

# Efektivita na GPU

Petr Holub

hopet@ics.muni.cz



PV197

2012-11-19

# Přehled přednášky

Vlastnosti CUDA

Metriky algoritmů

JPEG2000

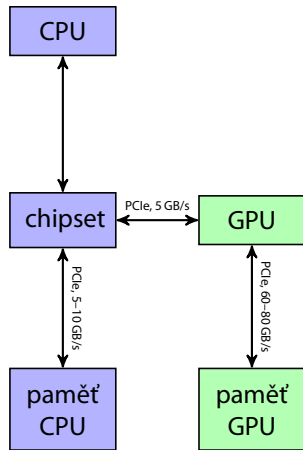
JPEG

# Literatura

- Park I. K., Singhal N., Lee M. H., Cho S., Kim C. W., “Design and Performance of Evaluation of Image Processing Algorithms on GPUs,” *IEEE Transactions on Parallel and Distributed Systems*, 2010 (zatím pouze v elektronické verzi)
- Cope B., Cheung P. Y. K., Luk W., Howes L., “Performance Comparison of Graphics Processors to Reconfigurable Logic: A Case Study”, *IEEE Transactions on Computers*, vol. 59, no. 4, April 2010
- Best Practices Guide – CUDA 2.2, 2009 [http://developer.download.nvidia.com/compute/cuda/2\\_3/toolkit/docs/NVIDIA\\_CUDA\\_BestPracticesGuide\\_2.3.pdf](http://developer.download.nvidia.com/compute/cuda/2_3/toolkit/docs/NVIDIA_CUDA_BestPracticesGuide_2.3.pdf)
- Wil Braithwaite, “The CUDA architecture: The Art of performance optimization”, Siggraph 2009, [http://developer.download.nvidia.com/presentations/2009/SIGGRAPH/asia/6\\_cuda\\_optimization.pdf](http://developer.download.nvidia.com/presentations/2009/SIGGRAPH/asia/6_cuda_optimization.pdf)

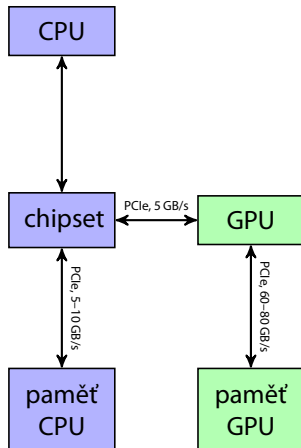
# Shrnutí práce s pamětí

- Práce s pamětí
  - omezení paměti
    - omezená šířka pásma mezi host a device (cca 6 GB/s pro PCI-e x16 Gen2)
    - latence globální paměti (teoreticky 141 GB/s, cca 400–600 cyklů latence)
 
$$\underbrace{1107 \times 10^6}_{\text{takt paměti [Hz]}} \times \underbrace{(512/8)}_{\text{Inteface paměti}} \times \underbrace{2}_{\text{DDR}} / 10^9 = 141,6 \text{ GB/s}$$
 (ev. 132 GB/s při dělení  $1024^3$ )
    - omezená velikost sdílené paměti
    - omezený počet registrů

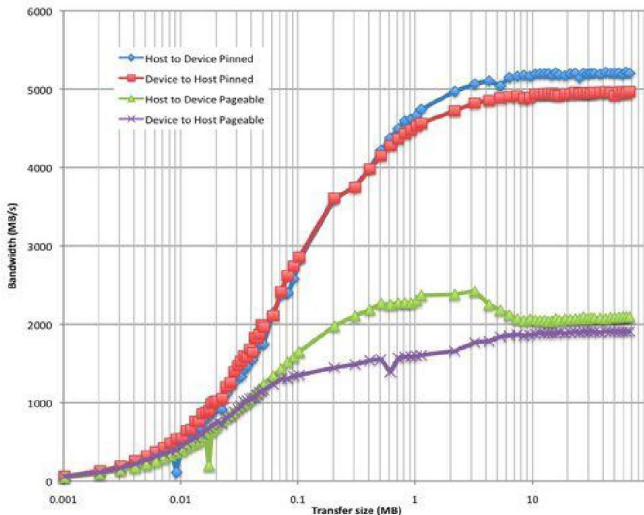


# Shrnutí práce s pamětí

- Práce s pamětí
  - optimalizace
    - ◆ maximální využití sdílené paměti a registrů
    - ◆ koalescentní přístup do globální paměti a vyhnutí se partition campingu
    - ◆ překrývání výpočtů a přístupu k datům
    - ◆ využití asynchronních přenosů: `cudaMemcpyAsync(...)`;
    - ◆ přiměřené používání page-locked paměti: `cudaHostAlloc(...)`;
    - ◆ co největší přenosy host ↔ device najednou



# Shrnutí práce s pamětí



Zdroj: Wil Braithwaite, "The CUDA architecture: The Art of performance optimization"

# Pokročilá práce s pamětí

- Další triky s pamětí
  - rozvrhnout, co cacheovat a co opakovaně počítat
  - mapování OpenGL bufferu do adresního prostoru zařízení (device)
    1. zaregistrujte si buffer pomocí CUDA-C  
`cudaGLRegisterBufferObject (GLuint buffObj);`
    2. namapujte zaregistrovaný buffer do globální paměti zařízení (vrátí adresu)  
`cudaGLMapBufferObject (void** devPtr, GLuint buffObj);`
    3. použijte adresu v kernelu
    4. odmapujte buffer  
`cudaGLUnmapBufferObject (GLuint buffObj);`
    5. odregistrujte buffer  
`cudaGLUnregisterBufferObject (GLuint buffObj);` (potřeba pouze pokud je buffer cíl rendrování)
    6. použijte buffer v OpenGL
      - ◆ může pomoci odstranit přenosy host ↔ device
      - ◆ automatické DMA mezi kartami Tesla a Quadro (momentálně přes host)
      - ◆ vykreslování z pixel buffer object pomocí `glDrawPixels` nebo `glTexImage2D`

## Pokročilá práce s pamětí

### Textura generovaná pomocí CUDA:

```

// setup code:
2  cudaGLRegisterBufferObject (pbo);
// CUDA texture generation code:
4  unsigned char *d_buffer;
   cudaGLMapBufferObject ((void**) &d_buffer, pbo);
6  prep_texture_kernel<<<...>>>(d_buffer);
   cudaGLUnmapBufferObject (pbo);
8  // OpenGL rendering code:
   glBindBuffer(GL_PIXEL_UNPACK_BUFFER_ARB, pbo);
10  glBindTexture(GL_TEXTURE_2D, tex);
   glTexSubImage2D(GL_TEXTURE_2D, 0, 0, 0, 256, 256, GL_BGRA, GL_UNSIGNED_BYTE, 0);

