

Semestrální projekt z predmetu A4M39GPU

Tema: Simulace satu (tkaniny)

Vypracoval: Michael Smutny (smutnmic@fel.cvut.cz)

1. Struktura dokumentu

Prace je rozdelena do 7 kapitol. Kapitola 1 (strana 1) popisuje strukturu dokumentu. Kapitola 2 (strana 1) obsahuje uvod, tedy co je tematem prace, proc je reseny problem zajimavy z hlediska pouzitelnosti, a nastinuje vytycene cile prace. V kapitole 3 (strana 1) je proveden teoreticky rozbor prace. Jsou zde vysvetleny metody, pomoci nichz se dany problem resi. Kapitola 4 (strana 4) se zabýva sekvencnim resenim na CPU. Je zde uveden popis trid a jejich metod. Kapitola 5 (strana 5) resi paralelni reseni na GPU. Nejprve je zde uveden postup, jak ulohu transformovat, aby se dala pocitat paralelně. Pote jsou reseny zpusoby ulozeni dat na GPU a jejich presuny mezi CPU a GPU. Na zaver kapitoly jsou nastineny i optimalizace. V kapitole 6 (strana 9) jsou uvedeny vysledky ve forme tabulky a grafu. Porovnavany jsou jak rozdily v nastaveni jednotlivych parametru, tak (coz je hlavni) zrychleni GPU verze oproti CPU verzi. Kapitola 7 (strana 11) obsahuje zaver a zhodnoceni. V priloze A (strana 13) je popsane ovladani aplikace a v priloze B (strana 14) jsou obrazky aplikace.

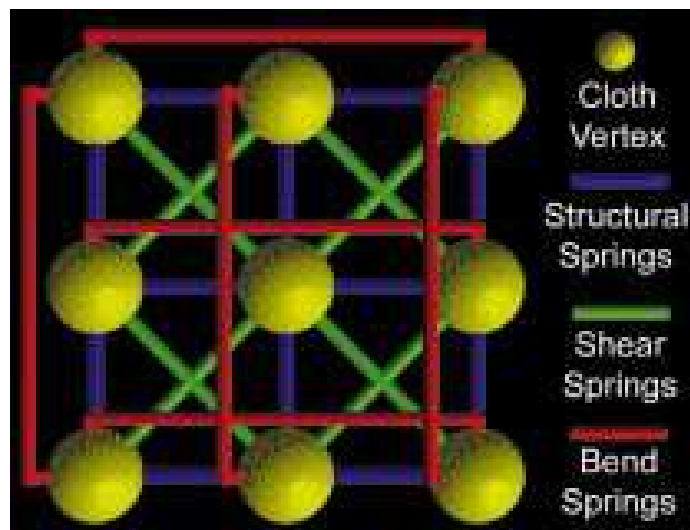
2. Uvod

Prace se zabýva simulaci tkanin pomoci casticoveho systemu. Simulace tkanin (satu) je problem, který nalezne sve uplatneni ve filmovem, hernim a modelarskem prumyslu. Jedna z nejpouzivanejsich metod pro simulaci satu, je pruzinovy model (Mass spring model), který je pouzit i zde. Cilem prace je implementovat tuto metodu na CPU a pote navrhnout reseni, které bude paralelizovatelne, a ulohu implementovat i na GPU, a to pomoci jazyka CUDA. Vysledkem GPU by melo byt znatelne zrychleni oproti CPU verzi. Prace se dale zabýva moznostmi optimalizace implementace na GPU. Soucasti prace je i mereni na ruznych scenach.

3. Teoreticky rozbor

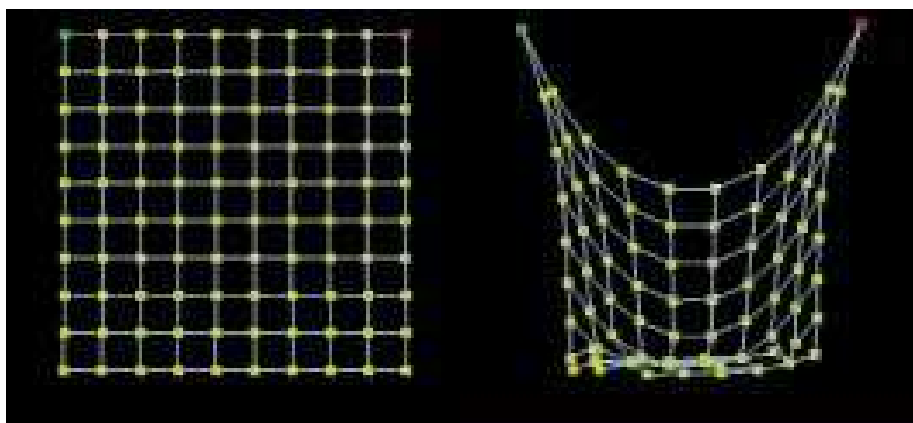
Mass spring model

Vychazel jsem ze zdroju [1], [2] a [3]. Tkanina je reprezentovana pomoci castic. Pouzil jsem ideu Mass-Spring modelu [1]. Kazda castice ma svou hmotnost (hmotnost je stejná pro vsechny castice ve scene) a vzajemne propojene castice jsou mezi sebou spojeny pruzinami. Propojeni (pruziny) jsou tri typu, jak ukazuje obrazek 1.



