítání je operace méně náročnější než operace násobení. Respektive dvě matice umíme sečíst nebo odečíst se složitostí $\mathcal{O}(n^2)$, ale vynásobit se složitostí $\mathcal{O}(n^3)$.

Volker Strassen tedy využil jisté symetrie [24] v násobení dvou matic A a B o velikosti dva a výslednou matici C seskládal pomocí sedmi pomocných matic.

Obrázek 2.1 ukazuje, z čeho se pomocné matice skládají a jak jsou do výsledné matice seskládány. V ilustračních maticích o velikosti čtyři ukazujeme, které sčítance pomocná matice do výsledku přičítá a které odečítá.



Obrázek 2.1: Strassenův algoritmus [12]

Zápis Strassenova algoritmu vypadá následovně:

$$\begin{aligned}
 A \cdot B &= \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix} \cdot \begin{pmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{pmatrix} \\
 M_1 &= (A_{1,1} + A_{2,2}) \cdot (B_{1,1} + B_{2,2}) \\
 M_2 &= (A_{2,1} + A_{2,2}) \cdot B_{1,1} \\
 M_3 &= A_{1,1} \cdot (B_{1,2} - B_{2,2}) \\
 M_4 &= A_{2,2} \cdot (B_{2,1} - B_{1,1}) \\
 M_5 &= (A_{1,1} + A_{1,2}) \cdot B_{2,2} \\
 M_6 &= (A_{2,1} - A_{1,1}) \cdot (B_{1,1} + B_{1,2}) \\
 M_7 &= (A_{1,2} - A_{2,2}) \cdot (B_{2,1} + B_{2,2}) \\
 C &= \begin{pmatrix} M_1 + M_4 - M_5 + M_7 & M_3 + M_5 \\ M_2 + M_4 & M_1 - M_2 + M_3 + M_6 \end{pmatrix} \quad (2.2)
 \end{aligned}$$

V pseudokódu používáme procedury **offset-add**, respektive **offset-sub**. Slouží ke sčítání, respektive odečítání bloku prvků o velikosti n v maticích od nějakého offsetu y a x . Parametry obou funkcí jsou: **offset-*(A, B, C, ay, ax, by, bx, cy, cx, n)**.

Algorithm 5 Strassenův algoritmus

```

1: procedure MMM-STRASSEN( $A, B, C, ay, ax, by, bx, cy, cx, n$ )
2:   if  $n = 1$  then
3:      $C[cy][cx] \leftarrow C[cy][cx] + A[ay][ax] \cdot B[by][bx]$ ;
4:     return;
5:   end if
6:    $h \leftarrow n/2$ ; ▷ čtvrtina
7:    $m[9] \leftarrow \text{init-matrices}(9, h)$ ; ▷ devět pomocných matic
8:   offset-add( $a, a, m[8], ay, ax, ay + h, ax + h, 0, 0, h$ ); ▷ M1
9:   offset-add( $b, b, m[9], by, bx, by + h, bx + h, 0, 0, h$ );
10:  MMM-strassen( $m[8], m[9], m[1], 0, 0, 0, 0, 0, 0, h$ );
11:  offset-add( $a, a, m[8], ay + h, ax, ay + h, ax + h, 0, 0, h$ ); ▷ M2
12:  MMM-strassen( $m[8], b, m[2], 0, 0, bx, by, 0, 0, h$ );
13:  offset-sub( $b, b, m[8], by, bx + h, by + h, bx + h, 0, 0, h$ ); ▷ M3
14:  MMM-strassen( $a, m[8], m[3], ay, ax, 0, 0, 0, 0, h$ );
15:  offset-sub( $b, b, m[8], by + h, bx, by, bx, 0, 0, h$ ); ▷ M4
16:  MMM-strassen( $a, m[8], m[4], ay + h, ax + h, 0, 0, 0, 0, h$ );
17:  offset-add( $a, a, m[8], ay, ax, ay, ax + h, 0, 0, h$ ); ▷ M5
18:  MMM-strassen( $m[8], b, m[5], 0, 0, by + h, bx + h, 0, 0, h$ );
19:  offset-sub( $a, a, m[8], ay + h, ax, ay, ax, 0, 0, h$ ); ▷ M6
20:  offset-add( $b, b, m[9], by, bx, by, bx + h, 0, 0, h$ );
21:  MMM-strassen( $m[8], m[9], m[6], 0, 0, 0, 0, 0, 0, h$ );
22:  offset-sub( $a, a, m[8], ay, ax + h, ay + h, ax + h, 0, 0, h$ ); ▷ M7
23:  offset-add( $b, b, m[9], by + h, bx, by + h, bx + h, 0, 0, h$ );
24:  MMM-strassen( $m[8], m[9], m[7], 0, 0, 0, 0, 0, 0, h$ );
25:  offset-add( $m[1], m[4], m[8], 0, 0, 0, 0, 0, 0, h$ ); ▷ c1,1
26:  offset-sub( $m[8], m[5], m[8], 0, 0, 0, 0, 0, 0, h$ );
27:  offset-add( $m[8], m[7], c, 0, 0, 0, 0, cy, cx, h$ );
28:  offset-add( $m[3], m[5], c, 0, 0, 0, 0, cy, cx + h, h$ ); ▷ c1,2
29:  offset-add( $m[2], m[4], c, 0, 0, 0, 0, cy + h, cx, h$ ); ▷ c2,1
30:  offset-sub( $m[1], m[2], m[8], 0, 0, 0, 0, 0, 0, h$ ); ▷ c2,2
31:  offset-add( $m[8], m[3], m[8], 0, 0, 0, 0, 0, 0, h$ );
32:  offset-add( $m[8], m[6], c, 0, 0, 0, 0, cy + h, cx + h, h$ );
33: end procedure

```

Výpočet asymptotické složitosti provedeme podobně jako u **MMM-recursive**, tedy mistrovskou metodou.

V každém kroku rekurze počítáme $\Theta(n^2)$ operací na vytvoření pomocných matic.

$$T(n) = \begin{cases} \Theta(1) & \text{if } n = 1 \\ 7T(n/2) + \Theta(n^2) & \text{if } n > 1 \end{cases}$$

Protože platí, že $a = 7, b = 2, r = \log_2 7, n^r = n^{\log_2 7} = \Omega(n^2)$, tak

asymptotická složitost podle mistrovské metody je $\text{MMM-strassen}(n) = \mathcal{O}(n^{\log_2 7}) \approx \mathcal{O}(n^{2.8})$.

Stejně jako v předešlém algoritmu i zde demonstrujeme výpočet levého horního prvku matice z násobení dvou matic o velikosti dva. Místo parametrické matice použijeme konkrétní desetinná čísla, abychom ukázali numerickou stabilitu Strassenova algoritmu. Pro ukázkou budeme uvažovat počítač, který u čísel ukládá pouze pět cifer, znaménko a desetinnou čárku.

Pomocí algoritmu podle definice by takový počítač vypočítal součin dvou matic následovně:

$$\begin{pmatrix} 30.234 & 0.5678 \\ 0.9123 & 10.456 \end{pmatrix} \cdot \begin{pmatrix} 0.8912 & 0.3456 \\ 0.7891 & 9.999 \end{pmatrix} = \begin{pmatrix} 27.392 & \dots \\ \dots & \dots \end{pmatrix}$$

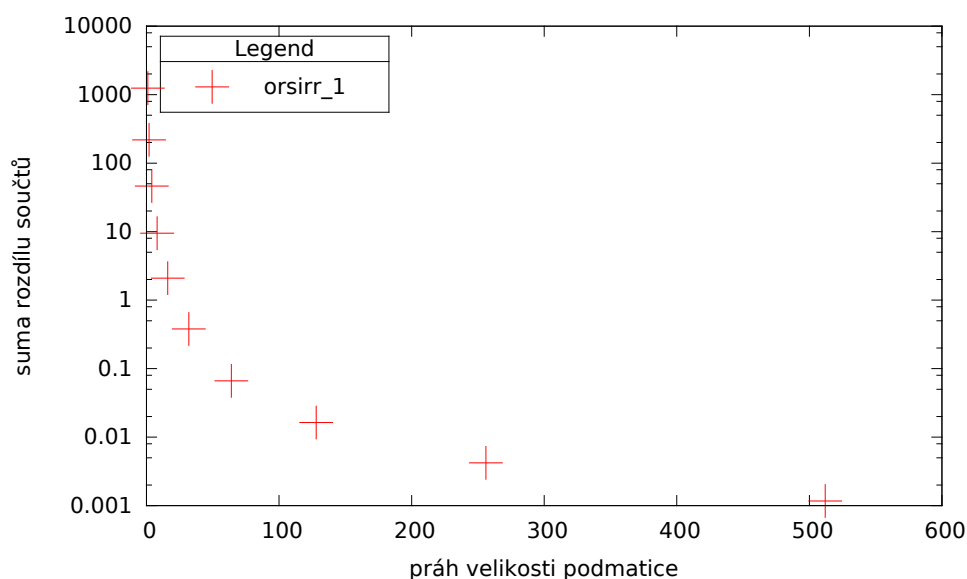
Správný výsledek je $30.234 \times 0.8912 + 0.5678 \times 0.7891 = 26.9445408 + 0.44805098 = 27.39259178$.

Nyní výpočet provedeme pomocí Strassenova algoritmu:

$$\begin{aligned} & \begin{pmatrix} 30.234 & 0.5678 \\ 0.9123 & 10.456 \end{pmatrix} \cdot \begin{pmatrix} 0.8912 & 0.3456 \\ 0.7891 & 9.999 \end{pmatrix} \\ M_1 &= (30.234 + 10.456) \cdot (0.8912 + 9.999) = 443.12 \\ & \dots \\ M_4 &= 10.456 \cdot (0.7891 - 0.8912) = -1.067 \\ M_5 &= (30.234 + 0.5678) \cdot 9.999 = 307.98 \\ & \dots \\ M_7 &= (0.5678 - 10.456) \cdot (0.7891 + 9.999) = -106.67 \\ & \begin{pmatrix} M_1 + M_4 - M_5 + M_7 & \dots \\ \dots & \dots \end{pmatrix} = \begin{pmatrix} 27.403 & \dots \\ \dots & \dots \end{pmatrix} \end{aligned}$$

Jak můžeme vidět, zatímco u algoritmu podle definice jsme pouze ztratili desetinnou přesnost, výsledek Strassenova algoritmu se lišil už v prvním desetinném čísle.

Pro reálnou představu stability Strassenova algoritmu jsme provedli experiment, ve kterém jsme vynásobili dvě stejné matice (matice `orsirr_1`, oříznuta na velikost 1024) algoritmem podle definice a Strassenovým algoritmem s různými prahy a sečetli všechny rozdíly mezi výsledky. Násobení probíhalo ve dvojité desetinné přesnosti, tedy v datovém typu `double` jazyka C99. Z grafu 2.2 je vidět exponenciální závislost mezi velikostí prahu a celkovou chybou.



Obrázek 2.2: Ukázka numerické stability Strassenova algoritmu

Strassenův algoritmus lze ještě vylepšit. Algoritmům na stejném principu se říká Strassen-like [9]. Pro sedm operací násobení je možné snížit počet sčítání a odečítání. Pro jednoduchost zde uvádíme pouze originální algoritmus.

2.6 Rychlé algoritmy

Potom, co Strassen ukázal, že existují rychlejší algoritmy než $\mathcal{O}(n^3)$, ještě rychlejší algoritmy než ten jeho na sebe nenechaly dlouho čekat. Složitost násobení matic pro jednoduchost označíme jako $\mathcal{O}(n^\omega)$.

Nejpomalejší algoritmus s $\omega = 3$ je podle definice. Strassenův algoritmus se sedmi násobeními má $\omega \approx 2.807354$. Jeden z nejrychlejších algoritmů je algoritmus Virginie Williamsové [45][44], pro který je $\omega < 2.3727$.

Hranicí nejlepší možné složitosti může být $\omega = 2$, protože každý prvek z matice musíme nějak započítat. Existují z velké části podložené domněnky na základě teorie grup [34], že $\omega = 2$ skutečně platí, ale přímý důkaz ještě neexistuje.

2.7 Algoritmus podle definice upravený pro řídké matice

Pokud násobíme řídké matice A a B o velikosti N , můžeme vynechat násobení takových dvou prvků, z nichž je alespoň jeden nulový. Označme $nnzr_{M,i}$

jako počet nenulových prvků v i -tém řádku matice M a $nnz_{C_{M,i}}$ jako počet nenulových prvků v i -tém sloupci matice M . Protože násobíme každý řádek s každým sloupcem, bude celkový počet operací násobení dán vzorcem:

$$\sum_{i=1}^n nnz_{A,i} \cdot nnz_{B,i} \quad (2.3)$$

Složitost tohoto algoritmu pro násobení dvou matic A a B o velikosti n tedy můžeme vyjádřit jako $O(mn)$, kde $m = \max(nnz(A), nnz(B))$. Pro husté matice bez jediného nulového prvku samozřejmě platí, že $m = n^2$. Nutno podotknout, že $O(mn)$ je nejhorší případ, kdy se všechny nenulové prvky matice A budou násobit se všemi nenulovými prvky matice B .