```

### Zpracování snímku pomocí CUDA:

```

// OpenGL rendering code:
2  // ...
// CUDA post-processing code:
4  unsigned char *d_buffer;
   cudaGLRegisterBufferObject (pbo);
6  cudaGLMapBufferObject ((void**) &d_buffer, pbo);
   post_process_kernel<<<...>>>(d_buffer);
8  cudaGLUnmapBufferObject (pbo);
   cudaGLUnRegisterBufferObject (pbo);

```

Zdroj: Wil Braithwaite, "The CUDA architecture: The Art of performance optimization"



# Pokročilá práce s pamětí

- Další triky s pamětí

- Write-Combining

```
cudaHostAlloc((void**) &h_data, num_bytes,  
              cudaHostAllocWriteCombined);
```

- ◆ paměť není cacheovaná ani cache koherentní
- ◆ PCI nedělá snooping
- ◆ podle CUDA 2.2 Pinned Memory APIs ([http://www.fcsc.es/download/Archivo%20Cursos/CUDA\\_Unileon\\_2009/CUDA2.2PinnedMemoryAPIs.pdf](http://www.fcsc.es/download/Archivo%20Cursos/CUDA_Unileon_2009/CUDA2.2PinnedMemoryAPIs.pdf)) může poskytnout až o 40 % větší výkon
- ◆ může zvýšit i výkon pro zápis procesorem (Write Combining Memory Implementation Guidelines, <http://download.intel.com/design/PentiumII/applnots/24442201.pdf>): agregací zápisu, obcházením L1/L2 cache
- ◆ problematické čtení – potřeba používat paměťové bariéry před čtením hodnot (`_mm_sfence` na Linuxu, `_WriteBarrier` na Windows, u SSE4 lze pro čtení použít instrukci `MOVNTDQA`); bariéry provádí CUDA driver a jsou relativně pomalé

# Pokročilá práce s pamětí

- Zero-copy mapování paměti
  - přímý přístup k datům v paměti CPU
  - schopnost tohoto se testuje pomocí pole `canMapHostMemory` dotazu `cudaDeviceProp`
  - automatický přenos dat po PCIe dle potřeby
  - relativně pomalé
    - ◆ pro jednorázově použitá malá data
    - ◆ pokud dokáže vysoký podíl výpočtu maskovat latenci
  - zajímavé zejména v kombinaci s integrovanými kartami – `integrated` pole `cudaDeviceProp`
    - ◆ na UMA architektuře odpadá přenos úplně
    - ◆ Nvidia ION

# Shrnutí práce s procesory

- Vlákna a multiprocesory
  - potřebujeme využít výpočetní výkon karty
  - kolik warpů potřebujeme k maskování latence globální paměti?
    - ◆ řekněme, že potřebujeme 100 aritmetických instrukcí k maskování latence (400 taktů latence / 4 takty na instrukci)
    - ◆ řekněme, že máme 8 aritmetických instrukcí ( $8 \times 4$  takty) na 1 přístup do globální paměti (400 taktů latence)
    - ◆  $100/8 \approx 13$  warpů
  - kolik warpů potřebujeme k maskování read-after-write latence registrů?
    - ◆ latence je cca 24 taktů
    - ◆  $24/4 \approx 6$  warpů

# Shrnutí práce s procesory

- Vlákna a multiprocesory
  - obsazení (occupancy)

$$O = \frac{\text{\# warpů běžících na MP v daný okamžik}}{\text{maximální \# souběžných warpů}}$$

- ◆ c. c. 1.2: maximálně 32 warpů  
c. c. 1.1: maximálně 24 warpů
- ◆  $13/32 = 40\%$  obsazení pro maskování latence globální paměti (c. c. 1.2)
- ◆  $6/32 = 18,75\%$  obsazení pro maskování latence registrů (c. c. 1.2)
- více vláken  $\stackrel{?}{=}$  větší výkon
  - ◆ záleží na zdrojích požadovaných vláknem
  - ◆ záleží na uspořádání vláken

# Shrnutí práce s procesory

- Limity multiprocesoru:

```

kepler$ ./deviceQuery
2
   CUDA Device Query (Runtime API) version (CUDA static linking)
4
Found 1 CUDA Capable device(s)
6
Device 0: "GeForce GTX 580"
8   CUDA Driver Version / Runtime Version          4.10 / 4.10
   CUDA Capability Major/Minor version number:    2.0
10  Total amount of global memory:                 1535 MBytes (1609760768 bytes)
   (16) Multiprocessors x (32) CUDA Cores/MP:     512 CUDA Cores
12  GPU Clock Speed:                               1.57 GHz
   Memory Clock rate:                             2010.00 Mhz
14  Memory Bus Width:                             384-bit
   L2 Cache Size:                                 786432 bytes
16  Max Texture Dimension Size (x,y,z)            1D=(65536), 2D=(65536,65535), 3D=(2048,2048,2048)
   Max Layered Texture Size (dim) x layers        1D=(16384) x 2048, 2D=(16384,16384) x 2048

```

# Shrnutí práce s procesory

- Limity multiprocesoru:

```
1  Total amount of constant memory:          65536 bytes
2  Total amount of shared memory per block:  49152 bytes
3  Total number of registers available per block: 32768
4  Warp size:                                32
5  Maximum number of threads per block:      1024
6  Maximum sizes of each dimension of a block: 1024 x 1024 x 64
7  Maximum sizes of each dimension of a grid: 65535 x 65535 x 65535
8  Maximum memory pitch:                     2147483647 bytes
9  Texture alignment:                        512 bytes
10 Concurrent copy and execution:           Yes with 1 copy engine(s)
11 Run time limit on kernels:                No
12 Integrated GPU sharing Host Memory:      No
13 Support host page-locked memory mapping: Yes
14 Concurrent kernel execution:             Yes
15 Alignment requirement for Surfaces:      Yes
16 Device has ECC support enabled:          No
17 Device is using TCC driver mode:         No
18 Device supports Unified Addressing (UVA): Yes
19 Device PCI Bus ID / PCI location ID:     1 / 0
20 Compute Mode:
21   < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously)
```

- Pozor na dostatečný počet vláken/warpů v bloku ( $\geq 96$ )
- Pozor na limity v počtu potřebných registrů