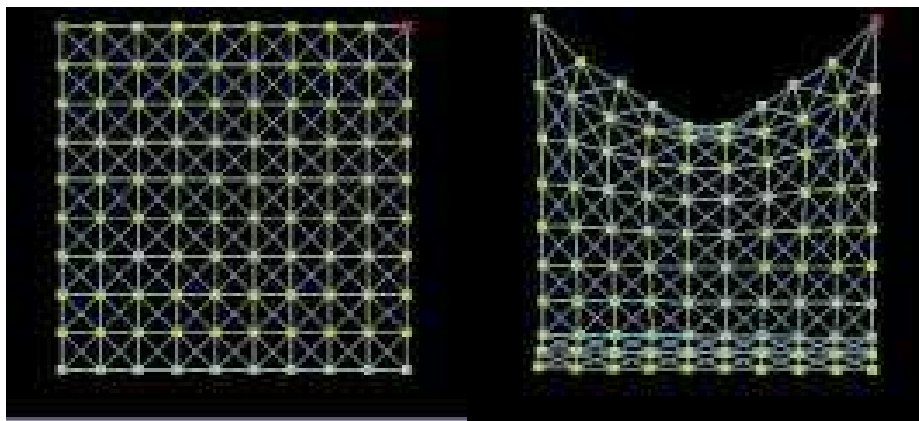
Obrazek 1: Pruzinova struktura tkaniny. (prevzato z [1])

Strukturní pružiny drží pohromadě sousední částice a jsou to tedy základní pilíře, bez kterých by se systém rozpadl. Nicméně jen strukturní pružiny nestačí, výsledek nevypadá moc realisticky. Viz obrazek 2.



Obrazek 2: Strukturní pružiny v počáteční poloze (vlevo) a při působení gravitace (vpravo). (prevzato z [1])

Z toho důvodu autor článku přidává další typ pružin, a to tzv. Shear Springs, tedy pružiny, pomocí nichž tkanina lépe drží svůj tvar. Tyto pružiny spojují „úhlopříčné sousedy“. S těmito typy pružin se daří již vystačit (obrázek 3).



Obrazek 3: Strukturní pružiny + Shear springs v počáteční poloze (vlevo) a při působení gravitace (vpravo). (prevzato z [1])

Autor však přidává další typ pružin, tzv. Bend Springs, které spojují částice ve vzdálenosti 2 (tzn. "obsouseda"). Výsledek vypadá ještě trochu lépe. Hlavním důvodem přidání těchto pružin však je to, že bez nich se při kolizi nechová tkanina zcela korektně a způsobí to, že částice utvoří takový "bezstrukturní chomáč".

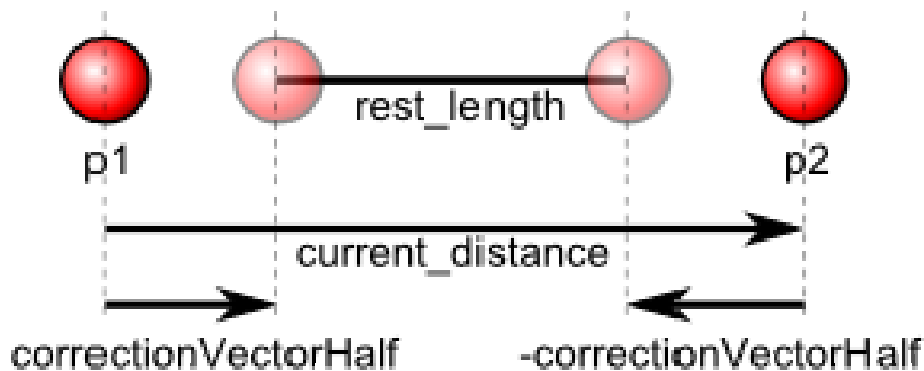
Verletova integrační metoda

Nyní, když je struktura vytvořena, je potřeba nějakým způsobem vypočítat aktuální pozice částic. Z Newtonova 2. zákona dostaneme zrychlení, které je rovno druhé časové derivaci pozice. Cílem pro stanovení pozice ze zrychlení je potřeba integrovat. Ve zdroji [2] použili tzv. Verletovu integrační metodu [3] kvůli lepší stabilitě než má Eulerova metoda a i v tomto řešení je použita stejná metoda.

```
Vec3 temp = pos;
pos = pos + (pos-old_pos) + acceleration*TIME_STEPSIZE2
old_pos = temp;
```

Do implementace je zahrzen i parametr odpor prostředí. Po jeho aplikaci se vypočítaný vektor posunutí částice o daný koeficient utlumi.