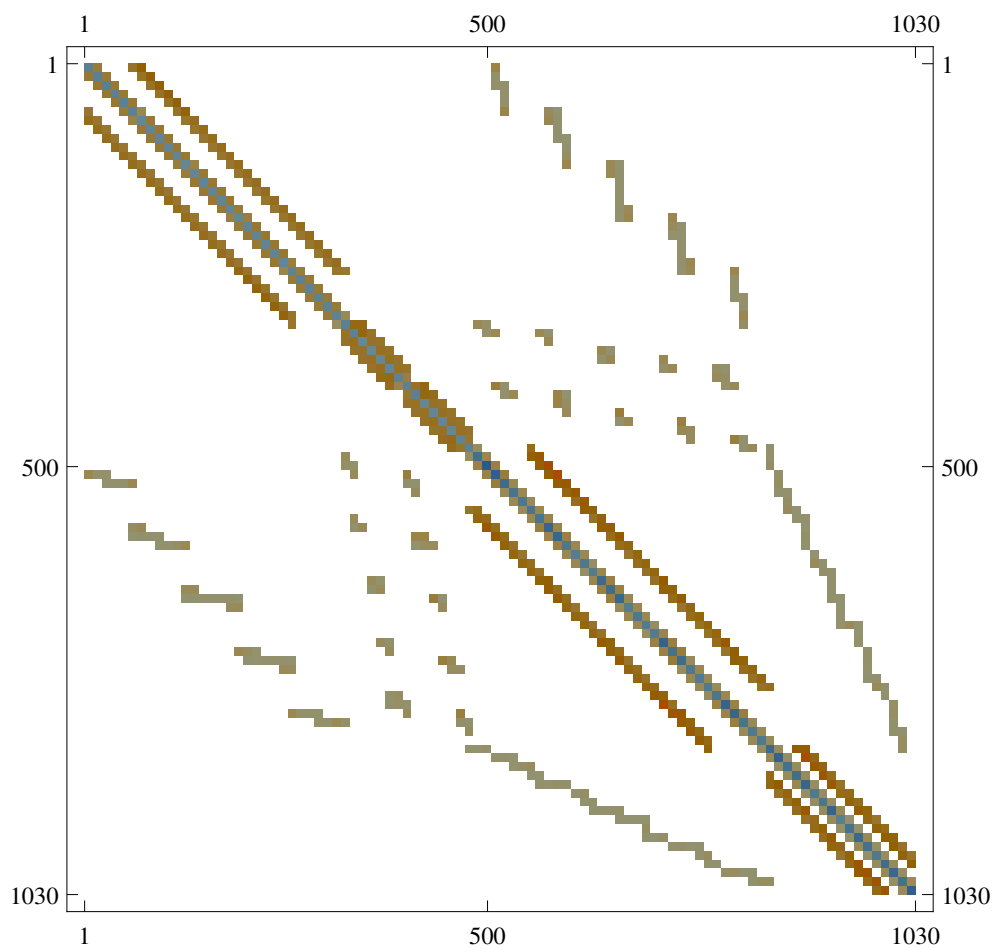
Pro reálnou představu demonstrujeme násobení dvou stejných matic. Opět se jedná o dvě matice `orsirr_1` [25]. Matice `orsirr_1` má velikost $n = 1030$ a $m = 6858$ nnz. Rozmístění nnz před respektive po vynásobení ukazují obrázky 2.3 respektive 2.4.

Kdyby nastal nejhorší možný případ, počet operací násobení při použití algoritmu podle definice pro řídké matice by byl $1030 \times 6858 = 7.063740 \times 10^6$. Pro tento případ ovšem stačí 4.6976×10^4 operací násobení. Oproti nejhoršímu možnému případu nastalo $7.063740 \times 10^6 - 4.6976 \times 10^4 = 7.016764 \times 10^6$ situací, kdy jeden ze dvou prvků byl nulový a operace násobení nemusela být provedena. Při násobení algoritmem podle definice pro husté matice by v 1.092680024×10^9 případech operace násobení nemusela být provedena, protože alespoň jeden ze dvou prvků by byl nula. Po vynásobení matice `orsirr_1` sama se sebou stoupl její počet nnz z 6858 na 23532. V tomto případě to bylo především z důvodu, že všechny prvky z diagonály matice jsou nenulové, každý prvek se tedy započítá.

2.8 Rychlé násobení řídkých matic

Strassenův algoritmus $\omega \approx 2.807354$ od $nnz < n^{1.8}$ a algoritmus Virginie Williamsové $\omega < 2.3727$ od $nnz < n^{1.37}$ jsou asymptoticky stejně rychlé jako algoritmus podle definice upravený pro řídké matice $O(mn)$. Například pro $n = 1000$ je tato hranice pro Strassenův algoritmus 251189 (40 %) nnz a pro algoritmus Virginie Williamsové 12882 (78 %) nnz z celkových možných 1000000 prvků.

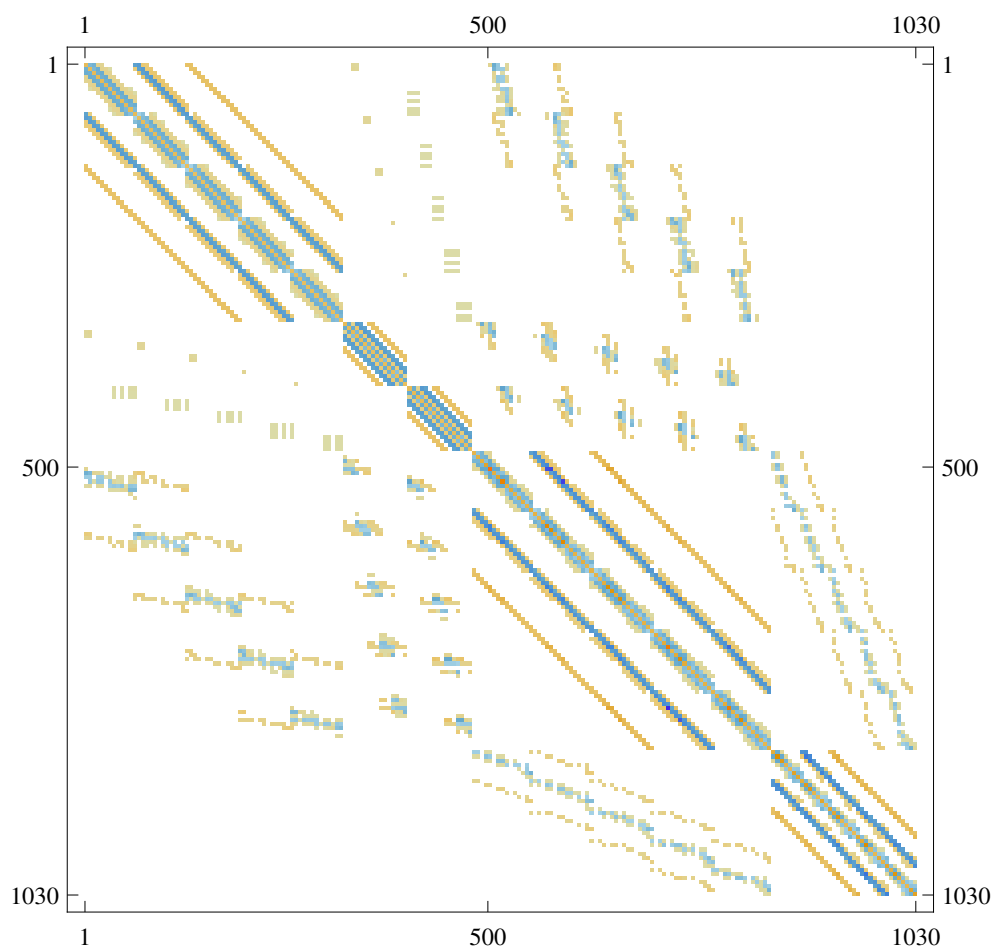
Raphael Yuster a Uri Zwick ukázali algoritmus [49] s asymptotickou složitostí $\mathcal{O}(m^{0.7}n^{1.2} + n^{2+o(1)})$, který rozdělí permutace řádků a sloupců na řídké a husté. Řídké permutace násobí algoritmem podle definice v úpravě pro řídké matice $\mathcal{O}(mn)$. Husté permutace vynásobí v té době nejznámějším nejrychlejším algoritmem, a to od Dona Coppersmitha a Shmuela Winograda s asymptotickou složitostí $\omega = 2.375477$ [10].



Obrázek 2.3: Matice orsirr_1 před vynásobením

2.9 Další algoritmy pro řídké matice

Algoritmů násobení řídkých matic je mnoho. Často se odvíjejí od typu řídkých matic a formátu, v jakém jsou uloženy. Příkladem může být formát, ve kterém se ukládají diagonály [40].



Obrázek 2.4: Matice `orsirr_1` po vynásobením sama se sebou

Formáty uložení řídkých matic

Formáty uložení řídkých matic obecně ukládají jednotlivé prvky zvlášť, a tedy nemusí ukládat ty nulové. To ale přináší řadu nevýhod. Za prvé se musí ukládat informace o souřadnicích a za druhé ztrácíme možnost přístupu k prvku na libovolných souřadnicích v čase $\Theta(1)$, protože prvky nemáme přímo indexované podle jejich umístění v řádku a sloupci.

Protože řídké matice můžeme rozdělit do mnoha kategorií a provádět nad nimi mnoho operací, existuje mnoho formátů, jak řídkou matici efektivně uložit a pracovat s ní. Formáty uložení řídkých matic můžeme také rozdělit podle toho, zdali je možné do nich přidávat nebo z nich odebírat prvky.

Při násobení matic $C = A \cdot B$ se matice A ani B nemění. V této práci budeme předpokládat, že matice C bude hustá a formáty umožňující přidávání nebo odebírání prvků nebudou součástí práce. Stejně tak při násobení matice A vektorem B je výsledek C vektor.

Dalším kritériem pro rozdělení formátů je uspořádanost nenulových prvků v řídké matici. Pro uspořádané matice bude efektivnější takový formát, který využije určitý vzor, v jakém jsou prvky matice uloženy. V řídkých maticích takovým vzorem může být například diagonála nebo blok prvků. Efektivně lze za vzor považovat i prvky v řádku nebo ve sloupci.

Uspořádáním může být také symetrie matice, kdy nám stačí uložit pouze polovinu matice. Často bývají řídké matice symetrické podle hlavní diagonály.

3.1 COO - Coordinate list

Formát COO, česky seznam souřadnic, je základní formát řídkých matic. Ke každému nenulovému prvku ukládá jeho souřadnice y a x . Implementovat tento formát můžeme například jako tři pole. Pole v s hodnotami prvků, pole r s y souřadnicemi a pole c s x souřadnicemi.

Pro ukázkou v tomto formátu uložíme matici o velikosti $n = 8$ s $nnz = 5$ nenulovými prvky.

3. FORMÁTY ULOŽENÍ ŘÍDKÝCH MATIC

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \mathbf{1} & \mathbf{2} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \mathbf{3} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \mathbf{4} & \mathbf{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$\begin{array}{l|lllll} \text{values}[5] & 1 & 2 & 3 & 4 & 5 \\ \text{y-coords}[5] & 1 & 1 & 4 & 7 & 7 \\ \text{x-coords}[5] & 0 & 1 & 1 & 5 & 6 \end{array}$$

Tabulka 3.1: Matice uložená ve formátu COO

Jak můžeme vidět, délka polí je závislá pouze na nnz . Pro velké matice s malým počtem neuspořádaných nenulových prvků je tento formát velmi efektivní. Pokud by bylo prvků velké množství, informace o uložení y nebo x souřadnic by byla často redundantní.

Paměťová náročnost formátu COO je $O(3 \cdot nnz)$.

Formát COO je velmi jednoduchý a přímočarý. Při procházení jeho prvků nám stačí jedna iterace přes tři stejně dlouhá pole. Tím je tedy například algoritmus násobení řídké matice ve formátu COO s vektorem velmi jednoduchý:

Algorithm 6 Násobení matice COO s vektorem

```

1: procedure COO-MVM(COO,V,C)
2:   for  $i \leftarrow 0$  to COO.nnz do
3:      $V.v[\text{COO}.r[i]] \mathrel{+}= \text{COO}.v[i] * V.v[\text{COO}.c[i]]$ ;
4:   end for
5: end procedure

```

Při násobení dvou matic narazíme na problém. Při této operaci se každý prvek násobí dvakrát. Potřebujeme způsob, jak se v matici vrátit zpátky na určité místo. Takový naivní algoritmus pro násobení dvou COO matic by byl složitý. Museli bychom si pamatovat začátky řádků a neustále kontrolovat, jestli jsme nepřesáhli další řádek. Lepším řešením je dopředu si předpočítat, kde který řádek začíná a končí. Předpočítáme si tedy pole, nazvané *row_pointers*, o délce $M.height + 1$, obsahující indexy začátků a konců řádek.

Při procházení matice se tak v poli s prvky mezi indexy $row_pointers[i]$ a $row_pointers[i + 1]$ nachází prvky na řádku i . Proto je pole právě o jedna delší než výška matice, abychom mohli určit konec posledního řádku.

Algorithm 7 Násobení dvou COO matic

```

1: procedure COO-MMM(A,B,C)
2:   arp  $\leftarrow$  ComputeRowBeginnings(A, A.nnz + 1);
3:   brp  $\leftarrow$  ComputeRowBeginnings(B, B.nnz + 1);
4:   for i  $\leftarrow$  0 to A.height do                                      $\triangleright$  násobení
5:     for ac $\leftarrow$ arp[i] to arp[i + 1] do
6:       for bc $\leftarrow$ brp[A.c[ac]] to brp[A.c[ac] + 1] do
7:         C.v[r][A.c[bc]] += A.v[ac] * B.v[bc];
8:       end for
9:     end for
10:  end for
11: end procedure

```

V části s násobením již pole s y souřadnicemi nepotřebujeme. Lepší formát uložení řídkých matic pro násobení by byl takový, který namísto pole s y souřadnicemi obsahuje předpočítané začátky a konce řádků. Takovým formátem je CSR.

3.1.1 Formát MatrixMarket

MatrixMarket [7] je internetová sada matic s vlastním formátem pro uložení řídkých matic `.mtx`. Tato sada obsahuje skoro pět set matic z různých oblastí. Obsahuje i generátory řídkých matic, jejichž výstupy jsou matice různých vlastností.

Formát MatrixMarket se velmi podobá formátu COO. Navíc umožňuje nastavit matici jako symetrickou a uložit pouze její polovinu.