## Shrnutí práce s procesory

- Příklad limitů pro 8132 registrů a 24 warpů:
  - 10 registrů na vlákno, 256 vláken na blok
    - ◆ každý blok použije 2650 registrů  $\implies$  mohou běžet 3 bloky (7680 registrů)
    - ◆  $256 \times 3/32 = 24$  warpů může běžet současně  $\implies$  může dosáhnout 100 % využití
  - 17 registrů na vlákno, 256 vláken na blok
    - ◆ každý blok použije 4352 registrů  $\implies$  může běžet 1 blok (4352 registrů)
    - ◆  $256 \times 1/32 = 8$  warpů může běžet současně  $\implies$  může dosáhnout jen 33 % využití
  - 17 registrů na vlákno, 128 vláken na blok
    - ◆ každý blok použije 2176 registrů  $\implies$  mohou běžet 3 bloky (6528 registrů)
    - ◆  $128 \times 3/32 = 12$  warpů může běžet současně  $\implies$  může dosáhnout 50 % využití

## Shrnutí práce s procesory

- Určení používaných zdrojů
  - přeložíme s `-cubin`
  - výsledný `.cubin` soubor obsahuje

```

1 architecture {sm_10}
  abiversion {0}
3 modname {cubin}
  code {
5     name = MyKernel
      lmem = 0    // lokální paměť per blok
7     smem = 68 // sdílená paměť per blok
      reg = 20   // počet registrů per vlákno
9     bar = 0
      bincode {
11         0xa0004205 0x04200780 0x40024c09 0x00200780

```

- nebo použijeme `-ptxas-options=-v`

```

ptxas info : Used 4 registers, 60+56 bytes lmem, 44+40 bytes
smem, 20 bytes cmem[1], 12 bytes cmem[14]

```



# Shrnutí práce s procesory

- Omezení tlaku na registry
  - kompilátor se snaží počet registrů minimalizovat
  - `-maxrregcount=<N>` umožňuje nastavit požadovaný maximální počet registrů na kernel
  - přetečení do lokální paměti způsobí zpomalení
- Heuristiky pro velikost mřížky
  - # bloků > # multiprocesorů
    - ◆ aby všechny multiprocesory měly alespoň jeden blok k vykonávání
  - # bloků / # multiprocesorů > 2
    - ◆ na jednom procesoru může běžet více bloků
    - ◆ bloky nečekající v `__syncthreads()`; udržují zátěž hardware
  - # bloků > 1000
    - ◆ rezerva do budoucnosti přes několik generací karet

# Shrnutí práce s procesory

- Heuristiky na velikost bloku
  - čím více vláken v bloku, tím méně může jedno vlákno používat registrů
  - čím více vláken v bloku, tím hlubší pipeline a lepší maskování latence
  - počty vláken v bloku by měly být násobky 64
    - ◆ minimalizace bank konfliktů na registrech
    - ◆ rozumné: 192 nebo 256

# Shrnutí práce s procesory

- **CUDA Occupancy Calculator**

[http://developer.download.nvidia.com/compute/cuda/CUDA\\_Occupancy\\_calculator.xls](http://developer.download.nvidia.com/compute/cuda/CUDA_Occupancy_calculator.xls)

**CUDA GPU Occupancy Calculator**

Your chosen resource usage is indicated by the red triangle on the graphs. The other data points represent the range of possible block sizes, register counts, and shared memory allocation.

Resource	Value
21 Enter your resource usage:	
9 Threads Per Block	256
10 Registers Per Thread	1024
11 Shared Memory Per Block (bytes)	
23 GPU Occupancy Data is displayed here and in the graphs:	
17 Active Threads per Multiprocessor	1536
18 Active Warps per Multiprocessor	48
19 Active Thread Blocks per Multiprocessor	6
20 Occupancy of each Multiprocessor	100%
23 Physical Limits for GPU Compute Capability:	
22 Threads per Warp	32
24 Warps per Multiprocessor	48
25 Threads per Multiprocessor	1536
26 Thread Blocks per Multiprocessor	6
27 Total # of 32-bit registers per Multiprocessor	32768
28 Register allocation and size	64
29 Register allocation granularity	warp
30 Shared Memory per Multiprocessor (bytes)	49152
31 Shared Memory allocation and size	528
32 Warp allocation granularity (for register allocation)	0
34 Allocation Per Thread Block	
35 Warps	6
36 Registers	2048
37 Shared Memory	1024
38 These data are used in computing the occupancy data in blue	
40 Maximum Thread Blocks Per Multiprocessor	Blocks
41 Limited by Max Warps Per Multiprocessor	6
42 Limited by Registers per Multiprocessor	16
43 Limited by Shared Memory per Multiprocessor	48
44 Thread Block Limit Per Multiprocessor highlighted	6
46 CUDA Occupancy Calculator	
47 Version	2.1
48	<a href="#">License and License</a>

**Varying Block Size**  
My Block Size: 256

**Varying Register Count**  
My Register Count: 1024

**Varying Shared Memory Usage**  
My Shared Memory: 1024

Zdroj: Wil Braithwaite, "The CUDA architecture: The Art of performance optimization"

# Metriky z pohledu obrazových algoritmů

1. Paralelní podíl
2. Poměr mezi operacemi v plovoucí čárce a přístupů do paměti
3. Počet operací v plovoucí čárce na pixel
4. Počet přístupů do paměti na pixel
5. Míra větvení
6. Závislost úloh
  - Srovnávané algoritmy implementovány na CPU
    - výpočet metrik před implementací na CUDA

# Metriky z pohledu obrazových algoritmů

- Paralelní podíl
  - odpovídá Amdahlovu zákonu se všemi důsledky
  - při složeném algoritmu odpovídá poměrům částí
- Poměr mezi operacemi v plovoucí čárce a přístupů do paměti
  - jakékoli přístupy (sdílená, lokální, globální paměť)
  - skrývání latence překrýváním výpočtem
- Počet operací v plovoucí čárce (FP) na pixel
  - při výpočtech v plovoucí čárce překonávají GPU cca 20× CPU
  - obrazové zpracování zahrnuje typicky hodně operací v plovoucí čárce
  - nepřímě koreluje i s přístupů do paměti
- Počet přístupů do paměti na pixel
  - GPU mají přibližně 10× větší kapacitu přístupu do paměti než CPU
  - per pixel charakteristika umožňuje často využití sdílené paměti
  - problémy se sekvenčními přístupů do paměti u některých algoritmů – omezují paralelismus