Samotný výpočet nové pozice všech částic ovšem nestaci. Výsledná pozice je ještě upravena, a to v závislosti na všech svých pružinách. Je možné určit maximální a minimální délku pružiny, kterou když pružina překročí, je délka pružiny nastavena na svou mez, nicméně lepší je vždy částice posunout do takové polohy, že jejich vzdálenost bude stejná jako originální délka pružiny (délka bez působení sil – viz obrázek 4). Jelikož každá částice je součástí až dvanácti pružin (kromě částic na hranici), tak v každém kroku "korekce vzdálenosti" se k sobě částice přiblíží na vzdálenost blízkou původní délce pružiny, než tomu bylo před tímto korekčním krokem. Čím více těchto korekčních kroků provedeme, tím stabilnější bude řešení a tím méně budou částice "poskakovat".



Obrazek 4: Korekce vzdalenosti. (prevzato z [2])

4. CPU reseni

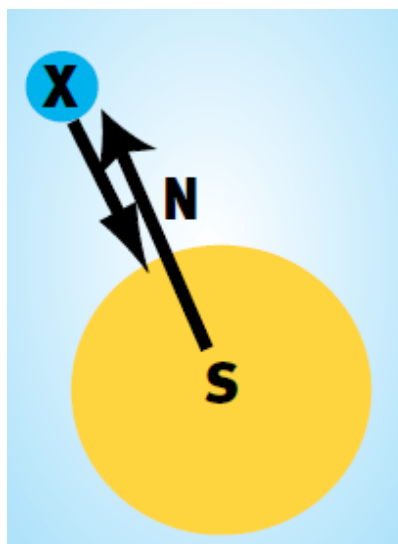
CPU verze 1

Pri seriovem reseni jsem pouzil 3 hlavni tridy, a to tridu Spring, Cloth a Particle.

Trida **Particle** predstavuje jednu castici a jejimi clenskymi daty jsou aktualni pozice, stara pozice, zrychleni, hmotnost a take informace, zda se jedna o pohyblivou castici. Některé castice museji byt ukotveny ve scene a tudiz byt nepohyblive. Trida obsahuje take metodu pro vypocet aktualni pozice castice, tedy vyse zminenou metodu integrace.

Trida **Spring** je velice jednoduchá a predstavuje pruzinu, na jejichz obou koncich je castice. Trida tedy obsahuje ukazatele na dve koncove castice a take metodu pro korekci vzdalenosti. Kvuli tomu je take nutne uchovavat v ni informaci o puvodni vzdelenosti (vzdalenosti v klidovem stavu).

Trida **Cloth** predstavuje tkaninu a jejimi clenskymi daty tedy nutne museji byt vsechny castice (`std::vector`) a pruziny (taktez `std::vector`). Trida obsahuje metody pro vykresleni satu jako trojuhelniku a i jako quadu. V konstruktoru se vyrobi vsechny castice a z prislusnych dvojic castic se vytvori pruziny. Castice v hornich rozich se museji fixovat. Trida dale obsahuje metodu `addForce()`, která vsem casticim prida silu, která na ne pusobi. Nejdulezitejsi metodou ovsem je metoda `NextState()`, která spocita aktualni pozice vseh castic, a to ve dvou krocich. Nejdrive je potreba spocitat integracni metodou nove pozice a pote korekcni metodou castice "vratit do rozumných pozic". Tato korekce se obvykle provadi vice nez jednou. Posledni metodou je metoda `computeSphereCollision()`, která zjistí, zda se nejaka castice nenachazi uprostred koule. Pokud ano, posune ji na hranici koule ve smeru normaly. Viz obrazek 5.



Obrazek 5: Posun castic ve smeru normaly na povrchu koule. (prevzato z [1])

Dalsi tridou, která se v projektu nachazi, je trida Camera. Jeji metoda Look() se stara o nastaveni kamery. Kod pro tuto tridu jsem prevzal ze sve prace z GSY.

CPU verze 2

Pri prechodu k GPU jsem pomerne znacne zmenil kod a zrusil jsem tridy Particle a Spring. Proto jsem se rozhodl, ze udelam druhou CPU verzi, která bude vychazet z GPU verze. Jen s tim rozdilem, ze misto kernelu budou "for-cykly". Zdokumentovana (doxygen) je prave tato druha verze, nicmene nazornejsi a lepe strukturovana ta verze prvni.