3.2 CSR - Compressed sparse row

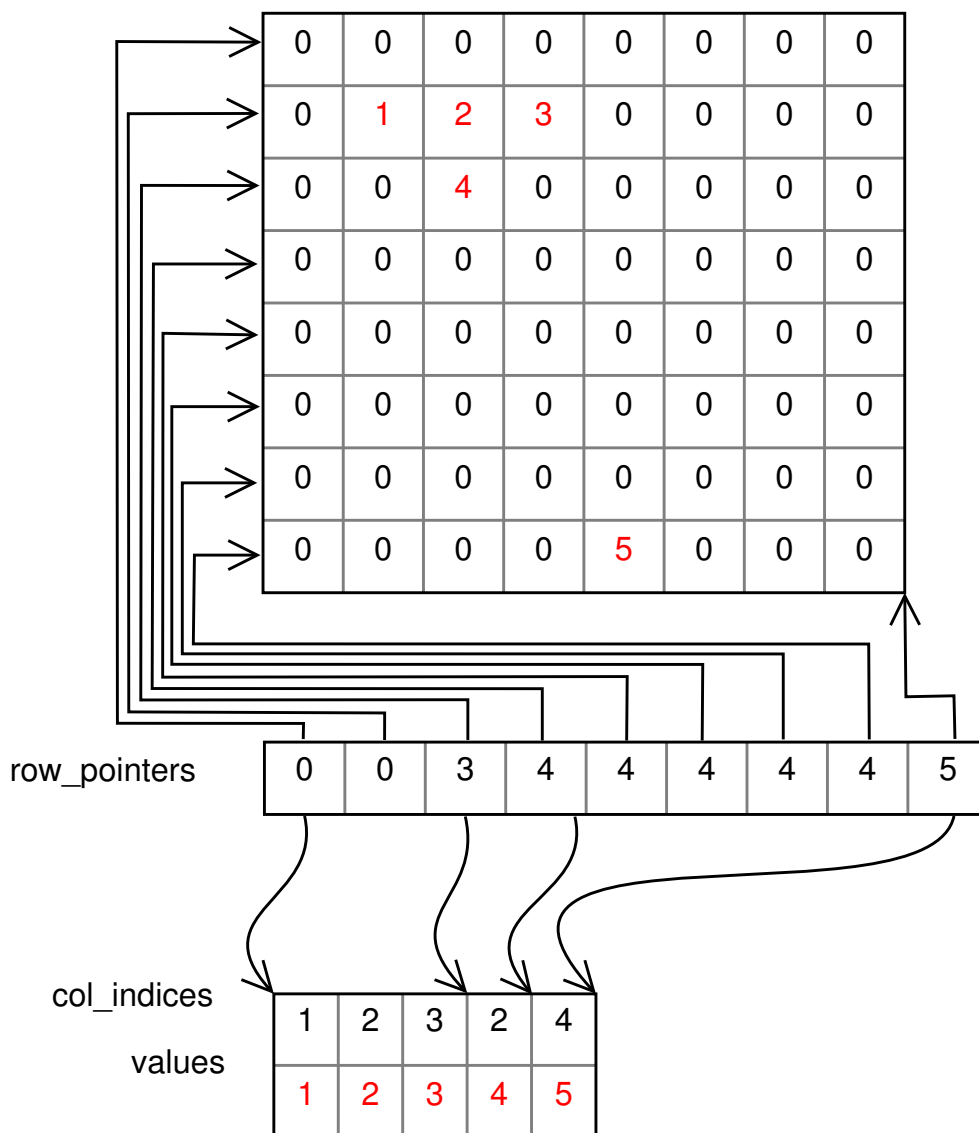
Problém efektivnosti formátu COO pro větší množství prvků řeší formát CSR, česky komprimované řídké řádky. Formát CSR obsahuje pole *row_pointers*, které ukládá informace o tom, kolik se v daném řádku nachází prvků, tedy to, co jsme si předpočítali v algoritmu 7. K poli s hodnotami je další pole *col_indices*, přiřazující ke každému prvku informaci o sloupci.

Jak je vidět z ilustrace 3.1, řádek s více prvky je uložen efektivně. Kvůli mnoha prázdným řádkům se ale pole *row_pointers* nezdá rozumně využité.

Při násobení matice CSR s vektorem potřebujeme o jeden for cyklus více než v případě násobení matice COO s vektorem. Důvodem je ztráta informace o řádku prvku.

Násobení dvou CSR matic je stejné jako v případě násobení dvou COO matic 3.1. Jediný rozdíl je, že předpočítané začátky a konce řádků jsou součástí formátu.

3. FORMÁTY ULOŽENÍ ŘÍDKÝCH MATIC



Obrázek 3.1: Matice uložená ve formátu CSR

Pro uložení matice ve formátu CSR potřebujeme pole o velikost $n + 1$ pro uložení informací o řádcích a $2 \cdot nnz$ pro prvky a jejich sloupce. Paměťová složitost formátu CSR je $\mathcal{O}(n + 1 + 2 \cdot nnz)$.

Existuje varianta tohoto formátu, nazvaná CSC - compressed sparse columns, která místo ukládání řádku ukládá sloupce.

Algorithm 8 Násobení matice CSR s vektorem

```

1: procedure CSR-MVM(CSR,V,C)
2:   for  $i \leftarrow 0$  to CSR.h do
3:     for  $ci \leftarrow \text{CSR.rp}[i]$  to  $\text{CSR.rp}[i + 1]$  do
4:        $C.v[r] += \text{CSR.v}[ci] * V.v[A.ci[ci]];$ 
5:     end for
6:   end for
7: end procedure

```

Algorithm 9 Násobení dvou CSR matic

```

1: procedure CSR-MMM(A,B,C)
2:   for  $i \leftarrow 0$  to A.height do ▷ násobení
3:     for  $ac \leftarrow A.rp[i]$  to  $A.rp[i + 1]$  do
4:       for  $bc \leftarrow B.rp[A.ci[ac]]$  to  $B.rp[A.ci[ac] + 1]$  do
5:          $C.v[r][B.ci[bc]] += A.v[ac] * B.v[bc];$ 
6:       end for
7:     end for
8:   end for
9: end procedure

```

3.3 BSR - Block Sparse Row

Jako formát CSR využívá uložení prvků v řádku, formát BSR [38][26] ještě navíc detekuje a ukládá prvky v blocích.

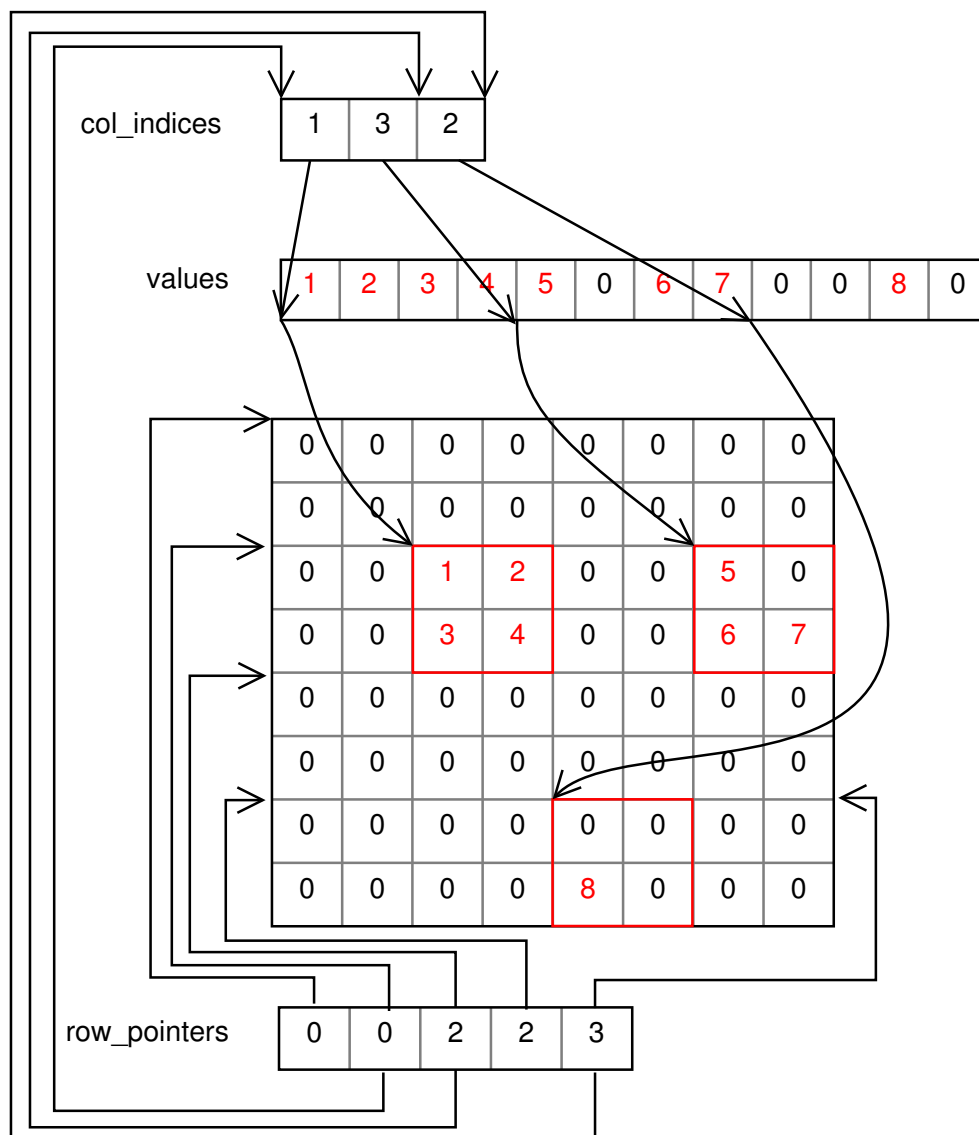
Protože při násobení matic násobíme každý prvek dvakrát, bylo by dobré tyto dvě operace provést co nejdříve po sobě, abychom při druhém znovunačtení prvku mohli sáhnout pro prvek do cache. Pokud jsou prvky procházené po menších blocích, dostaneme se k prvku podruhé dříve, než jej z cache přemaže jiný prvek.

Matici A ve formátu BSR ukládáme velice podobně jako u formátu CSR pomocí tří polí. Je potřeba i jedna proměnná, která uchovává velikost bloku. Tuto proměnnou nazveme *block_size*. Pole *col_indices* označuje sloupec, ve kterém se blok nachází. Sloupcem rozumíme blok prvků o šířce $A.width/A.block_size$.

Uložení matice ve formátu BSR ilustruje obrázek 3.2. Pole *row_pointers* obsahuje informace o tom, na kterém řádku je kolik bloků. Na řádku i je $row_pointers[i + 1] - row_pointers[i]$ bloků. Ve kterém sloupci se blok prvků nachází, udává pole *col_indices*. První blok z řádku i je ve sloupci $col_indices[row_pointers[i]]$ a poslední blok je ve sloupci $col_indices[row_pointers[i + 1]]$. Pole *values* obsahuje prvky v blocích, včetně nulových prvků.

Pokud vezmeme algoritmus násobení dvou CSR matic respektive CSR matice s vektorem, jen místo prvků násobíme bloky, vznikne nám algoritmus pro násobení dvou BSR matic, respektive BSR matice s vektorem. Za povšimnutí

3. FORMÁTY ULOŽENÍ ŘÍDKÝCH MATIC



Obrázek 3.2: Matice uložená ve formátu BSR

stojí větší počet for cyklů s menším rozsahem iterace než u CSR. To nám v nejvnitřnějších cyklech dovoluje lépe využívat cache. Jedná se tedy o přístup podobný optimalizační technice loop tiling[46].

Algorithm 10 Násobení matice BSR s vektorem

```

1: procedure BSR-MVM(A,B,C)
2:    $bs \leftarrow A.block\_size;$ 
3:   for  $i \leftarrow 0$  to  $A.height / bs$  do
4:     for  $ac \leftarrow A.rp[i]$  to  $A.rp[i + 1]$  do
5:       for  $l \leftarrow 0$  to  $A.bs$  do ▷ násobení bloku
6:         for  $m \leftarrow 0$  to  $bs$  do
7:            $C.v[(i * bs) + l] += A.v[ac * (bs * bs) + (l * bs) + m] * B.v[A.ci[ac] * bs + m];$ 
8:         end for
9:       end for
10:    end for
11:  end for
12: end procedure

```

Algorithm 11 Násobení dvou BSR matic

```

1: procedure BSR-MMM(A,B,C)
2:    $bs \leftarrow A.block\_size;$ 
3:   for  $i \leftarrow 0$  to  $A.height / bs$  do
4:     for  $ac \leftarrow A.rp[i]$  to  $A.rp[i+1]$  do
5:       for  $bc \leftarrow B.rp[A.ci[ac]]$  to  $B.rp[A.ci[ac]+1]$  do
6:         for  $l \leftarrow 0$  to  $A.bs$  do ▷ násobení bloku
7:           for  $m \leftarrow 0$  to  $bs$  do
8:             for  $n \leftarrow 0$  to  $bs$  do
9:                $C.v[(i * bs) + l][(B.ci[bc] * bs) + m]$ 
10:               $+= A.v[ac * (bs * bs) + (l * bs) + n] * B.v[bc * (bs * bs) + (n * bs) + m];$ 
11:             end for
12:           end for
13:         end for
14:       end for
15:     end for
16: end procedure

```

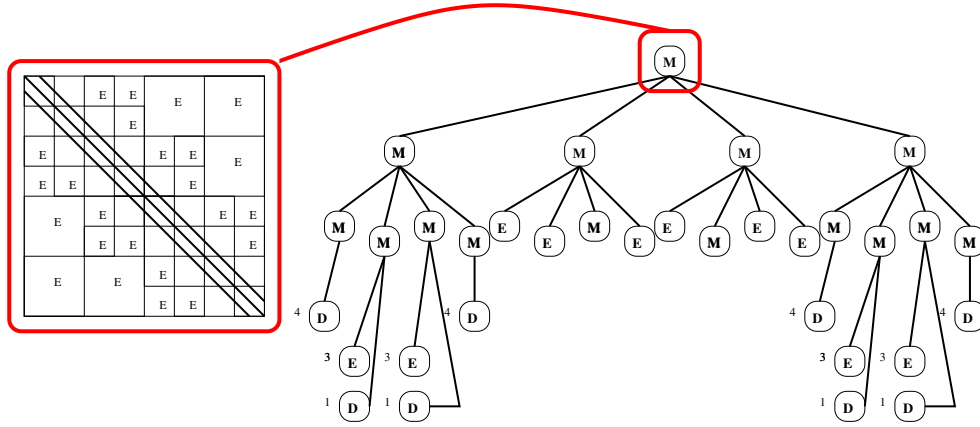
Paměťová složitost je závislá na počtu bloků b , velikosti bloků $block_size$ a počtu řádek matice n . Pole *row_pointers* má velikost $n + 1$, pole *col_indices* má velikost b a pole hodnot *values* má velikost $b \cdot sm_size^2$. Celková velikost matice uložené ve formátu BSR je $\mathcal{O}(n + b + b \cdot sm_size^2)$.

3.4 Quadtree

Předchozí popsané formáty uložení řídkých matic jsou velmi přímočaré. Neumožňují rekursivní přístup, který je potřebný například pro Strassenův algoritmus. Tento problém řeší formát Quadtree [32][1]. Jedná se o matici uloženou v 4-árním stromě.

Tento formát dělí matici na čtvrtiny do té doby, než se dosáhne velikosti podmatice sm_size , nebo ve čtvrtině nejsou žádné prvky. Pokud je celá čtvrtina prázdná, je uzel označen jako prázdný, E . Pokud obsahuje nějaké prvky, je list označen jako hustý, D . Vnitřní uzly jsou označeny jako smíšené, M .

Obrázek 3.3 ukazuje, že kořen stromu Quadtree představuje celou matici. Jeho synové představují čtvrtinu matice, jak je v obrázku 3.4. Dělení matice se zastaví, pokud je blok k dělení prázdný, jako například v 3.5, nebo byla dosažena velikost sm_size .