# Metriky z pohledu obrazových algoritmů

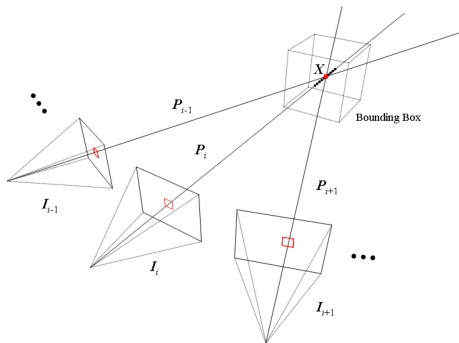
- Míra větvení
  - větvení přes `if`, `switch`, `do`, `for`, `while`
  - algoritmy pro zpracování obrázků často používají větvení na základě výsledku bitových operací
  - počítá se na základě rozptylu délky běhu jednotlivých vláken
  - berou se vlákna po blocích velikosti 32
- Závislost úloh
  - při zpracování obrázků se závislosti většinou řeší sekvenčním spouštěním CUDA kernelů
  - sleduje se počet bariér

# Relativní důležitost metrik

- Paralelní podíl > Větvení > FP operace per pixel > přístupy do paměti per pixel > poměr mezi FP operacemi a přístupy do paměti > závislost úloh
  1. Amdahl limituje vše
  2. omezení SIMT modelu
  3. algoritmy na zpracování obrazu mají obecně hodně FP operací
  4. poměr je závislý na předchozích dvou hodnotách (proč ho zavádět?)
  5. závislost úloh určuje obtížnost implementace (opravdu?)

# Studované pilotní algoritmy

- Stereo srovnávání obrázků (Multiview Stereo Matching, MVS)





# Studované pilotní algoritmy

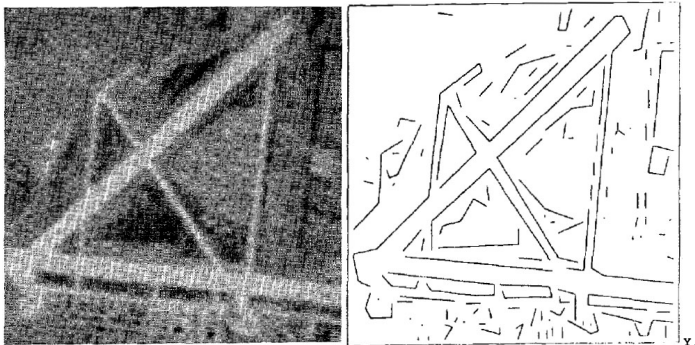
- Stereo srovnávání obrázků (Multiview Stereo Matching, MVS)
  - popis algoritmu
    - ◆ vstup:
      - zkalibrované obrazy  $I = I_0, \dots, I_{N-1}$
      - projekční matice  $P = P_0, \dots, P_{N-1}$
    - ◆ výstup:
      - 3D body  $X = X_0, \dots, X_{M-1}$
    - ◆ porovnávání lokálních oken mezi jednotlivými obrazy
    - ◆ hloubka se pro bod  $(x, y)$  v referenčním obrazu určí přeložením příslušné oblasti z  $I_i$  na oblasti v obraz  $I_{i-1}$
    - ◆ bod s minimálním součtem absolutních odchylek (SAD) a normalizovanou cross-korelací (NCC) je uložen jako best-match
    - ◆ opakujeme pro ostatní sousedící obrazy (např.  $I_{i+1}$ )
    - ◆ bod je korektně určen, pokud počet best-match je nad hranici MIN\_COUNT
    - ◆ opakujeme pro všechny body a všechny referenční obrazy
  - složitost:  $\mathcal{O}(N^2 WHL)$ 
    - $N$  ... počet vstupních obrazů,  $W$  resp.  $H$  ... vodorovné resp. svislé rozlišení,  $L$  ... velikost ohraničujícího boxu

# Studované pilotní algoritmy

- Stereo srovnávání obrázků (Multiview Stereo Matching, MVS)
  - mapování na GPU
    - ◆ porovnávání lokálních oken je nezávislé – dobře mapuje na GPU
    - ◆ vlákno  $\sim$  pixel, tj.  $W \times H$  vláken
    - ◆  $\mathcal{O}(N)$  volání kernelu pro výpočet hloubky pro jeden referenční obraz
    - ◆  $\mathcal{O}(L)$  volání uvnitř kernelu
    - ◆ obrázky se nakopírují do globální paměti
    - ◆ koeficienty lokálního okna jdou do sdílené paměti – častý přístup
    - ◆ složitost:  $\mathcal{O}\left(\frac{N^2 WHL}{T_{max}}\right)$   
 $T_{max}$ ... maximální počet vláken na GPU (např. 12288 na G80)

# Studované pilotní algoritmy

- Získávání lineárních charakteristik (Linear Feature Extraction)
  - rozpoznávací aplikace: budovy, silniční pruhy, ...



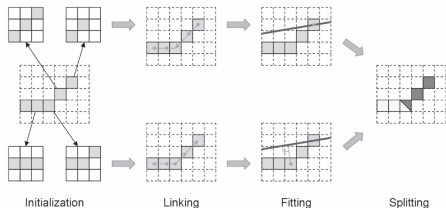
T. Zhou, "Linear Feature Extraction Based on an AR Model Edge Detector"  
<http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=01169687>

# Studované pilotní algoritmy

- Získávání lineárních charakteristik (Linear Feature Extraction)
  - popis algoritmu
    - ◆ vstup: obrázek
    - ◆ výstup: obrázek s detekovanými hranami
    - ◆ použijeme algoritmus Nevatia-Babu s Cannyho detektorem hran
  - 1. detekce hran
  - 2. zúžení hran
  - 3. vytvoření řetízků hran na základě spojitosti v 8 směrech
  - 4. fitting čar na řetízky
  - 5. pokud segment čáry dává větší chybu než hraniční hodnota, je čára rozdělena na 2
  - 6. poslední krok iterativně opakujeme

# Studované pilotní algoritmy

- Získávání lineárních charakteristik (Linear Feature Extraction)
  - mapování na GPU