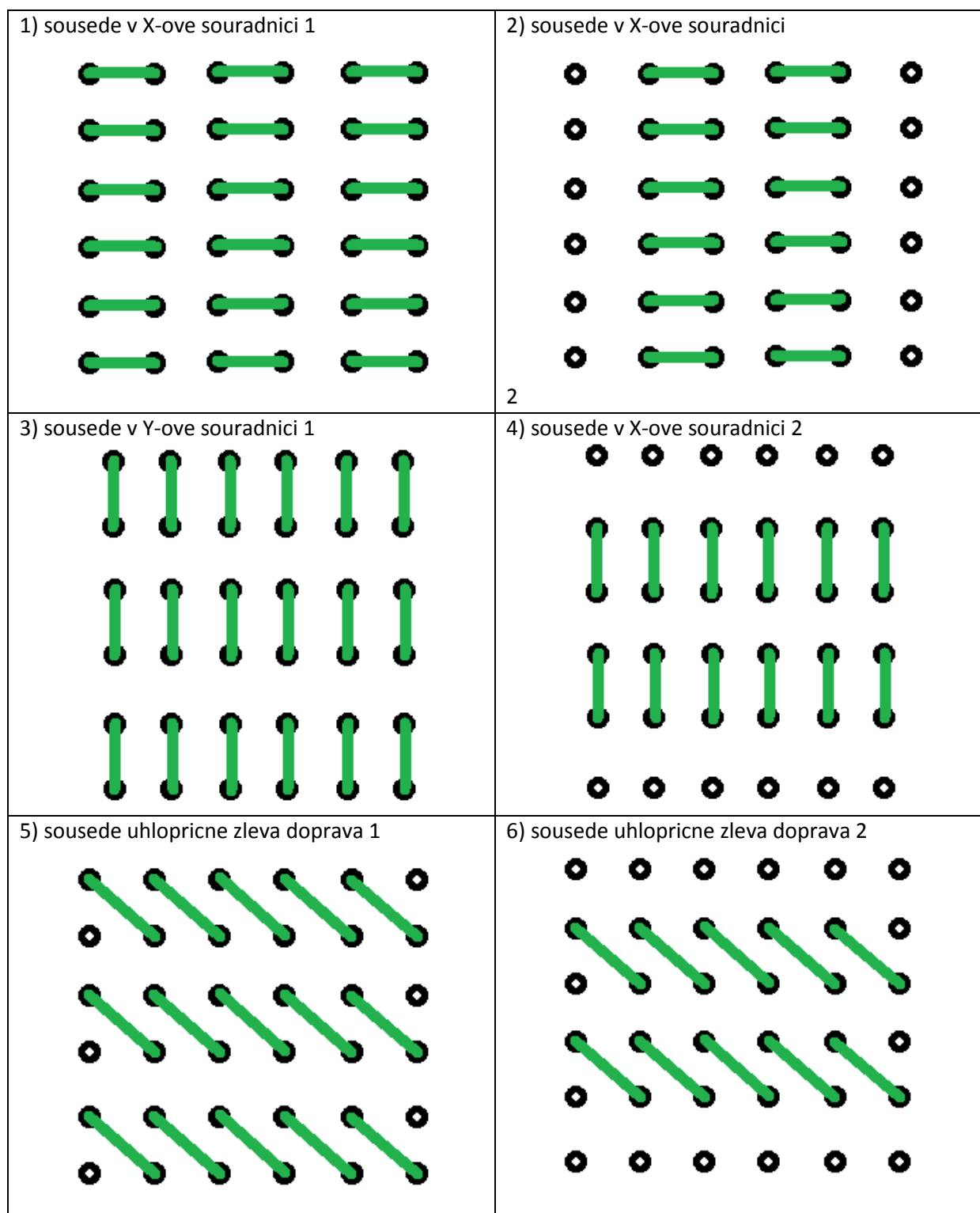
5. GPU reseni

Nejprve bylo potreba zbavit se sloziteho konceptu, kdy trida Cloth obsahovala vektor castic a take vektor pruzin, které v sobe mely ukazatele na castice. Pro CUDA reseni bylo potreba pracovat pouze s casticemi, lepe jeste pouze s jejich pozicemi (ci jinymi atributy). Odstranil jsem tedy tridy Particle a Spring a misto nich trida Cloth nyní obsahuje pouze pole atributu castic. Jmenovite je to pole aktualnich pozic castic, pole starych pozic castic, pole zrychleni castic, pole hmotnosti castic a pole boolean informaci, zda je castice fixni. Nad temito poli se daji provadet potrebne operace, které jsou shodne pro kazdou castici (->priprava pro CUDA reseni). Abych nahradil i vektor pruzin, rozhodl jsem se, ze kernel funkce si sama urci, které dve castice budou predstavovat pruzinu.