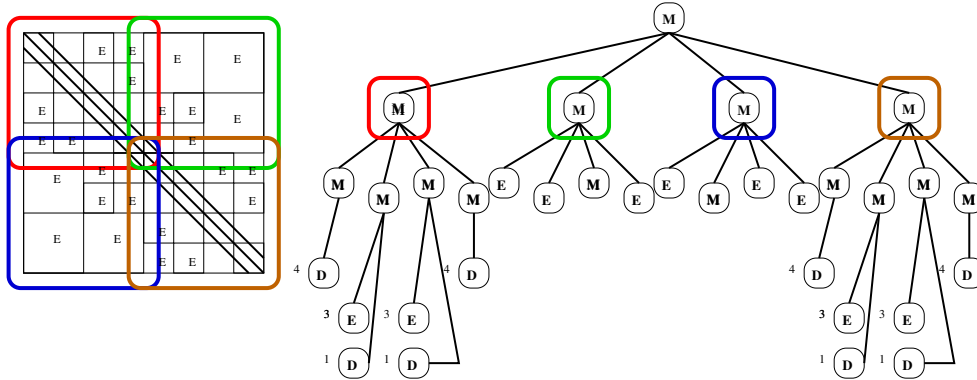
Obrázek 3.3: Rozdělení matice

Výška stromu Quadtree je při reprezentaci matice o velikosti n a při velikosti podmatice sm_size :

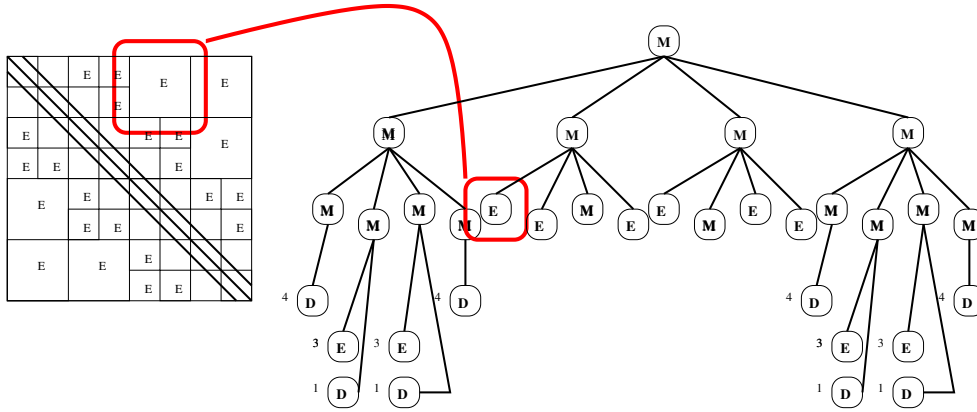
$$\left\lceil \log_4 \frac{n^2}{sm_size^2} \right\rceil \quad (3.1)$$

Například pro matici o velikosti $n = 16384$ s velikostí bloku $sm_size = 128$ je výška stromu 7.

Algoritmy pro práci s formátem Quadtree budou ukázány v kapitole o obecnějším formátu 3.5.



Obrázek 3.4: Rozdělení matice



Obrázek 3.5: Rozdělení matice

3.5 Modifikace formátu quadtree

V této práci navrhujeme místo 4-árního stromu obecný k -ární strom. Výška stromu se tím snížší. Pro předchozí příklad 3.4 by tedy pro 16-ární strom byla výška pouze 4. Tento formát nazveme *k-ary tree matrix* se zkratkou KAT. Pro přehlednost nebudeme uvádět k , ale $KAT.n = \text{sqrt}(k)$. Formát Quadtree můžeme prohlásit za matici KAT s $KAT.n = 2$. Pro snadnější rekursivní přístup budeme jako $KAT.n$ uvažovat pouze mocniny dvou.

Výška stromu pro KAT matici označíme jako $KAT.height$ a vypočítáme ji jako výšku Quadtree, jen s parametrem k :

$$KAT.height = \left\lceil \log_k \frac{n^2}{sm_size^2} \right\rceil \quad (3.2)$$

Pokud to bude z kontextu jasné, budeme výšku KAT stromu značit pouze *height*.

3.5.1 Typy listů

KAT matice obsahuje proměnnou *KAT.dense_threshold*, která udává maximální počet prvků, aby se s listem pracovalo jako s některým z řídkých formátů, tedy COO, CSR, BSR. Při překročení této hranice je list uložen ve formátu husté matice, tedy bude obsahovat i nulové prvky.

3.5.2 Násobení

Pro násobení budeme používat algoritmus rekurzivního násobení popsaného v 2.4. Popsaný algoritmus dělí matice na čtvrtiny, jde tedy aplikovat na formát Quadtree. Při násobení matice uložené v *k*-árním stromě budeme násobit matici podmatic o velikosti *k*:

Algorithm 12 Násobení matice KAT s vektorem

```

1: procedure KAT-MVM(KAT, KAT_node, VB, VC)
2:   for i ← 0 to KAT.n do
3:     for j ← 0 to KAT.n do
4:       if KAT_node.childs[i][j] ≠ NIL then
5:         if KAT_node.childs[i][j].type = "submatrix" then
6:           multiplyNode(KAT_node.childs[i][j], VB, VC);
7:           continue;
8:         end if
9:         if KAT_node.childs[i][j].type = "inner" then
10:          KAT-MVM(KAT, KAT_node.childs[i][j], VB, VC);
11:          continue;
12:        end if
13:      end if
14:    end for
15:  end for
16: end procedure

```

Přesnější asymptotická složitost násobení matic v KAT formátu je vyšší než u předchozích formátů, protože je potřeba připočítat průchod stromem. Pokud převedeme hustou matici do formátu KAT se složitostí násobení dvou listů $\mathcal{O}(sm_size^3)$, bude asymptotická složitost násobení dvou KAT matic $\mathcal{O}((\frac{n^2}{sm_size} \cdot (height + sm_size^3)))$.

3.5.3 Paměťová složitost

Paměťová složitost celé KAT matice záleží na typu a hustotě uložení dat. Uvedeme zde pouze paměťovou složitost stromu bez listů. Velikost uzlu je

Algorithm 13 Násobení dvou KAT matic

```

1: procedure KAT-MMM(KATa, KATa_node, KATb, KATb_node, C)
2:   for i  $\leftarrow$  0 to KAT.n do
3:     for j  $\leftarrow$  0 to KAT.n do
4:       for k  $\leftarrow$  0 to KAT.n do
5:         if KATa_node.chilids[i][k]  $\neq$  NIL and
KATb_node.chilids[k][j]  $\neq$  NIL then
6:           if KAT_node.chilids[i][j].type = "submatrix" then
7:             multiplyNode(KAT_node.chilids[i][j],
KAT_node.chilids[i][j], VC);
8:             coninue;
9:           end if
10:          if KAT_node.chilids[i][k].type = "inner" then
11:            KAT-MMM(KATa, KATa_node.chilids[i][k],
KATb, KATb_node.chilids[k][j], C);
12:            coninue;
13:          end if
14:        end if
15:      end for
16:    end for
17:  end for
18: end procedure

```

pole k ukazatelů na syny. Počet vnitřních uzlů spočítáme pomocí geometrické posloupnosti. Paměťová složitost stromu KAT matice tedy je $\mathcal{O}(\frac{k^{height-1}-1}{k-1} \cdot k)$ po zkrácení $\mathcal{O}(k^{height-2} - k - \frac{1}{k})$.

Analýza a návrh

V této kapitole představíme některé známé systémy, které umí pracovat s řídkými maticemi. Popíšeme možnosti implementace a zvolené řešení.

Pokud potřebujeme rychle vynásobit dvě malé matice, můžeme použít obecný nástroj, se kterým se snadno pracuje. Čím je ale nástroj obecnější, tím spíš bude pomalejší. Některé nástroje jsou ale specificky optimalizované a i přes jejich obecnost mají dobrou výkonnost. Maximální efektivnosti dosáhneme, pokud co nejvíce využijeme hardware a náš program bude řešit konkrétní úkol.

4.1 Software pro práci s řídkými maticemi

4.1.1 SciPy.sparse

SciPy je knihovna skriptovacího jazyka Python [30]. Obsahuje modul pro základní práci s řídkými maticemi nazvaný sparse [39]. Samotné výpočty nad řídkými maticemi jsou v této knihovně z důvodu rychlosti napsány v jazyce C++.

4.1.2 Boost

Pro práci s řídkými maticemi v jazyce C++ je možné použít knihovnu uBLAS [43] ze sady knihoven Boost[22]. Knihovna uBLAS obsahuje algoritmy pro řešení základních úkolů z lineární algebry. Protože je kód napsaný pomocí C++ šablon, je vygenerovaný kód výkonný.

4.1.3 ATLAS

ATLAS [42] je v překladu zkratkou pro automatické generování vyladěného kódu pro software řešící úlohy z lineární algebry. Tohoto je dosaženo jak pomocí obecných optimalizací kódu, tak optimalizací pro konkrétní architektury.

Podobných projektů je více, konkrétně pro násobení matice s vektorem je to například Sparsity [41] nebo pro násobení matice s maticí je to například PHiPAC [36][6].

4.1.4 Wolfram Mathematica

Wolfram Mathematica [47] je mocný výpočetní software. Práce v tomto software je snadná, interaktivní a přitom dokáže efektivně rozdělovat práci nejen na více procesorových jader, ale také využít grafické karty. Nechybí tedy ani implementace operací pro řídké matice.

V této práci pomocí Wolfram Mathematicy vizualizujeme řídké matice z formátu `.mtx` 3.1.1. Mathematica umí pracovat i s maticemi v souborech `.mtx.gz`, tedy komprimovanými soubory. Příkaz pro importování řídké matice a její vizualizaci je `Import["/var/tmp/matrix.mtx", "Graphics"]`, následné uložení matice lze provést například příkazem `Export["/var/tmp/matrix.pdf", %]` [48].

4.2 Řešení implementace

Jednou z možných variant bylo pro porovnání všech zvolených formátů byla doimplementovat formát Quadtree, respektive KAT do knihovny SciPy.sparse. Další možností bylo pokračovat v semestrální práci z předmětu Efektivní implementace algoritmů s tématem násobení matic ve formátu CSR.

Protože se v této práci zabýváme pouze vlivem uložení řídké matice pro operaci násobení, rozhodli jsme se vytvořit program pouze pro tuto operaci. Obecná knihovna pro práci s řídkými maticemi, jakým například SciPy.sparse je, nemůže využít některých vlastností matic při jejich násobení. Hlavně zacházení s násobenými maticemi jako s konstantami. Navíc pokud nebudou použity žádné jiné implementace, budou všechny algoritmy stejně optimalizované. Navážeme tedy na semestrální práci. Pro každý formát vytvoříme ve zdrojovém kódu třídu, v níž implementujeme operaci násobení s maticí a násobení s vektorem.

Realizace

5.1 Prostředí

Popsané algoritmy jsme implementovali v jazyce C99 [27]. Implementace probíhala v operačním systému Lubuntu [29], v IDE Eclipse CDT [37]. Zdrojový kód byl kompilován pomocí překladače GNU C [2], s optimalizačním přepínačem `-Ofast`.

K ladění chyb pro práci s pamětí jsme použili nástroj Valgrind [5] s přepínači `--leak-check=full --show-reachable=yes`. Běžící program jsme krokovali pomocí nástroje GNU Debugger [3]. Zdrojový kód pro ladění chyb byl kompilován s přepínači `-Wall -pedantic -Og -g`.

5.2 Design implementace

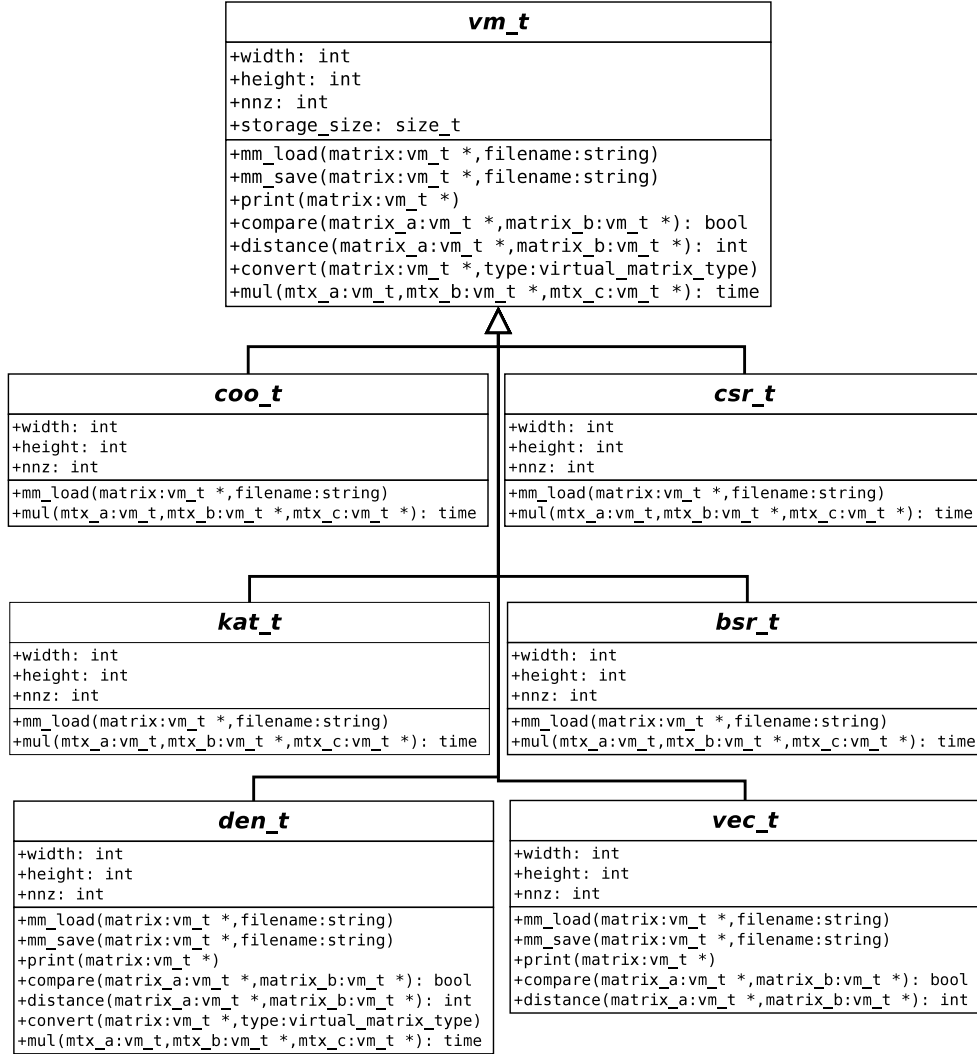
Po různých formátech požadujeme stejné operace, proto v programu zavádíme abstraktní matici a jednoduchý objektově orientovaný model [31], která formáty dědí z virtuální matice. Obrázek 5.1 ukazuje jednoduchý návrh tříd v UML.

Program běží v příkazovém řádku bez GUI a je jednovláknový. Načte jednu nebo dvě matice v určitém formátu a vynásobí je. Pomocí přepínačů lze zobrazit nebo uložit výslednou matici, popřípadě vypsát některé údaje o načtených maticích. Podrobnější popis nastavení programu je v příloze.