- ◆ 6 per-pixel kernelů

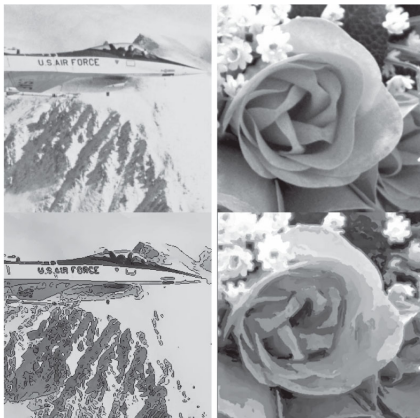
1. Cannyho detektor hran jako per-pixel filtr
2. klasifikace pixelů podle okolí  $3 \times 3$  – Initialization
3. nalezení co nejdelších souvislých částí, určení počátečního a koncového bodu – Linking
4. výpočet odchylky od úsečky proložené mezi počáteční a koncový bod – Fitting
5. pokud odchylka přesahuje  $D_{max}$ , rozdělí se souvislý řetězec na dva
6. iterujeme poslední 3 kroky, dokud vznikají nové segmenty

# Studované pilotní algoritmy

- Získávání lineárních charakteristik (Linear Feature Extraction)
  - mapování na GPU
    - ◆ počítá se pro všechny pixely – i ty, co hrany neobsahují  $\implies$  omezení
    - ◆ redundantní hledání nejdelších souvislých čar
    - ◆ významný podíl over-computation

# Studované pilotní algoritmy

- Nefotorealistický rendering



- Cartoon-style NPR
- Oily-style NPR

# Studované pilotní algoritmy

- Nefotorealistický rendering
  - popis algoritmu: Cartoon-style NPR
    - ◆ bilaterální filtrování (plochy)
    - ◆ Cannyho detekce hran (čáry)
    - ◆ přeložení čar přes plochy
  - mapování na GPU: Cartoon-style NPR
    - ◆ per pixel pro všechny operace
    - ◆ vstupní obrázky jsou uloženy ve 2D texturové paměti



# Studované pilotní algoritmy

- Nefotorealistický rendering
  - popis algoritmu: Oily-style NPR
    - ◆ rozdíl původního obrazu a gaussovsky rozostřeného obrazu per pixel (míra rozostření udává tloušťku štětce)
    - ◆ tah se generuje, pokud oblast v součtu rozdílů dosahuje nad stanovenou hranici
    - ◆ tah začíná od lokálního maxima a vede po gradientu
    - ◆ tah se uloží pouze je-li dost dlouhý
    - ◆ uložené tahy se „obtáhnou štětcem“ dané tloušťky
    - ◆ proces začíná s čistým pozadím a opakuje se od nejtlustšího štětce k nejtenčímu
  - mapování na GPU: Oily-style NPR
    - ◆ paralelizace per tah: problém s překryvem tahů
    - ◆ paralelizace per pixel: problém jak vybrat pořadí tahů
    - ◆ reformulace/heuristika: světlejší štětec se použije později než tmavší  
⇒ pro daný pixel vybereme nejsvětlejší barvu překrývajících se tahů
    - ◆ musíme hledat, kterými tahy bude pixel ovlivněn (maximální vzdálenost od trajektorie tahu)
    - ◆ vstupní obrázky jsou uloženy ve 2D texturové paměti

# Přehled pilotních algoritmů

Characteristics Algorithms	Processing Domain	Degree of Concurrency	Branching Diversity of Parallel Threads
Multiview Stereo Matching	Pixel	Massive	Low
Linear Feature Extraction	Pixel / Feature	Medium	Low (Edge) High (Linking, Fitting)
JPEG2000 Encoding	Pixel / Bitplanes	Massive (DWT) Low (EBCOT)	High
Non-Photorealistic Rendering	Pixel	Massive	Low

Characteristics Algorithms	Floating-Point Arithmetic Intensity	Other Features
Multiview Stereo Matching	Low	-
Linear Feature Extraction	High	Overcomputation occurs
JPEG2000 Encoding	High (DWT) High (EBCOT)	-
Non-Photorealistic Rendering	High	Only pixelwise convolutions

# Výsledky pilotních algoritmů

- charakteristika algoritmů

Algorithms \ Characteristics	Parallel Fraction (↑)	FP Computation to Memory Access Ratio (↑)	FP Computation Per Pixel (↑)
Multiview Stereo Matching	0.994	1.838	4,900
Linear Feature Extraction	0.706	5.183	1,083
JPEG2000 Encoding (DWT)	0.983	3.630	552
JPEG2000 Encoding (Tier-1)	0.834	7.630	579
Cartoon-Style NPR	0.987	8.043	45,118
Oily-Style NPR	0.992	13.363	4,972

- MVS kompenzuje malý poměr FLOP/mem vysokou mírou paralelismu
- DWT kompenzuje málo FLOP/pixel vysokou mírou paralelismu
- LFE má problém s nízkou mírou paralelismu

# Výsledky pilotních algoritmů

- charakteristika algoritmů

Algorithms \ Characteristics	Memory Access Per Pixel (↑)	Branching diversity (↓)	Task dependency (↓)
Multiview Stereo Matching	2,665	0.117	1
Linear Feature Extraction	209	0.113	11
JPEG2000 Encoding (DWT)	152	0.138	12
JPEG2000 Encoding (Tier-1)	76	0.307	1
Cartoon-Style NPR	5,609	0.156	6
Oily-Style NPR	372	0.121	34

- task depenednecy: implementace bude stát hodně úsilí (schoval se nám díky zvolenému přístupu EBCOT Tier-1)
  - ⇒ Oily-Style NPR má mnoho iterací, během nichž aktualizuje buffer
- EBCOT Tier-1 má problém s větvením, nižším paralelismem a malým počtem FLOP/pixel

# Výsledky pilotních algoritmů

- benchmarky CUDA

Algorithms	Characteristics	Global memory coalescing (↑)	SM to GM access ratio (↑)	Global memory Transfer (ms) (↓)	GPU occupancy (↑)
Multiview Stereo Matching		0.0013	0.020	107.28	33%
Linear Feature Extraction		0.0256	0.090	3.382	41.65%
JPEG2000 Encoding (DWT)		0.000	5.000	3.069	33.33% (vertical) 83.33% (horizontal)
JPEG2000 Encoding (Tier-1)		0.0167	1.903	0.307	17.00%
Cartoon-Style NPR		0.1439	0.000	0.435	40.18%
Oily-Style NPR		0.1517	0.823	0.435	52.19%

- DWT: agregace dlaždic sniží potřebný počet registrů (  $\implies$  GPU occupancy)

# Výsledky pilotních algoritmů

- benchmarky CUDA

Algorithms	Characteristics	Data dependency (↓)	Source lines (Host)	Kernel lines (Device)
Multiview Stereo Matching		2	184	187
Linear Feature Extraction		36	620	318
JPEG2000 Encoding (DWT)		2	132	112 (vertical) 90 (horizontal)
JPEG2000 Encoding (Tier-1)		2	>1500	406
Cartoon-Style NPR		0	223	232
Oily-Style NPR		4	514	635