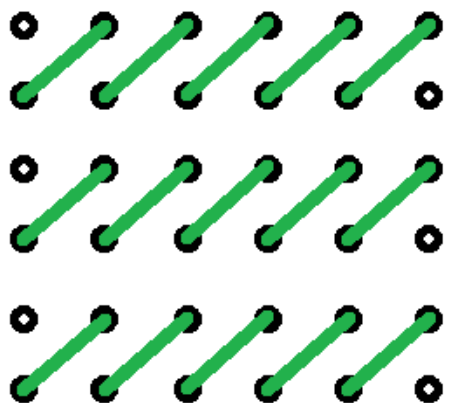
Problem prvni – vylucne sekvencni iteracni reseni korekce vzdalenosti

To, jakym zpusobem probiha korekce vzdalenosti castic, tak v teto podobě nejde paralelizovat. V kazdem kroku se pracuje s pozici castice, která byla vypoctena v kroku predeslem. Je potreba to uchopit troche z jineho pohledu. Kazda castice je soucasti dvanacti ruznych pruzin. Kazda pruzina koriguje vzdalenosti castic v korekcnim kroku prave jednou. Je jasne, ze s jednou castici se v korekcnim kroku pracuje dvanactkrat, cili zcela paralelizovat tento krok nejde. Pri vhodnem rozdeleni

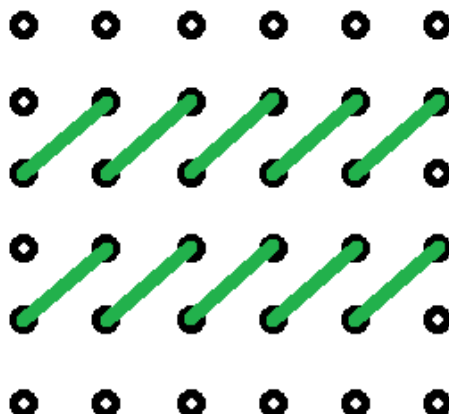
je možné v jednom okamžiku nad polem částic pracovat jako s pružinami, jejichž částice se nepřekrývají. Takovýchto konfigurací je dvanáct.