5.3 Implementace KAT

U formátů COO, CSR, BSR je implementace přímočará podle pseudokódů. U implementace KAT máme více možností, proto zde popíšeme naši implementaci.

Jeden z důležitých parametrů je k , tedy maximální počet synů vnitřních uzlů. V naší implementaci je tento parametr uložen v konstantě `KAT.n`, pro



Obrázek 5.1: UML diagram tříd programu

připomenutí $KAT.n = \sqrt{k}$. Překladač poté cykly, kde iterujeme do této konstanty, rozbalí. Překladači i explicitně sdělujeme, ať cykly rozbalí přes atribut `__attribute__((optimize("unroll-loops")))`.

V bakalářské práci Rozšíření implementace formátu kvadrantového stromu od Tomáše Karabely [28] je Quadtree implementován jako strom, jehož listy tvoří virtuální matice podobným těm v naší práci. Protože jeho formát byl určený pro algoritmus LU rozklad, jeho virtuální matice musely umět přijímat i další prvky. Protože v naší práci se s formáty uložení řídké matice zachází jako s konstantami, rozhodli jsme se hodnoty prvků a informace v polích `row_pointers` a `col_indices` uložit mimo listy stromu. V listech se nachází

ukazatel do velkého pole na příslušné místo pro daný list. Ušetříme tak práci knihovně `libc` s mnohonásobným volání funkce `malloc`.

5.3.1 Tvoření stromu

Při tvoření matice v KAT formátu pro každý prvek procházíme strom a hledáme správný list. Protože výška stromu je daná třemi parametry, tedy velikostí matice n , velikostí podmatice sm_size a počtem větvení uzlu k , neefektivnější využití bude, pokud je n beze zbytku dělitelné $k * sm_size$. Pokud toto neplatí, rodiče listů budou mít část synů nevyužitých i při uložení husté matice. Protože se prvky do matice načítají po řádcích zleva doprava, je velká pravděpodobnost, že další načtený prvek bude patřit do stejného listu. Z tohoto důvodu si pamatujeme poslední list, a pokud prvek patří do něj, vrátíme jej. Algoritmus 14 ukazuje, jakým způsobem vyhledáváme list pro prvek.

5.3.2 Násobení listů matice KAT

Protože dovoluujeme dva druhy listů, tedy hustý a řídký ve formátu CSR, bylo potřeba implementovat následující algoritmy:

1. hustý list \cdot hustý list
2. hustý list \cdot CSR list
3. hustý list \cdot vektor
4. CSR list \cdot hustý list
5. CSR list \cdot CSR list
6. CSR list \cdot vektor

Protože tyto algoritmy jsou velice podobné již popsaným algoritmům 2, ukážeme zde pouze násobení CSR listu s hustým listem.

5.4 Testování

Pro ověření správnosti algoritmů je potřeba testovací software. Strategie testování v našem programu spočívá ve výběru testovacích matic, vynásobení v hustém formátu uložení, vynásobení v některém z řídkých formátů uložení a porovnání výsledků. Tento proces ukazuje algoritmus 16.

Při testování jsme i na relativně malých maticích narazili na problém s numerickou stabilitou. Část z testovacích matic má velký desetinný rozvoj, a protože čísla jsou v různých formátech násobena v jiném pořadí, dochází k rozdílu mezi výsledky. Při kompilaci programu lze nastavit přesnost výpočtu pomocí typu proměnné pro uchování hodnot buď na `float`, `double` nebo `long`

Algorithm 14 Vyhledání listu pro KAT matici

```
1: procedure KAT-GETNODE(KAT, y, x)
2:   if  $\{y, x\} \in \text{lastLeaf.area}$  then
3:     return lastLeaf;
4:   end if
5:   tmpNode  $\leftarrow$  KAT.root;
6:   blockY  $\leftarrow$  0; ▷ výřez matice
7:   blockX  $\leftarrow$  0;
8:   blockS  $\leftarrow$  KAT.sm_size;
9:   ▷ až do předposledního vnitřního uzlu traverzujeme podle velikosti
   KAT.n
10:  while blockS > (KAT.n * KAT.sm_size) do
11:    nodeY  $\leftarrow$  (y - blockY) / (blockS / KAT.n);
12:    nodeX  $\leftarrow$  (x - blockX) / (blockS / KAT.n);
13:    tmpNode  $\leftarrow$  tmpNode.childs[nodeY][nodeX];
14:    blockY += nodeY * (blockS / KAT.n);
15:    blockX += nodeX * (blockS / KAT.n);
16:    blockS /= KAT.n;
17:  end while
18:  ▷ do posledního vnitřního uzlu traverzujeme podle velikosti
   KAT.sm_size
19:  nodeY  $\leftarrow$  (y - blockY) / KAT.sm_size;
20:  nodeX  $\leftarrow$  (x - blockX) / KAT.sm_size;
21:  tmpNode  $\leftarrow$  tmpNode.childs[nodeY][nodeX];
22:  ▷ nyní je tmpNode list
23:  tmpNode.y  $\leftarrow$   $\lfloor y / \text{KAT.sm\_size} \rfloor * \text{KAT.sm\_size}$  ;
24:  tmpNode.x  $\leftarrow$   $\lfloor x / \text{KAT.sm\_size} \rfloor * \text{KAT.sm\_size}$  ;
25:  lastLeaf  $\leftarrow$  tmpNode;
26:  return tmpNode;
27: end procedure
```

double. Při porovnávání dvou matic neporovnáváme hodnoty přímo, ale sledujeme velikost jejich rozdílů. V defaultní konfiguraci je povolena odchylka 0.001, při kompilaci ji lze změnit. Pokud spustíme testy pouze s přesností float, část testů selže právě kvůli numerické stabilitě.

5.5 Měření

Měření probíhalo na serveru `star.fit.cvut.cz`. Měřené parametry byly:

1. Procentuální zrychlení výpočtu oproti formátu CSR
2. Počet uzlů ve stromě KAT matice

Algorithm 15 Násobení hustého KAT listu s CSR listem

```

1: procedure KAT-MMM-DEN-CSR( $ka, kb, na, nb, c$ )  $\triangleright ka, kb = \text{KAT}$ 
   matice;  $na, nb = \text{listy}$ ,  $c = \text{hustá matice}$ 
2:   for  $i \leftarrow 0$  to  $ka.sm\_size$  do
3:     for  $j \leftarrow na.rp[i]$  to  $na.rp[i]$  do
4:       for  $k \leftarrow 0$  to  $ka.sm\_size$  do
5:          $C.v[na.y + i][nb.x + k] += na.v[j] * nb.v[ka.sm\_size * na.ci[j] + j];$ 
6:       end for
7:     end for
8:   end for
9: end procedure

```

Algorithm 16 Testování

```

1: procedure TESTFORMATS(PairList, FormatList)
2:   for all  $pair \in \text{PairList}$  do
3:      $denseA \leftarrow vm\_load(pair.a, DENSE);$ 
4:      $denseB \leftarrow vm\_load(pair.b, DENSE);$ 
5:      $denseC \leftarrow vm\_mul(denseA, denseB, denseC);$ 
6:     for all  $format \in \text{FormatList}$  do
7:        $sparseA \leftarrow vm\_load(pair.a, format);$ 
8:        $sparseB \leftarrow vm\_load(pair.b, format);$ 
9:        $sparseC \leftarrow vm\_mul(sparseA, sparseB, sparseC);$ 
10:      if  $vm\_compare(denseC, sparseC) = \text{NOT\_SAME}$  then
11:         $print("Error:", pair.a, pair.b, format);$ 
12:      end if
13:    end for
14:  end for
15: end procedure

```

3. Velikost datových struktur matic

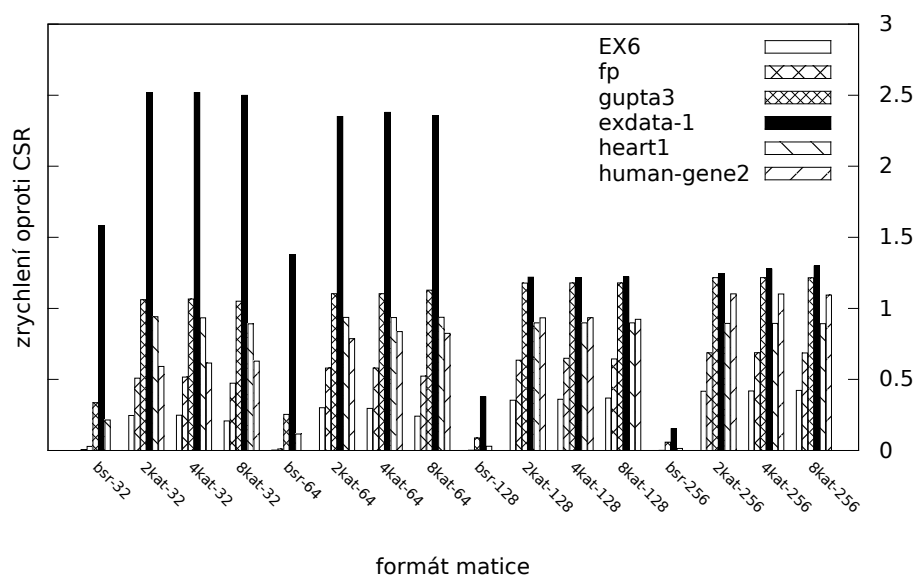
V grafech značíme matici KAT jako $KAT.nkatblock_size$. Matici BSR jako $bsrblock_size$.

Pro měření výsledků byly použity následující matice:

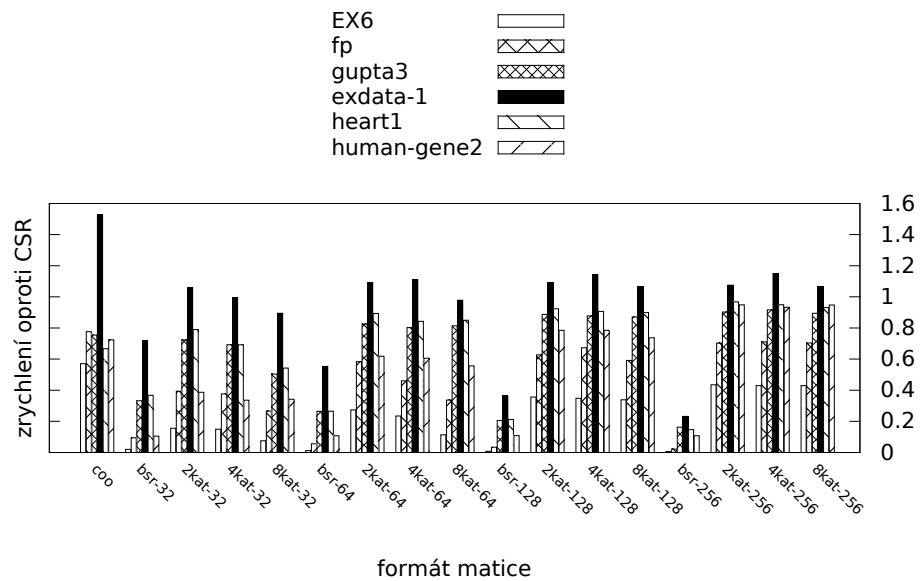
Skupina	Název matice	Velikost	Nenulových prvků	Oblast
GHS_indef	exdata_1 [15]	6001	$2.269500 \cdot 10^6$	optimalizace
Norris	heart1 [19]	3557	$1.385317 \cdot 10^6$	2D/3D
Gupta	gupta3 [16]	16783	$9323427 \cdot 10^6$	optimalizace
JGD_SPG	EX6 [17]	6545	$2.95680 \cdot 10^5$	kombinatorika
MKS	fp [18]	7548	$8.34222 \cdot 10^5$	elektromagnetika
Belcastro	human-gene2 [14]	14340	$1.8068388 \cdot 10^7$	graf

5.5.1 Rychlost výpočtu

Graf 5.2 ukazuje zrychlení výpočtu na ukázkových maticích oproti násobení ve formátu CSR. Graf 5.3 ukazuje stejný případ pro operaci násobení matice vektorem.



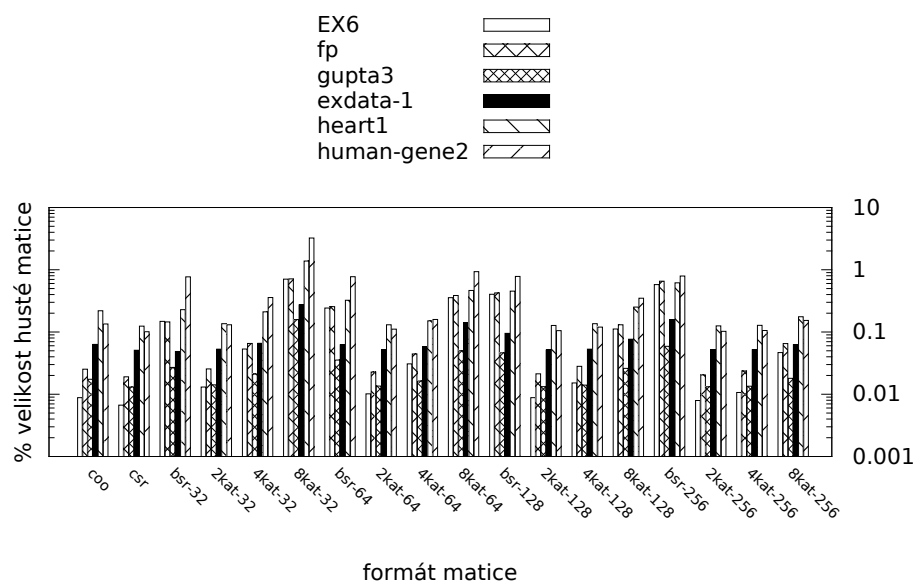
Obrázek 5.2: Zrychlení výpočtu matice \cdot matice oproti formátu CSR



Obrázek 5.3: Zrychlení výpočtu matice · vektor oproti formátu CSR

5.5.2 Velikost datových struktur

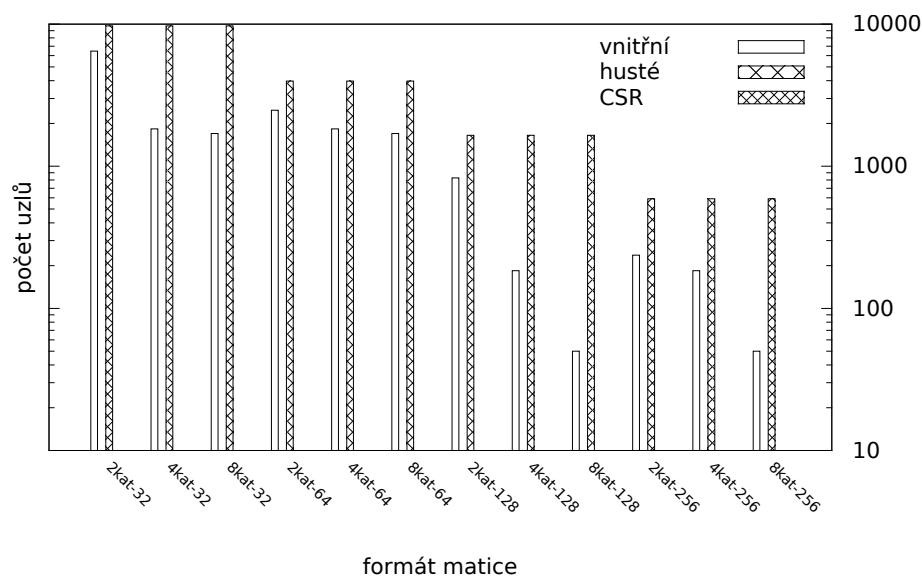
Graf 5.4 ukazuje v bytech přesnou velikost datových struktur matic v porovnání proti uložení matice v hustém formátu. Pro každý prvek bylo potřeba 8B dat, respektive prvky byly uloženy v proměnné typu double.