- Data Dependency: celkový počet volání `__syncthreads()`; v rámci bloku vláken

# Výsledky pilotních algoritmů

- zrychlení

Algorithms	Data Resolution	Speedup		GPU Scalability
		G92/CPU	GX200/CPU	GX200/G92
Multiview Stereo Matching	TempleRing (47 Images)	54.18x	167.47x	3.09x
Cartoon-Style NPR	512 × 512	93.71x	149.84x	1.61x
	1024 × 768	97.51x	168.46x	1.73x
	1280 × 1024	76.79x	131.34x	1.71x
	1200 × 1800	117.30x	201.27x	1.72x
	2288 × 1712	126.64x	219.03x	1.73x
Oily-Style NPR	512 × 512	81.71x	129.72x	1.59x
	1024 × 768	69.00x	112.52x	1.63x
	1280 × 1024	105.47x	159.17x	1.51x
	1200 × 1800	83.16x	130.00x	1.56x
	2288 × 1712	85.27x	139.01x	1.63x

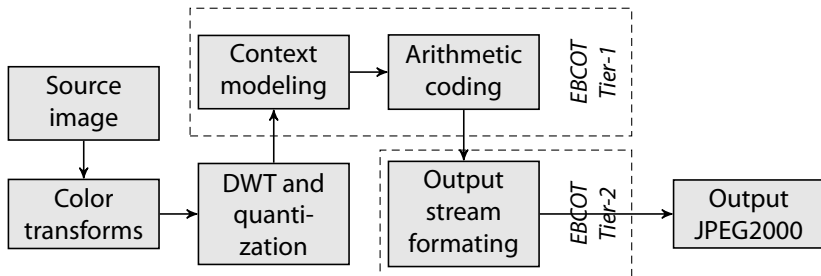
# Výsledky pilotních algoritmů

- zrychlení

Algorithms	Data Resolution	Speedup		GPU Scalability
		G92/CPU	GX200/CPU	GX200/G92
Linear Feature Extraction	512 × 512	1.99x	2.42x	1.22x
	1024 × 768	2.53x	2.90x	1.14x
	1280 × 1024	2.43x	2.72x	1.11x
	1200 × 1800	2.65x	3.22x	1.22x
	2288 × 1712	2.33x	3.00x	1.29x
JPEG2000 Encoding (DWT)	512 × 512	4.06x	6.94x	1.71x
	1024 × 768	7.27x	12.92x	1.78x
	1280 × 1024	5.28x	8.98x	1.70x
	1200 × 1800	5.14x	9.13x	1.78x
	2288 × 1712	5.18x	9.26x	1.79x
	3024 × 2089	5.30x	9.34x	1.76x
JPEG2000 Encoding (Tier-1)	512 × 512	0.46x	0.35x	0.77x
	1024 × 768	0.60x	0.70x	1.17x
	1280 × 1024	0.68x	0.75x	1.11x
	1200 × 1800	1.21x	1.12x	0.92x
	2288 × 1712	1.12x	1.61x	1.44x
	3024 × 2089	0.73x	0.98x	1.35x



## JPEG2000 – přehled procesu

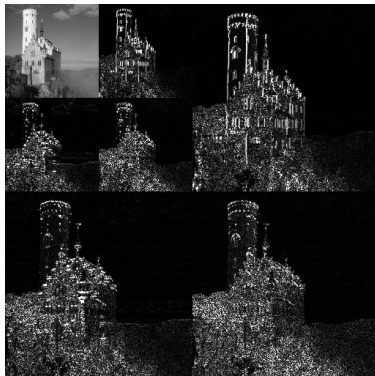
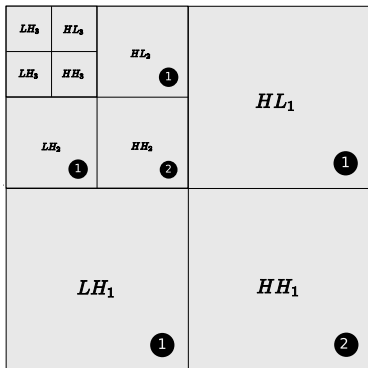


- Literatura z naší zahrádky :-))

- MATELA, Jiří. GPU-Based DWT Acceleration for JPEG2000. In Annual Doctoral Workshop on Mathematical and Engineering Methods in Computer Science. Brno : NOVAPRESS s.r.o., 2009. od s. 136-143, 8 s. ISBN 978-80-87342-04-6.
- MATELA, Jiří - RUSŇÁK, Vít - HOLUB, Petr. Efficient JPEG2000 EBCOT Context Modeling for Massively Parallel Architectures. In Storer, James A. and Marcellin, Michael W.. Data Compression Conference (DCC), 2011. Washington, DC, USA : IEEE Computer Society, 2011. od s. 423-432, 10 s. ISBN 978-0-7695-4352-9.
- MATELA, Jiří - ŠROM, Martin - HOLUB, Petr. Low GPU Occupancy Approach to Fast Arithmetic Coding in JPEG2000. In Z. Kotásek et al.. MEMICS 2011, LNCS 7119 - to appear. Heidelberg : Springer, 2011. od s. 136-145.

# Diskrétní vlnková transformace

- Rozložení obrazu do rekurzivně se opakujících pásů LL, HL, LH, HH



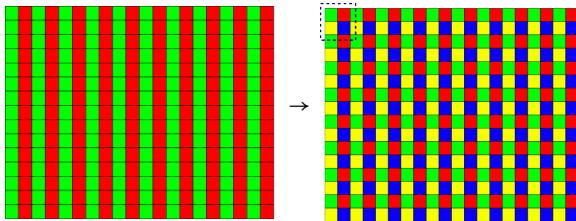
# Diskrétní vlnková transformace

- Lifting schéma
  - low-pass a high-pass filtry

$$d_i^1 = d_i^0 - \frac{1}{2}(s_i^0 + s_{i+1}^0)$$

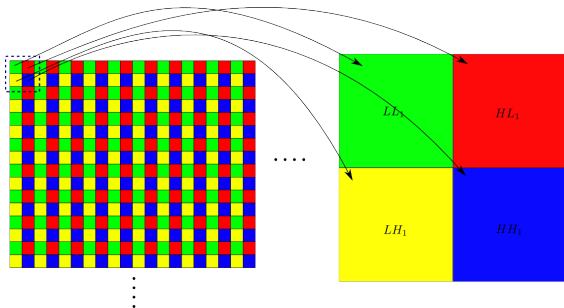
$$s_i^1 = s_i^0 + \frac{1}{4}(d_{i-1}^1 + d_i^1)$$