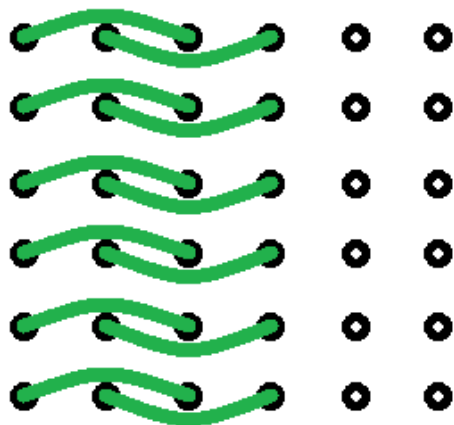
7) sousede uhlopricne zprava doleva 1



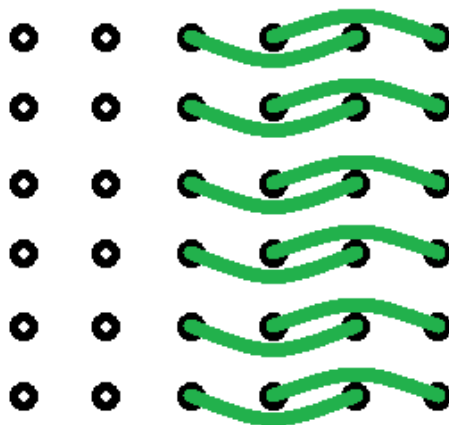
8) sousede uhlopricne zprava doleva 2



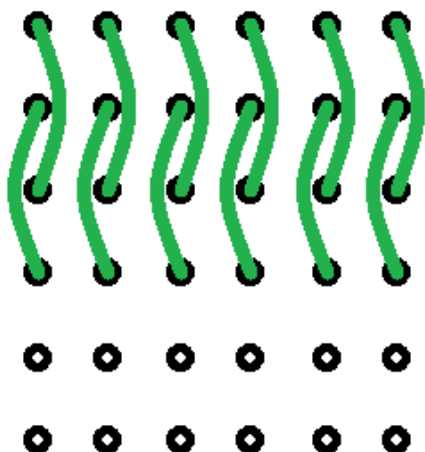
9) ob-sousedse v X-ove souradnici 1



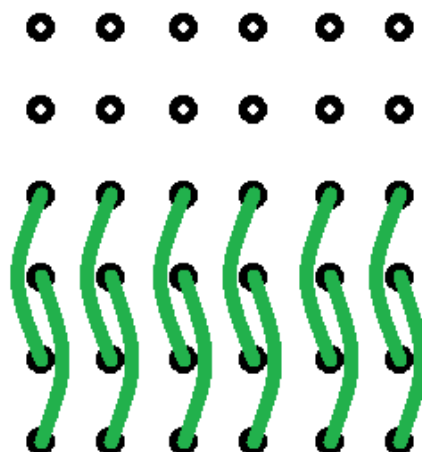
10) ob-sousedse v X-ove souradnici 2



11) ob-sousedse v Y-ove souradnici 1



12) ob-sousedse v X-ove souradnici 2



Problem druhy – vsechny vypocty presunout na GPU

Aby se zabránilo zbytečnému přesouvání dat mezi GPU a CPU, tak se ze začátku inicializují data na CPU, zkopírují se na GPU a zde již probíhají všechny výpočty. To znamená, že kromě paralelizace výpočtu pozice (integrace) a korekce vzdálenosti, jsem musel na GPU přesunout také výpočty pro přidání síly všem částicím (metoda `addForceGPU()`) a také výpočet kolize s kouli (`computeSphereCollisionGPU()`). Jediným přesunem dat tak zůstává po každém celkovém výpočtu přesun dat z GPU na CPU, aby mohla být tkanina zobrazena na výstupní zařízení (do spolupráce CUDA \leftrightarrow OpenGL jsem se nepouštěl).

Problem třetí – rozměry mřížky

Rozhodl jsem se, že počet vláken v bloku bude 16x16. Zároveň musejí být rozměry tkaniny (počet částic na dimenzi) násobky 32. Pokud uživatel zadá jiné číslo, je toto číslo zaokrouhleno na nejbližší nižší číslo 32k (kde $k=1,2,\dots$). Rozměry mřížky jsou potom dvojce, a sice:

A) $(\text{num_particles_x}/32, \text{num_particles_y}/16)$

B) $(\text{num_particles_x}/16, \text{num_particles_y}/32)$

První případ je pro konfigurace 1,2,9 a 10, druhý případ pro ostatní konfigurace. Nazornější je to na předeslých 12ti obrazech.

Jednotlivé kernely

Kernely pro výpočet pozice (integrace), pro přidání síly a pro výpočet kolize s kouli jsou velice jednoduché. Kernel je zavolan pro každou částici a požadovaná veličina je vypočtena nezávisle na ostatních částicích. Problematickejší je výpočet korekce vzdálenosti. Jak již bylo zmíněno, jedná se o 12 sekvencních kroků (12 po sobě jdoucích volání RŮZNYCH kernelů). Každý kernel dělá toto: vypočítá dva indexy do pole pozic částic a zavola funkci `CorrectPair()`, která upraví vzdálenosti. Je zde 12 kernelů pro 12 konfigurací, každý kernel vypočítá jiné indexy.

Optimalizace

Výsledné GPU řešení je zhruba 10x rychlejší než CPU řešení (viz výsledky v další kapitole). Zajímalo mě také o možnost zrychlení. Uvažoval jsem tři možnosti, maximalizovat sdružený přístup, použít texturování paměti a použít sdílenou paměť. Maximalizovat sdružený přístup bylo prakticky nemožné, protože každý kernel přečetl dvě PŘEDEM URCENÉ hodnoty z globální paměti. Sdružený přístup do paměti se tedy odvíjel od návrhu konfigurací. Ke sdruženému přístupu nedochází v konfiguracích 1,2,9 a 10. V ostatních konfiguracích k němu dochází (někde částečně). Použití texturovací paměti jsem zavrhl, protože každá hodnota z globální paměti byla přečtena (je jedno, jakým vláknem) JEN JEDNOU, tudíž zde by CACHEování (jedna z výhod texturovací paměti) nepomohlo. Použití sdílené paměti se jako perspektivní z toho důvodu, že by čtení z globální paměti mohlo být výhradně sdruženým přístupem (stejně tak zápis). Jelikož v rámci kernelu dochází pouze k jednomu čtení z globální paměti, tak by se tato výhoda proměnila spíše v nevýhodu (další zbytečná režie). Proto jsem využil rozměru mřížky jednotlivých volání kernelů a místo do 12 skupin jsem je rozdělil do dvou skupin. První skupinu tvoří konfigurace 1,2,9 a 10 a mají rozměry mřížky $(\text{num_particles_x}/32,$

num_particles_y/16). Druha skupina jsou ty ostatní konfigurace. Idea je tedy mít pouze dva kernely, nacíst do sdílené paměti potřebný kus globální paměti, nad kterým bude prováděn výpočet a jakoby provést několik původních kernelů v jednom kernelu naraz. Samozřejmě je potřeba nacíst do sdílené paměti více dat, že je potřeba (data “za hranici” bloku), takže nakonec blok prvního typu bude mít 18x16 vláken a blok druhého typu 18x18 vláken (obsahuje i uhlopříčné pružiny).