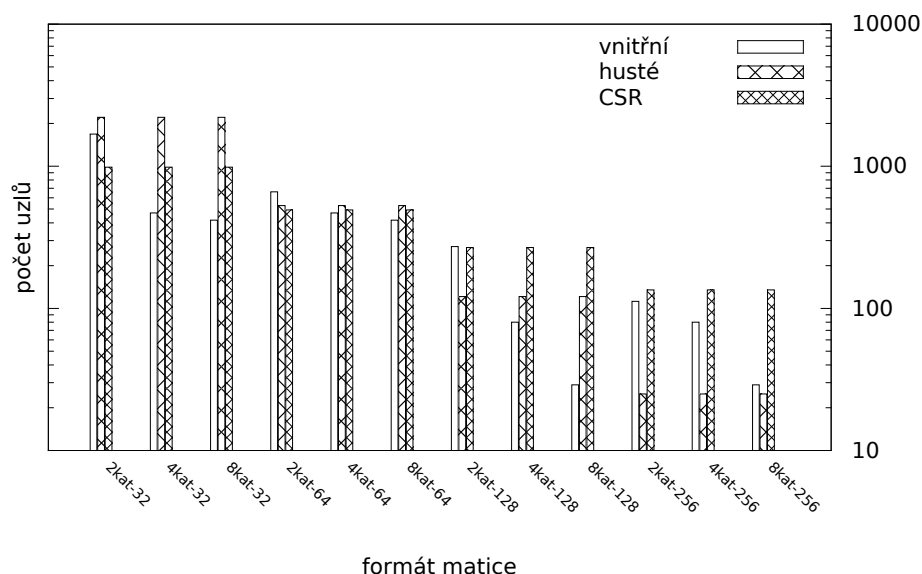
Obrázek 5.4: Velikost datových struktur

5.5.3 Počet uzlů v matici KAT

Grafy 5.5 5.6 5.7 5.8 5.9 5.10 ukazují počet uzlů ve stromě a jejich druh.



Obrázek 5.5: Počet uzlů v matici KAT pro uložení matice EX6



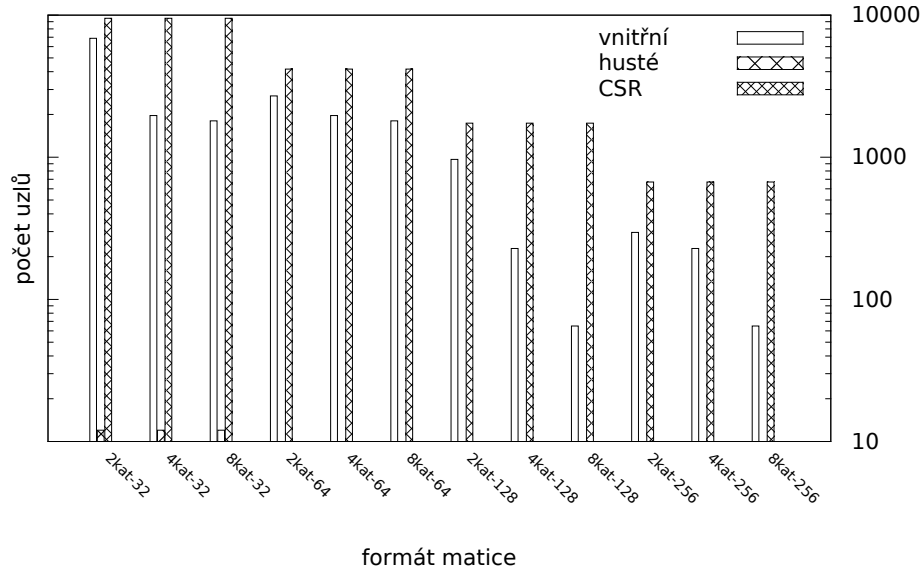
Obrázek 5.6: Počet uzlů v matici KAT pro uložení matice exdata-1

5.6 Porovnání výkonnosti s předpoklady

Protože při násobení řídkých matic velmi záleží na rozložení jejich prvků, tak abychom porovnali teoretické předpoklady s naměřenými hodnotami, provedeme měření na hustých maticích o velikosti 256, 512, 1024 a 2048. Protože jsme použili pouze základní algoritmy násobení, časová složitost by měla být $\mathcal{O}(n^3)$ pro násobení matice s maticí a $\mathcal{O}(n^2)$ pro násobení matice s vektorem. Paměťová složitost by měla být nejmenší u CSR a BSR, kde ukládáme u prvků pouze jejich hodnoty a sloupce. U formátu KAT by měla být paměťová složitost velká z důvodu ukládání struktury stromu.

Z měření vyplývá, že se teoretické předpoklady střetly s naměřenými hodnotami. Graf pro násobení matice maticí 5.11 ukazuje kubickou závislost mezi velikostí matice a časem. Graf pro násobení matice vektorem 5.12 ukazuje kvadratickou závislost. Velikost datových struktur ukazuje graf 5.13. Lze z něj vyčíst, že COO potřebuje více paměti než CSR, protože ukládá o každém prvku i řádek, kdežto CSR si pro každý řádek pamatuje rozsah prvků a tedy pro každý prvek lze dopočítat, v jakém je řádku. Formát BSR má nejmenší datovou velikost, jelikož bloky ukládá jako husté matice. Formát KAT ukazuje, že uložení stromu sice zvýší celkovou velikost matice, ale ne tak výrazně, aby byl formát nepoužitelný.

Tabulka 5.6 shrnuje časové a tabulka 5.6 paměťové složitosti.



Obrázek 5.7: Počet uzlů v matici KAT pro uložení matice fp

Formát	Časová složitost MMM	Časová složitost MVM
COO	jako CSR	$\mathcal{O}(nnz)$
CSR	$\mathcal{O}(\sum_{i=1}^n nnzr_{A,i} \cdot nnz_{CB,i})$ 2.7	$\mathcal{O}(nnz)$
BSR	$\mathcal{O}(blks \cdot sm_size^3)$	$\mathcal{O}(blks \cdot sm_size^2)$
KAT	$\mathcal{O}(\frac{n^2}{sm_size} \cdot (height + sm_size^3))$	$\mathcal{O}(\frac{n^2}{sm_size} \cdot (height + sm_size^2))$

Tabulka 5.1: Přehled časových složitostí násobení matic

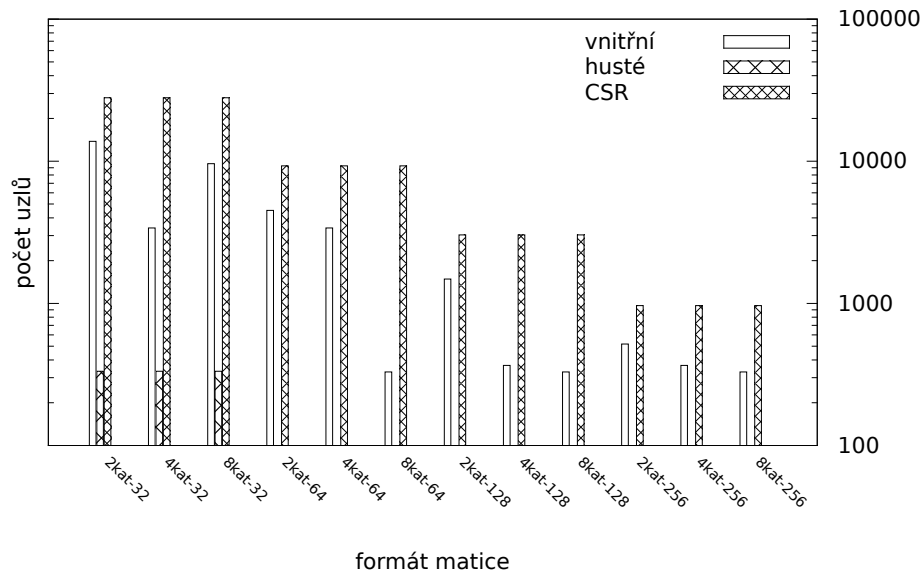
Formát	Paměťová složitost
COO	$\mathcal{O}(nnz \cdot 3)$
CSR	$\mathcal{O}(nnz \cdot 2 + m + 1)$
BSR	$\mathcal{O}(blks \cdot sm_size^2)$
KAT	$\mathcal{O}(k^{height-2} - k - \frac{1}{k})$

Tabulka 5.2: Přehled paměťových složitostí násobení matic

5.7 Zhodnocení výsledků

Matice ve formátu BSR dosahovaly velmi špatných výsledků, protože jejich bloky jsou pouze husté. V případě větších bloků s menším počtem nulových prvků tak dochází jak k ukládání těchto nulových prvků zbytečně, tak se tím zpomalí výpočet.

Největším zlepšením byla matice KAT pro matici exdata-1 [15], kde se



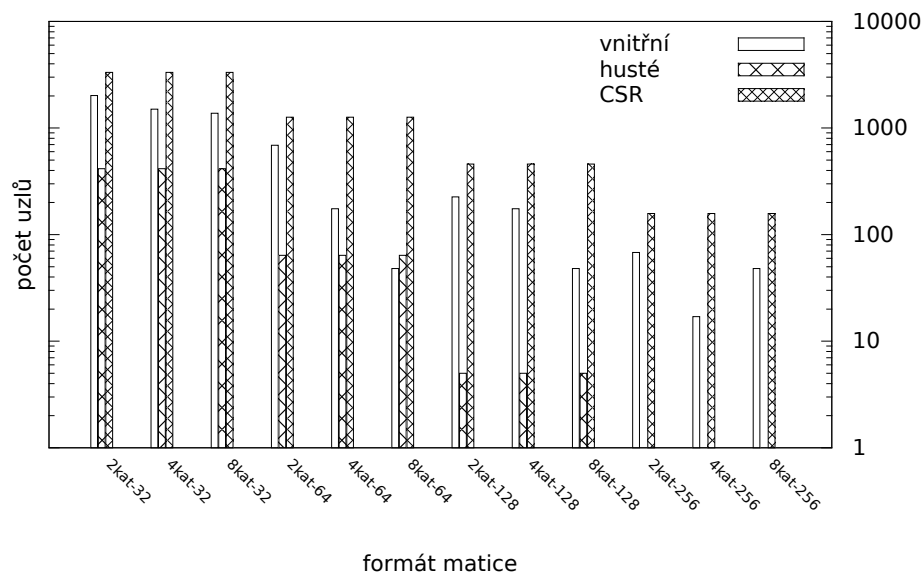
Obrázek 5.8: Počet uzlů v matici KAT pro uložení matice gupta3

výpočet zrychlil o více než dvojnásobek. Důvodem bylo efektivní rozdělení matice na husté a řídké části. I paměťová velikost byla srovnatelná s formátem CSR.

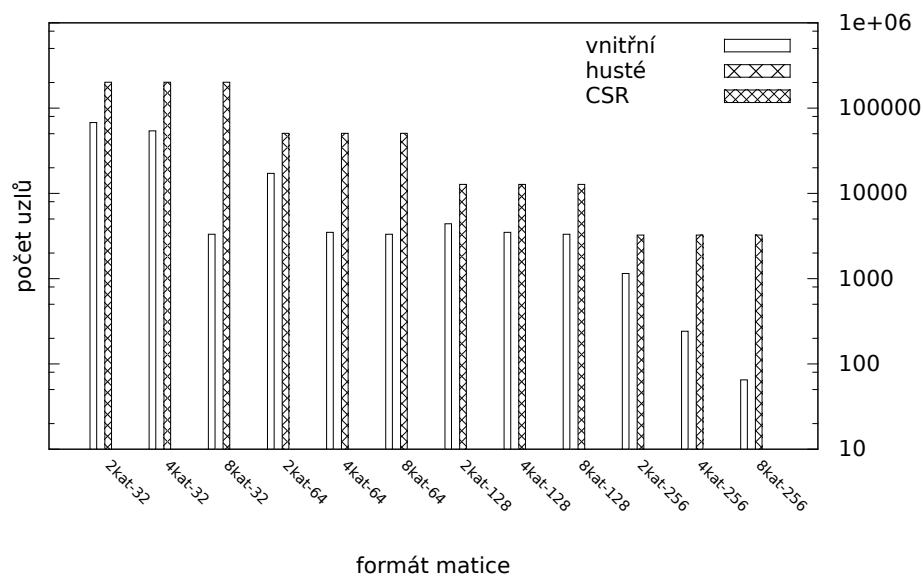
U matic, kde není nějaký vzor ve formě bloků, například u matice EX6 [17], byl formát CSR rychlejší ve všech případech. Za povšimnutí také stojí to, že v základním nastavení programu nebyl v matici KAT zvolen pro tuto matici ani jeden hustý list.

Také ani jeden hustý list nebyl zvolen pro matici human_gene2 [14]. Vzhledem k její hustotě a velikosti ale pro velké CSR bloky 128 a 256 díky časové a prostorové lokalitě [33] byly výsledky lepší než v obyčejném CSR.