- Aplikace na řádky, poté na sloupce



# Diskrétní vlnková transformace

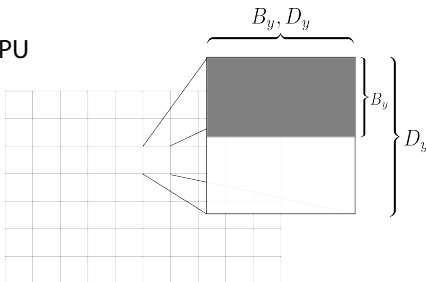
- Lifting schéma
  - přeuspořádání výsledného obrazu



- Mapování na GPU
  - 2D thread block
  - každé vlákno zpracovává jeden lichý a jeden sudý prvek

# Diskrétní vlnková transformace

- Mapování na GPU



1. načtení z globální do sdílené paměti
2. aplikace lifting schéma na řádky

$$s[T_x][2T_y + 1] = s[T_x][2T_y + 1] + p \cdot (s[T_x][2T_y] + s[T_x][2T_y + 2])$$

$$s[T_x][2T_y] = s[T_x][2T_y] + u \cdot (s[T_x][2T_y - 1] + s[2T_y + 1])$$

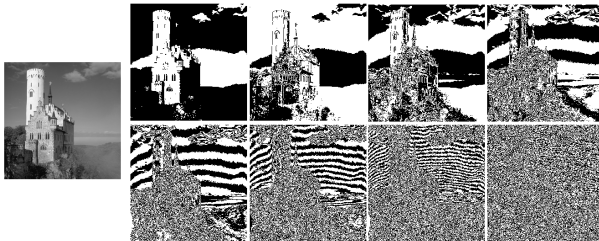
3. aplikace lifting schéma na sloupce (prohození  $T_x$  a  $T_y$ )
4. uložení výsledků do globální paměti s přeuspořádáním
  - ◆ uložení sudých pak lichých řádků
  - ◆ v rámci řádku první polovina vláken čte LL (resp. LH), druhá polovina HL (resp. HH)  $\implies$  koalescentní přístup do globální paměti

# Diskrétní vlnková transformace

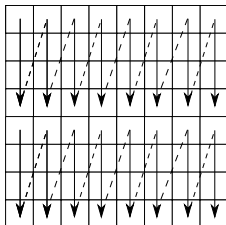
Implementation	512×512	1920×1080	Speedup
JasPer	6ms	55ms	N/A×
CUDA DWT	0.12ms	0.81ms	67.9×

# Modelování kontextu

- Hledání kontextu pro kompresi aritmetickým adaptivním kóděm
- Jde se po jednotlivých bitplanech

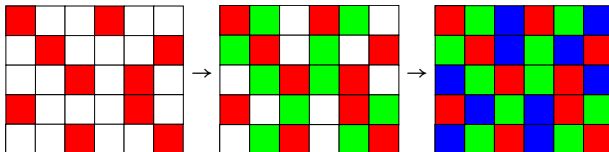


- Bitplane skenován podle vzoru (scan-pattern)



# Modelování kontextu

- 3 fáze:
  1. Significance Propagation Pass (SPP)
  2. Magnitude Refinement Pass (MRP)
  3. Cleanup Pass (CUP)





# Modelování kontextu

- 4 kódovací operace
  1. Zero Coding (ZC)
  2. Run-Length Coding (RLC)
  3. Magnitude Refinement Coding (MRC)
  4. Sign Coding (SC)
- každý bit je zakódován jednou nebo více operacemi právě v jedné fázi
- výstupem kódovacích operací je pár  $CX,D$

# Modelování kontextu

- Určení fáze pro daný bit
  - záleží na jeho stavu
  - záleží na stavu jeho sousedů
  - to vše se vyvíjí (proměnné  $\sigma, \sigma', \eta$ ), jak bity prochází jednotlivými fázemi podle scan-pattern – :-(((((((((((
- Dle definice vysoce sekvenční proces
  - umožňuje paralelizmus jedině na úrovni code block
  - ... s jistými se speciálními výjimkami (causal mode)

## Modelování kontextu

- Reformulace problému

- náhrada původních stavových proměnných  $\sigma, \sigma', \eta$  za proměnné  $\rho, \tau$
- nové proměnné lze předpočítat paralelně pro každý bit
- výpočet  $\rho$

Bitplane index  $1 \dots n$

Value of a pixel at  $x, y$

$$\rho_{x,y}^p = \begin{cases} 1, & \text{if } \gamma_{x,y} > 2^p \\ 0, & \text{otherwise} \end{cases}$$

Pixel position in the code block

# Modelování kontextu

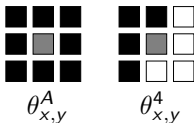
- Reformulace problému

- výpočet  $\tau$

- $\tau_{x,y}^p = 1 \Leftrightarrow \rho_{x,y}^{p+1} = 0 \wedge \gamma_{x,y}^p = 1 \wedge \bigvee_{(i,j) \in \theta_{x,y}^A} (\rho_{i,j}^{p+1} = 1)$

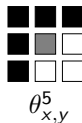
- $\tau_{x,y}^p = 1 \Leftrightarrow \rho_{x,y}^{p+1} = 0 \wedge \gamma_{x,y}^p = 1 \wedge \bigvee_{(i,j) \in \theta_{x,y}^4} (\tau_{i,j}^{p+1} = 1)$

- Step 2 is repeated until there is no new  $\tau_{x,y}^p = 1$  found



## Modelování kontextu

- Reformulace problému
  - pomocná proměnná  $\delta$ 
    - $\delta_{x,y}^p = 1$  indicates a position  $\gamma_{x,y}^p$  is in PN in SPP
    - $\delta_{x,y}^p = \bigvee_{(i,j) \in \theta_{x,y}^A} (\rho_{i,j}^{p+1} = 1) \vee \bigvee_{(i,j) \in \{\theta_{x,y}^5, \theta_{x,y}^4, \theta_{x,y}^3\}} (\tau_{i,j}^p = 1)$
    - Where selection of surrounding  $\theta_{x,y}$  depends on position of the bit  $\gamma_{x,y}^p$  on y-axis.



# Modelování kontextu

## • Důkazy ekvivalence

### ■ Původní:

Bit je kódován v SPP, pokud není signifikantní ( $\sigma = 0$ ) a je v preferovaném okolí (preferred neighborhood).