Pozn.: Podarilo se mi implementovat pouze kernel prvního typu, u kernelu druhého typu jsem nebyl schopen správně spočítat indexy při jednotlivé konfiguraci, což vedlo k nekorektnímu a nestabilnímu řešení. V aplikaci je možnost přepínání mezi řešením pomocí globální a sdílené paměti. Pokud je nastaven výpočet pomocí sdílené paměti, tak kernel prvního typu dělá výpočet konfigurací 1,2,9 a 10, zatímco pro ostatní konfigurace jsou použity příslušné kernely z řešení pomocí globální paměti. I tak je ale patrné zrychlení oproti verzi s globální paměti.

6. Výsledky

Měření proběhlo na PC s konfigurací: Core2Duo (2 x 2,9GHz), 2GB RAM, GF GTX 260 (898MB)

V aplikaci se dají nastavit různé parametry, ovšem ty, které ovlivňují rychlost výpočtu, jsou dva. Počet částic a počet iterací (počet kroku korekce vzdálenosti). V tabulkách je zaznamenán čas výpočtu 500 snímků (výsledný objekt nebyl zobrazován, aby zobrazování nezkreslovalo měření). Do tabulky je také vynesena hodnota zrychlení GPU oproti CPU verzi.

| | | | | | |
|---------------|-------|------------|-----|------------|------------|
| Pocet castic | 32x32 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 407 | 248 | 281 |
| | | ratio | 1,0 | 1,6 | 1,4 |

| | | | | | |
|---------------|-------|------------|------|------------|------------|
| Pocet castic | 64x64 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 1702 | 265 | 251 |
| | | ratio | 1,0 | 6,4 | 6,8 |

| | | | | | |
|---------------|---------|------------|------|------------|------------|
| Pocet castic | 128x128 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 7748 | 438 | 409 |
| | | ratio | 1,0 | 17,7 | 18,9 |

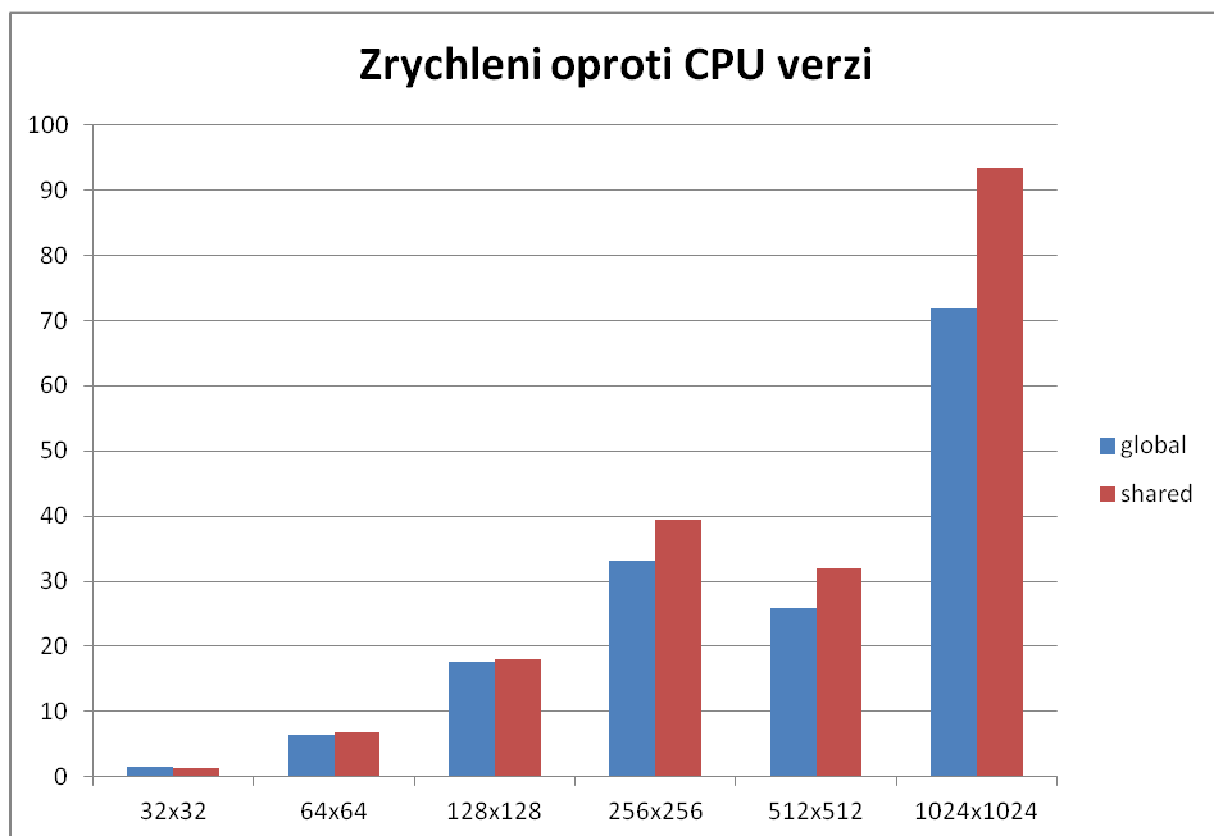
| | | | | | |
|---------------|---------|------------|-------|------------|------------|
| Pocet castic | 256x256 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 43110 | 1299 | 1093 |
| | | ratio | 1,0 | 33,2 | 39,4 |