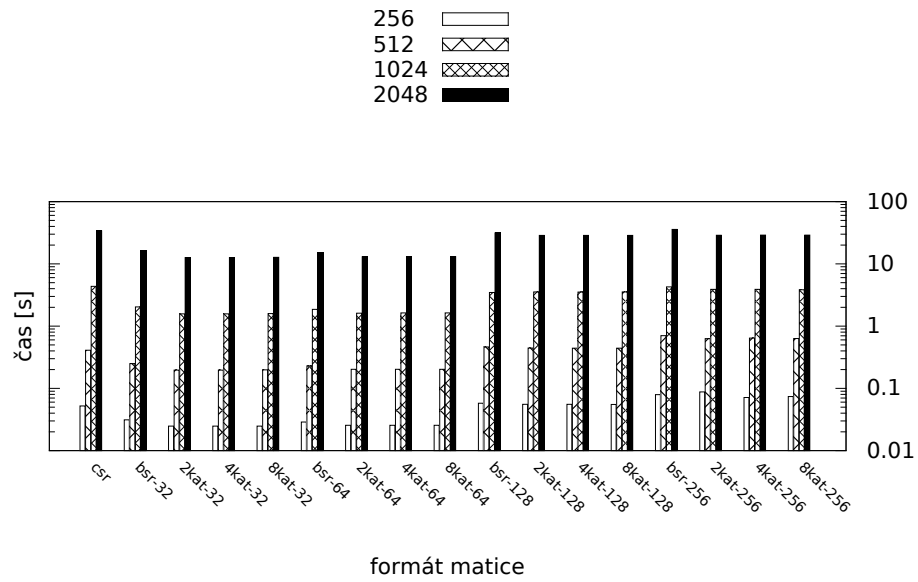
Přestože se s větším k u KAT matic výrazně snížil počet vnitřních uzlů, tak se celková datová velikost matice výrazně zvýšila. Důvodem je velikost uzlu, která obsahuje k^2 ukazatelů na syny.



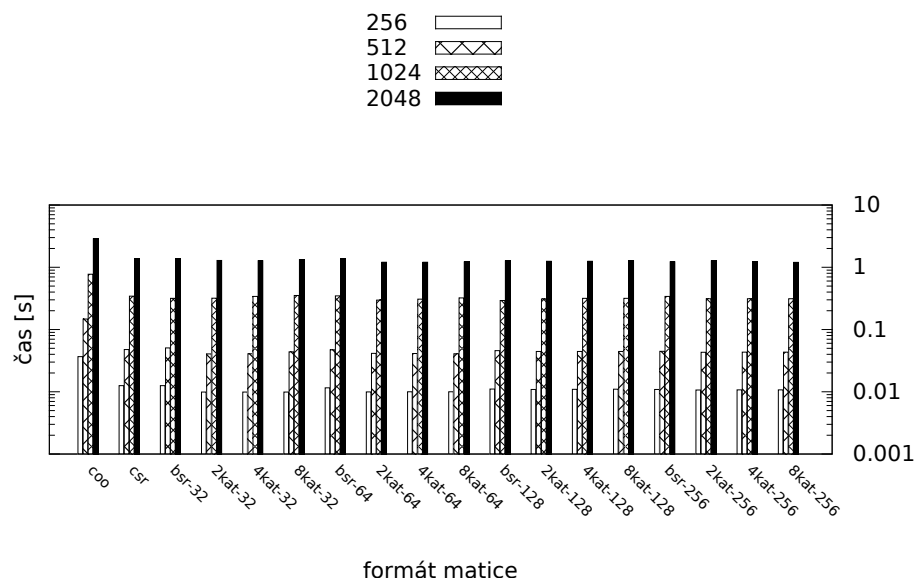
Obrázek 5.9: Počet uzlů v matici KAT pro uložení matice heart1



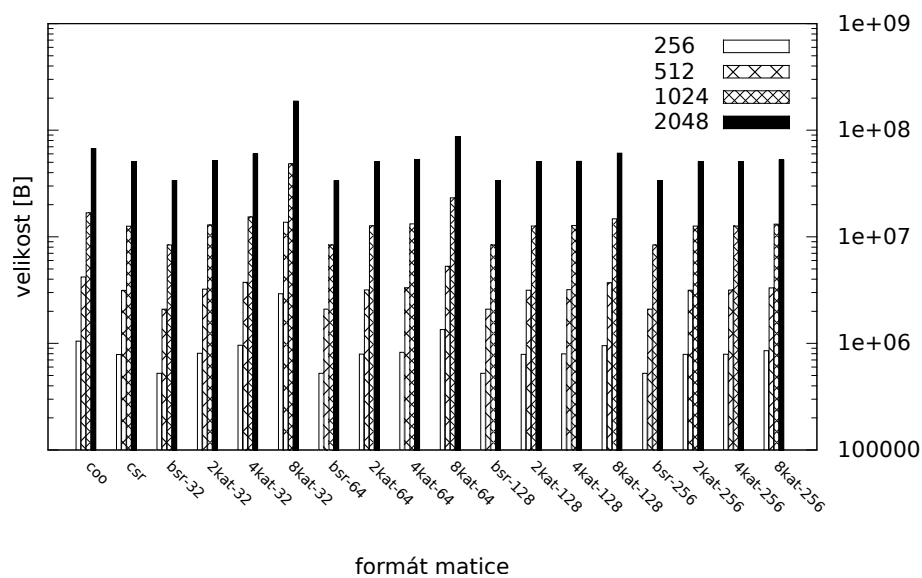
Obrázek 5.10: Počet uzlů v matici KAT pro uložení matice human-gene2



Obrázek 5.11: Čas násobení husté matice s hustou maticí v řídkém formátu



Obrázek 5.12: Čas násobení husté matice s vektorem v řídkém formátu



Obrázek 5.13: Datová velikost hustých matic v řádkých formátech

Závěr

Násobení matic je jednoduchá početní úloha, se kterou se setkáme již na střední škole. Pro násobení velkých řídkých matic je pro vyšší výkonnost potřeba sofistikovanějších algoritmů a rozdělení matic na části.

V této práci jsme popsali základní, ale i pokročilé algoritmy pro násobení matic. Popsali jsme formáty pro uložení řídkých matic COO, CSR, BSR, Quadtree a jeho modifikaci snížením stromu, pojmenovaným KAT. Tyto formáty a algoritmy pro násobení v těchto formátech jsme implementovali a změřili jsme jejich výkonnost na vybraných maticích z praxe.

Měření ukázalo, že násobení matic ve formátu Quadtree a jeho modifikaci do formátu KAT je vhodným formátem, pokud máme dostatek operační paměti a matice buď obsahuje hustší bloky nebo tolik prvků, že začne využívat cache více než formát CSR.

Další práce

Paralelizace

Díky paralelizaci by se mohlo u formátu KAT zvýšení rychlosti oproti CSR ještě umocnit, protože by jednotlivá vlákna řešila jednotlivé bloky a využívala by tak cache efektivněji než CSR.

Strassenův algoritmus

CSR není vhodný formát pro Strassenův algoritmus. Tento algoritmus by bylo možné implementovat do formátu Quadtree. Experimentem by mohla být implementace rozbaleného Strassenova algoritmu, tedy místo násobení matic o velikosti $n = 2$ by se násobily rovnou matice o velikosti například $n = 4$. Změněním pořadí operací by mohlo dojít k zefektivnění tohoto algoritmu.

Řídké vnitřní uzly

V naší implementaci byly uzly matice KAT implementovány jako husté matice ukazatelů na syny. Pro velká k by bylo vhodnější, kdyby i uzly tohoto stromu byly řídké matice.

Další typy listů

Dělit listy na husté a řídké se ukázalo jako správné řešení. Přítěží bylo, pokud blok obsahoval pouze diagonálu, kde formát CSR není příliš efektivní. Řešením by bylo podporovat listy, které budou ukládat prvky v diagonálách.

Literatura

- [1] *A Quadtree Based Storage Format and Effective Multiplication Algorithm for Sparse Matrix [online]*. [cit. 2014-04-27]. Dostupné z: <http://quardtree.sourceforge.net/>
- [2] Free Software Foundation, Inc.: GCC, the GNU Compiler Collection. 2014. Dostupné z: <http://gcc.gnu.org/>
- [3] Free Software Foundation, Inc.: GDB: The GNU Project Debugger. 2014. Dostupné z: <http://www.gnu.org/software/gdb/documentation/>
- [4] Free Software Foundation, Inc.: GNU Make. 2014. Dostupné z: <http://www.gnu.org/software/make/>
- [5] ValgrindTM Developers: Valgrind. 2014. Dostupné z: <http://valgrind.org/>
- [6] Bilmes, J.; Asanović, K.; Demmel, J.; aj.: PHiPAC: A Portable, High-Performance, ANSI C Coding Methodology and its application to Matrix Multiply. LAPACK working note 111, University of Tennessee, 1996.
- [7] Boisvert, R. F.; Pozo, R.; Remington, K.; aj.: Matrix Market: A Web Resource for Test Matrix Collections. In *Proceedings of the IFIP TC2/WG2.5 Working Conference on Quality of Numerical Software: Assessment and Enhancement*, London, UK, UK: Chapman & Hall, Ltd., 1997, ISBN 0-412-80530-8, s. 125–137. Dostupné z: <http://dl.acm.org/citation.cfm?id=265834.265854>
- [8] Borsic, A.; Halter, R.; Wan, Y.; aj.: Sensitivity study and optimization of a 3D electric impedance tomography prostate probe. *Physiological Measurement*, ročník 30, č. 6, 2009: str. S1. Dostupné z: <http://stacks.iop.org/0967-3334/30/i=6/a=S01>

- [9] Cenk, M.; Hasan, M. A.: On the Arithmetic Complexity of Strassen-Like Matrix Multiplications. *IACR Cryptology ePrint Archive*, ročník 2013, 2013: str. 107.
- [10] Coppersmith, D.; Winograd, S.: Matrix Multiplication via Arithmetic Progressions. *J. Symb. Comput.*, ročník 9, č. 3, 1990: s. 251–280.
- [11] Cormen, T. H.; Stein, C.; Rivest, R. L.; aj.: *Introduction to Algorithms*. McGraw-Hill Higher Education, druhé vydání, 2001, ISBN 0070131511.
- [12] Cyp: File:Strassen_algorithm.svg. 2008. Dostupné z: http://en.wikipedia.org/wiki/File:Strassen_algorithm.svg
- [13] Davis, T.: SPARSE MATRIX ALGORITHMS AND SOFTWARE. 2014. Dostupné z: <http://www.cise.ufl.edu/~davis/research.html>
- [14] Davis, T.; Hu, Y.: Matrix: Belcastro/human_gene2. 2014. Dostupné z: http://www.cise.ufl.edu/research/sparse/matrices/Belcastro/human_gene2.html
- [15] Davis, T.; Hu, Y.: Matrix: GHS_indef/exdata_1. 2014. Dostupné z: http://www.cise.ufl.edu/research/sparse/matrices/GHS_indef/exdata_1.html
- [16] Davis, T.; Hu, Y.: Matrix: Gupta/gupta3. 2014. Dostupné z: <http://www.cise.ufl.edu/research/sparse/matrices/Gupta/gupta3.html>
- [17] Davis, T.; Hu, Y.: Matrix: JGD_SPG/EX6. 2014. Dostupné z: http://www.cise.ufl.edu/research/sparse/matrices/JGD_SPG/EX6.html
- [18] Davis, T.; Hu, Y.: Matrix: MKS/fp. 2014. Dostupné z: <http://www.cise.ufl.edu/research/sparse/matrices/MKS/fp.html>
- [19] Davis, T.; Hu, Y.: Matrix: Norris/heart1. 2014. Dostupné z: <http://www.cise.ufl.edu/research/sparse/matrices/Norris/heart1.html>
- [20] Davis, T.; Hu, Y.: Visualizing Sparse Matrices. 2014. Dostupné z: <http://www.cise.ufl.edu/research/sparse/matrices/synopsis/>
- [21] Davis, T. A.; Hu, Y.: The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Softw.*, ročník 38, č. 1, Prosinec 2011: s. 1:1–1:25, ISSN 0098-3500, doi:10.1145/2049662.2049663. Dostupné z: <http://www.cise.ufl.edu/research/sparse/matrices>
- [22] Dawes, B.; Abrahams, D.; Rivera, R.: Boost. 2007. Dostupné z: <http://www.boost.org/>

-
- [23] El-Kurdi, Y.; Gross, W.; Giannacopoulos, D.: Sparse Matrix-Vector Multiplication for Finite Element Method Matrices on FPGAs. In *Field-Programmable Custom Computing Machines, 2006. FCCM '06. 14th Annual IEEE Symposium on*, April 2006, s. 293–294, doi:10.1109/FCCM.2006.65.
- [24] Gates, A. Q.; Kreinovich, V.: Strassen's Algorithm Made (Somewhat) More Natural: A Pedagogical Remark. *Bulletin of the European Association for Theoretical Computer Science (EATCS)*, ročník 73, 2001: s. 142–145.
- [25] Grimes, R.: ORSIRR 1: Oil reservoir simulation - generated problems oil reservoir simulation for 21x21x5 irregular grid. 2014. Dostupné z: http://math.nist.gov/MatrixMarket/data/Harwell-Boeing/oilgen/orsirr_1.html
- [26] Intel: Intel® Math Kernel Library 11.0.5 Reference Manual. 2014. Dostupné z: https://software.intel.com/sites/products/documentation/doclib/mkl_sa/11/mklman/GUID-9FCEB1C4-670D-4738-81D2-F378013412B0.htm
- [27] ISO: ISO C Standard 1999. Technická zpráva, 1999, iSO/IEC 9899:1999 draft. Dostupné z: <http://www.open-std.org/jtc1/sc22/wg14/www/docs/n1124.pdf>
- [28] Karabela, T.: *Rozšíření implementace formátu kvadrantového stromu*. Bakalářská práce, České vysoké učení technické v Praze, Fakulta informačních technologií, 2013.
- [29] Lubuntu: Lubuntu. 2014. Dostupné z: <http://lubuntu.net/>
- [30] Python Software Foundation: Python. 2014. Dostupné z: <https://www.python.org/>
- [31] Schreiner, A.: *Objektorientierte Programmierung mit ANSI C*. Hanser, 1994, ISBN 9783446174269. Dostupné z: <http://books.google.cz/books?id=sgA2AAAACAAJ>
- [32] Šimeček, I.: Sparse Matrix Computations with Quadrees. In *Seminar on Numerical Analysis*, Liberec: Technical University, 2008, ISBN 978-80-7372-298-2, s. 122–124.
- [33] Snir, M.; Yu, J.: On the Theory of Spatial and Temporal Locality. 2005.
- [34] Stothers, A. J.: *On the Complexity of Matrix Multiplication*. diploma thesis, University of Edinburgh, 2010.