| | | | | | |
|---------------|---------|------------|--------|------------|------------|
| Pocet castic | 512x512 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 122000 | 4717 | 3809 |
| | | ratio | 1,0 | 25,9 | 32,0 |

| | | | | | |
|---------------|-----------|------------|---------|------------|------------|
| Pocet castic | 1024x1024 | 500 snimku | CPU | GPU global | GPU shared |
| Pocet iteraci | 2 | cas [ms] | 1375000 | 19137 | 14706 |
| | | ratio | 1,0 | 71,9 | 93,5 |

| | | | | | |
|---------------|---------|------------|-------|------------|-----------|
| Pocet castic | 128x128 | 500 snimku | CPU | GPU global | GPU Sharp |
| Pocet iteraci | 20 | cas [ms] | 70999 | 4404 | 3782 |
| | | ratio | 1,0 | 16,1 | 18,8 |

| Zrychleni | | | | | | |
|--------------|-------|-------|---------|---------|---------|-----------|
| Pocet castic | 32x32 | 64x64 | 128x128 | 256x256 | 512x512 | 1024x1024 |
| GPU global | 1,6 | 6,4 | 17,7 | 33,2 | 25,9 | 71,9 |
| GPU shared | 1,4 | 6,8 | 18 | 39,4 | 32 | 93,5 |



Diskuse

Da se predpokladat, ze pocet iteraci na celkovem pomeru zrychleni (casu) nebude mit moc vliv. To se take potvrdilo (viz priklad pro 128x128 castic). Dale je dobre si povsimnout, ze zrychleni se zvetsuje se vzrustajicim pocetem castic (az na zajimavy pripad 512x512). To take odpovida predpokladum. Kyzeneho zrychleni oproti zakladni GPU verzi se dockalo take reseni za pouziti sdilene pameti (v radu jednotek az desitky procent).

7. Zaver

Implementoval jsem metodu [1], [2] pro simulaci satu na GPU pomoci jazyka CUDA. Zrychleni oproti CPU verzi se pohybuje v radu desitek. Zobrazovani je provedenou pouze velice jednoduse (wireframe nebo trojuhelniky/quady bez normal). Vice o program viz priloha A.

Zdroje

- [1] Lander J., Devil in the Blue Faceted Dress: Real-time Cloth Animation, Game Developer, 1999, stranka 17-22, <http://www.darwin3d.com/gamedev/articles/col0599.pdf>
- [2] [webova stranka]Mosegaards Cloth Simulation Coding Tutorial
<http://cg.alexandra.dk/2009/06/02/mosegaards-cloth-simulation-coding-tutorial/>
- [3] [webova stranka] Verletova integracni metoda
http://en.wikipedia.org/wiki/Verlet_integration

Priloha A – Cloth Simulation aplikace

Ovladani

| | | |
|---------|---|------------------------------------|
| W,S,A,D | - | pohyb kamery |
| Mys | - | rozhizeni |
| E | - | wireframe on/off |
| M | - | pouziti sdilene/globalni pameti |
| Q | - | kreslit trojuhelniky/quady |
| T/G | - | pocet iteraci korekcniho kroku +/- |
| Y/H | - | polomer koule +/- |
| U/J | - | rychlost koule +/- |
| I/K | - | odpor prostredi +/- |
| PG_UP | - | fullscreen |
| PG_DOWN | - | window |
| ESC | - | konec |

Spusteni

ClothSim_*.exe <cfg_file> (* CPU/GPU)

Config File (mezera nutna u rovnitka a mezi FLOATy)

| | |
|---------------------------|---|
| width = INT | //sirka |
| height = INT | //vyska |
| num_x = INT | //pocet castic na sirku |
| num_y = INT | //pocet castic na vysku |
| draw_cloth = BOOL | //zda se maji nebo nemaji kreslit saty (pro mereni se nekresli) |
| force = FLOAT FLOAT FLOAT | //sila;muze jich byt vice; gravitace je v programu automaticky |

Priloha B - Obrázky