- [35] Strassen, V.: Gaussian Elimination is not Optimal. *Numerische Mathematik*, ročník 13, č. 4, Prosinec 1967: s. 354–355.
- [36] of Tennessee, U.: PHiPAC. 2014. Dostupné z: <http://www1.icsi.berkeley.edu/~bilmes/hipac/>
- [37] The Eclipse Foundation: Eclipse CDT (C/C++ Development Tooling). 2014. Dostupné z: <http://www.eclipse.org/cdt/>
- [38] The Scipy community: scipy.sparse.bsr_matrix. 2014. Dostupné z: http://docs.scipy.org/doc/scipy-0.13.0/reference/generated/scipy.sparse.bsr_matrix.html
- [39] The Scipy community: Sparse matrices (scipy.sparse). 2014. Dostupné z: <http://docs.scipy.org/doc/scipy/reference/sparse.html>
- [40] Tvrdík, P.; Šimeček, I.: A new diagonal blocking format and model of cache behavior for sparse matrices.
- [41] University of California at Berkeley: Sparsity. 2014. Dostupné z: <http://www.cs.berkeley.edu/~yelick/sparsity/>
- [42] University of Tennessee: Automatically Tuned Linear Algebra Software. 2014. Dostupné z: <http://acts.nersc.gov/atlas/>
- [43] Walter, J.; Koch, M.; Winkler, G.; aj.: Basic Linear Algebra Library. 2010. Dostupné z: http://www.boost.org/doc/libs/1_55_0/libs/numeric/ublas/doc/index.htm
- [44] Williams, V. V.: Breaking the Coppersmith-Winograd barrier. 2011.
- [45] Williams, V. V.: Multiplying matrices faster than coppersmith-winograd. In *STOC*, editace H. J. Karloff; T. Pitassi, ACM, 2012, ISBN 978-1-4503-1245-5, s. 887–898.
- [46] Wolf, M. E.; Lam, M. S.: A Data Locality Optimizing Algorithm. In *Proceedings of the ACM SIGPLAN 1991 Conference on Programming Language Design and Implementation*, PLDI '91, New York, NY, USA: ACM, 1991, ISBN 0-89791-428-7, s. 30–44, doi:10.1145/113445.113449. Dostupné z: <http://doi.acm.org/10.1145/113445.113449>
- [47] Wolfram: Wolfram Mathematica 9. 2014. Dostupné z: <http://www.wolfram.com/mathematica/>
- [48] Wolfram: Wolfram Mathematica 9 Documentation Center: MTX. 2014. Dostupné z: <http://reference.wolfram.com/mathematica/ref/format/MTX.html>
- [49] Yuster, R.; Zwick, U.: Fast sparse matrix multiplication. *ACM Transactions on Algorithms*, ročník 1, č. 1, 2005: s. 2–13.

Přílohy

A.1 Překlad programu

Překlad řídí nástroj GNU Make [4]. Pro kompilaci pouze hlavní části programu je potřeba spustit příkaz `make`. Pro překlad i testů a generátoru matic je nutné zavolat `make tests`. Pro nastavení kompilace se předávají proměnné z příkazového řádku: `make PROMENNA=HODNOTA tests`. Zde je seznam nastavitelných hodnot:

- `KAT_DENSE_TRESHOLD` kolik procent z maximálního počtu prvků musí list stromu matice KAT obsahovat, aby s ním bylo zacházeno jako s hustým. Například pro `make KAT_DENSE_TRESHOLD=0.5 tests` a velikost `sm_size = 10` bude potřeba alespoň 50 prvků, aby byl list hustý. Defaultní hodnota je 0.75.
- `KAT_N` tímto nastavíme k v KAT matici. Například pro `make KAT_N=8` bude $k = 64$. Defaultní hodnota je 2, tedy Quadtree.
- `PRECISION` Nastavení přesnosti výpočtu, tedy v jaké proměnné budou uloženy prvky.
 - 1 - float (4B)
 - 2 - double (8B)
 - 3 - long double (alespoň 8B)
- `DIFF_TRESHOLD` Tolerance odchylky u testů při porovnávání dvou hodnot. Defaultní hodnota je 0.001.
- `DEBUG` Pokud je tato hodnota nastavena, program je zkompileován pomocí přepínačů `-Og -ggdb`, tedy zapnou se optimalizace, které neohroží ladění a přidají se do binárního výstupu debugovací symboly.

Díky přísným závislostem se pokaždé překládají všechny zdrojové kódy. Je tomu tak z důvodu, aby se vždy promítla tato nastavení do překladu.

A.2 Překlad textu práce

Překlad práce do formátu pdf se provede příkazem `make thesis`. Potřebné balíčky pro překlad jsou `latexmk`, `texlive-lang-czechslovak`, `texlive-generic-extra`, `texlive-science`.

A.3 Práce s programem

Binární soubor se jmenuje `main` a má následující přepínače:

- `-f [coo|csr|bsr|kat]` Volba formátu.
- `-a <matice .mtx>` Nastavení matice A.
- `-b <matice .mtx>` Nastavení matice B.
- `-o [<soubor>|stdout]` Nastavení výstupu výsledné matice. Buď do souboru, nebo na standardní výstup.
- `-s <velikost bloku>` Pro matice BSR a KAT nastavení `sm_size`.
- `-v` Vypíše informace o výpočtu. Hlavně dobu násobení a velikost matice v bytech.
- `-V` Matice B je vektor.

Pokud je zadána pouze matice A, program vynásobí matici A s maticí A.

A.4 Práce s formátem MatrixMarket

Náš program pracuje s formátem `.mtx` 3.1.1. Podporujeme více druhů tohoto formátu. Nesymetrický s banerem `%%MatrixMarket matrix coordinate real symmetric` a symetrický s banerem `%%MatrixMarket matrix coordinate real general`. Místo reálných čísel podporujeme i druh `pattern`, kdy nulové prvky nabývají pouze hodnoty jedna.

Přestože náš program neumí pracovat s komprimovanými soubory, v prostředí unixového systému předáváme programu pojmenovanou rouru, do které komprimovanou matici rozbalujeme příkazem: `<(gzip -cd matrix.mtx.gz)`

A.5 Generátor řídkých matic

Generátory z MatrixMarketu běží v internetovém prohlížeči a v jazyce Java. Protože takto není jednoduché matice generovat ve skriptu, abychom při distribuci našeho programu nemuseli přikládat velké testovací matice, implementovali jsme jednoduchý generátor řídkých matic. Parametry předáváme programu informace o výsledné matici a seznam objektů, tedy buď `podmatic` nebo

diagonál, které mají být do matice zahrnuty, jak ukazuje algoritmus 17. Manuál ke generátoru je možné vypsát zavoláním `./tests/bin/matrix_generator -h`.

Algorithm 17 Generování řídkých matic

```
1: procedure SPARSEMATRIXGENERATOR(file, width, height, ItemList)
2:   MtxWrapper  $\leftarrow$  InitMtxWrapper();
3:   MtxWrapper.PositionVector  $\leftarrow$  InitVector();
4:   for all Item  $\in$  ItemList do
5:     MtxWrapper.addItem(Item.y, Item.x, Item.properties);
6:     if Item.type == Mirrored then
7:       MtxWrapper.addItem(Item.x, Item.y, Item.properties);
8:     end if
9:   end for
10:  MtxWrapper.PositionVector.sort();
11:  MtxWrapper.PositionVector.removeDuplicates();
12:  MtxWrapper.write(file);
13: end procedure
```

Generátor řídkých matic byl implementován v jednom souboru. Lze ho spouštět s následujícími parametry:

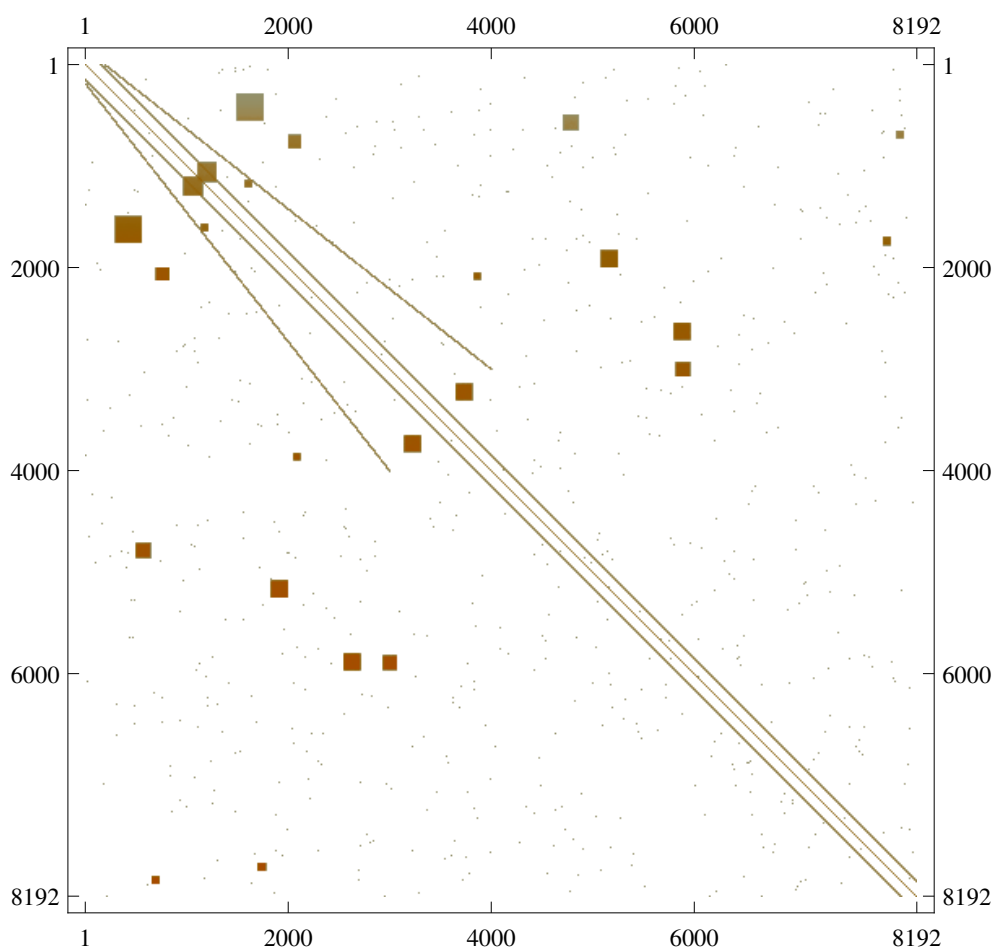
- `-c` matice nebude obsahovat hlavní diagonálu
- `-H <celé číslo>` výška matice
- `-i <typ,a,b,c,d,...>` seznam objektů, které se do matice přidají
 - `diagonal,ay,ax,by,bx,sparsity` prvky v přímce od bodu `[ax,ay]` do bodu `[bx,by]` s řídkostí `sparsity`
 - `block,ay,ax,by,bx,sparsity` blok prvků v obdelníku ohraničeného body `[ax,ay]` a `[bx,by]` s řídkostí `sparsity`
 - `mblock, rblocks, diablocks, mdiagonal` varianty bloků a diagonál. `r` znamená náhodný a `m` zrcadlový
- `-n <celé číslo>` velikost matice
- `-o <soubor>` cílový soubor (lze použít i `stdout`)
- `-s <desetinné číslo>` řídkost matice (`sparsity`)
- `-S <desetinné číslo>` startovací číslo
- `-W <celé číslo>` šířka matice
- `-h` zobraz nápovědu

A. PŘÍLOHY

- `-o <soubor>` cílový soubor (lze použít i `stdout`)
- `-v` vypisuj průběh generování (`verbose`)

Například A.1 byla vygenerována následujícím příkazem:

```
./tests/bin/matrix_generator -n 8192 -s 0.00001  
-i mdiagonal,150,0,8192,8042,0.95,  
mdiagonal,200,10,4000,3000,0.75,  
mblockwh,300,1500,256,256,0.95,  
mblockwh,700,2000,128,128,0.95,  
mrblocks,10,128,64,64,0.75 -o tmpmatrix2.mtx
```



Obrázek A.1: Matice vygenerovaná generátorem

A.6 Spuštění testů

Po překladu `make tests` je hlavní testovací spustitelný soubor `./tests/bin/test_mat_vec`.

Seznam použitých zkratk

COO	Coordinate
BSR	Block sparse row
CSC	Column sparse column
CSR	Column sparse row
KAT	k-ary tree
GNU	GNU is not UNIX
GUI	Graphical User Interface
MMM	Matrix-Matrix Mulitplication
MVM	Matrix-Vector Mulitplication

Seznam obrázků

1.1	3D neorientovaný graf ve tvaru helikoptéry, jeho reprezentace v řídké matici a výsledek simulace	4
-----	--	---

2.1	Strassenův algoritmus [12]	11
2.2	Ukázka numerické stability Strassenova algoritmu	14
2.3	Matice orsirr_1 před vynásobením	16
2.4	Matice orsirr_1 po vynásobením sama se sebou	17
3.1	Matice uložená ve formátu CSR	22
3.2	Matice uložená ve formátu BSR	24
3.3	Rozdělení matice	26
3.4	Rozdělení matice	27
3.5	Rozdělení matice	27
5.1	UML diagram tříd programu	34
5.2	Zrychlení výpočtu matice · matice oproti formátu CSR	38
5.3	Zrychlení výpočtu matice · vektor oproti formátu CSR	39
5.4	Velikost datových struktur	40
5.5	Počet uzlů v matici KAT pro uložení matice EX6	40
5.6	Počet uzlů v matici KAT pro uložení matice exdata-1	41
5.7	Počet uzlů v matici KAT pro uložení matice fp	42
5.8	Počet uzlů v matici KAT pro uložení matice gupta3	43
5.9	Počet uzlů v matici KAT pro uložení matice heart1	44
5.10	Počet uzlů v matici KAT pro uložení matice human-gene2	44
5.11	Čas násobení husté matice s hustou maticí v řídkém formátu	45
5.12	Čas násobení husté matice s vektorem v řídkém formátu	45
5.13	Datová velikost hustých matic v řídkých formátech	46
A.1	Matice vygenerovaná generátorem	56

Seznam algoritmů

1	Násobení matic podle definice	7
2	Násobení transponovanou maticí	8
3	Násobení po řádcích	9
4	Rekurzivní násobení	9
5	Strassenův algoritmus	12
6	Násobení matice COO s vektorem	20
7	Násobení dvou COO matic	21
8	Násobení matice CSR s vektorem	23
9	Násobení dvou CSR matic	23
10	Násobení matice BSR s vektorem	25
11	Násobení dvou BSR matic	25
12	Násobení matice KAT s vektorem	28
13	Násobení dvou KAT matic	29
14	Vyhledání listu pro KAT matici	36
15	Násobení hustého KAT listu s CSR listem	37
16	Testování	37
17	Generování řídkých matic	55
*		

Obsah přiloženého CD

	readme.txt.....	stručný popis obsahu CD
	barevne_grafy.....	barevné grafy použité v práci
	test_matrices.....	soubor testovacích matic
	sparse-matrices.....	zdrojové kódy
	src.....	zdrojové kódy implementace
	tests.....	zdrojové kódy testů
	thesis_cz.....	zdrojová forma práce ve formátu \LaTeX
	BP_Nesrovnal_Tomas.pdf.....	text práce