### ■ Paralelní:

Bit je kódován v SPP, pokud není signifikantní ( $\rho^{p+1} = 0$ ) a je v preferovaném okolí ( $\delta^p = 1$ ).

### ■ ... následuje důkaz ekvivalence indukcí ;-))

is processed in the four consecutive steps: SEB; BEC; BE; EC. Note that each coding operation is executed for all bits in a bitplane in parallel. The only constraint on the coding bitplanes comes from the binary number of bits coded by the BEC operation. The BE is defined to code one to four bits in a plane and a prediction of the result is available as a separate bit in the BEC coding itself. The only operations allowed for bits in EC are: we can choose to perform BEC operations to correct higher bitplanes; following the proposed design allows for prediction across bit planes too, we do not require a feature of limited channel capacity over an uncompressible GPU's.

The modified hierarchical parallel algorithm allows for processing individual bits in parallel threads, enabling a high utilization of multi-processors on GPUs. Depending on chosen word length, the data may be processed entirely on the first channel element.

#### V. POWER OF REPRESENTATION

The reproducibility of the original algorithm and the proposed parallel algorithm is based on reproducibility of information provided by the original state variables  $\sigma^i$  and  $\rho^i$  and by the newly defined  $\rho^i$  and  $\delta^i$ . The proof in Section 3 uses reproducibility of state  $\rho$  and  $\delta$  already presented in Section 3, and reproducibility of  $\sigma$  and  $\rho$  in an SPP on current bit plane (Lemma 2). The reproducibility of information provided by a set  $\rho$  is weakens in the description of SEB operation. Reproducibility of information  $\rho$  after  $\sigma^i = \rho^i$  has also been described in SEB.

**Definition 1.** A state variable  $\sigma_{i,j}$  or  $\rho_{i,j}$  or  $\delta_{i,j}$  is called *redundant* (because redundant) in the original algorithm right after the most significant bitplane  $\rho_{i,j}$  is processed in the current SPP on GPU device. A bit  $\rho_{i,j}$  becomes redundant in SPP if right before the bit is processed, it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \delta_{i,j} = 1 \wedge \rho_{i,j} = 1$ . Otherwise  $\rho_{i,j} = 0 \wedge \rho_{i,j} = 1$ , the bit becomes significant in SPP.

**Lemma 1.** Value of  $\sigma_{i,j}$  is equal to  $\rho_{i,j}$  prior to  $\rho_{i,j}$  bit processing.

**Proof.** Value of  $\sigma_{i,j}$  is set to 1 if each coefficient  $\rho_{i,j}$  was the most significant bit of the coefficient has been processed in the original sequential algorithm. I.e. right after processing bit  $\rho_{i,j}$ ,  $\sigma_{i,j} = 1$  if one of bits  $\rho_{i,j}^k = \rho_{i,j}$  was equal to 1, then  $\sigma_{i,j} = \bigvee_k \rho_{i,j}^k$ . Following from Def. 1,  $\rho_{i,j} = \rho_{i,j}$ . In the sequential algorithm, this has to occur prior to processing  $\rho_{i,j} = 1$  bitplane and thus also before  $\rho_{i,j}^k = 1$  is processed.  $\square$

Note, that  $\sigma_{i,j}$  in the original algorithm may change from 1 to 0 right after  $\rho_{i,j}^k$  is processed and this is not reflected by  $\rho_{i,j}$  but will be reflected by  $\delta_{i,j}^k$ .

**Lemma 2.** For a bit plane  $\rho$  processed on the first step (SEB) if holds that  $\sigma_{i,j} = 1 \wedge \rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$ , where  $k$  denotes SEB operation.

**Proof.** Using mathematical induction we want to show that in each step of the SPP algorithm  $\rho_{i,j}$  is a bit plane  $\rho$ .

Power of representation  
Theorem 1. The power of representation of the original SPP algorithm is based on the fact that the state variables  $\sigma$  and  $\rho$  are reproduced across bit planes and thus can be processed entirely on the first channel element.

1.  $\sigma_{i,j} = 1$  or  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$
  2.  $\sigma_{i,j} = 1$  or  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$
- Proof.** Let us assume that values of all state variables are equal to zero for bitplane  $\rho$  beyond the current bit plane.
1. We want to show that if  $\rho_{i,j} = 1$  then  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$ .
    - A. Let  $\rho_{i,j} = 1$  before  $\rho_{i,j}$  is processed. Then  $\rho_{i,j} = \rho_{i,j}^k = 1$  (Lemma 1). Following from Def. 1,  $\rho_{i,j} = 0$  because  $\rho_{i,j}^k = 1$ .
    - B. Let  $\rho_{i,j} = 1$  right after bit before last  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$ , and thus right before last  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .
  2. We want to show that if  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  then  $\rho_{i,j} = 1$ .
    - A. Let  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  before  $\rho_{i,j}$  is processed. Then  $\rho_{i,j} = 1$  (Lemma 1). Following from Def. 1,  $\rho_{i,j} = 0$  because  $\rho_{i,j}^k = 1$ .
    - B. Let  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  right after bit before last  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$ , and thus right before last  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .

**Induction step.** We assume that  $\sigma_{i,j} = 1$  or  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  for each  $\rho_{i,j}^k$  in the current bitplane  $\rho$  according to the presented state variables.

1. Let  $\rho_{i,j} = 1$ . Then:
  - A. We want to show that  $\rho_{i,j} = 1$  becomes significant in case of previous bit plane  $\rho$  and thus before  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$  and thus before  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .
  - B. Let  $\rho_{i,j} = 1$  right after bit before last  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$ , and thus right before last  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .
2. Let  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$ . Then:
  - A. We want to show that  $\rho_{i,j} = 1$  becomes significant in case of previous bit plane  $\rho$  and thus before  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$  and thus before  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .
  - B. Let  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  right after bit before last  $\rho_{i,j}$  is processed. This coefficient  $\rho_{i,j}$  becomes significant in current bit plane  $\rho$ , and thus right before last  $\rho_{i,j}$  is processed it holds that  $\sigma_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1 \wedge \rho_{i,j} = 1$ .

Power of representation  
Theorem 1. The power of representation of the original SPP algorithm is based on the fact that the state variables  $\sigma$  and  $\rho$  are reproduced across bit planes and thus can be processed entirely on the first channel element.

Both  $\rho_{i,j}^k$  can be combined together with respect to remaining  $k$ . Hence

$$\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1 \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 0 \right) \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 1 \right)$$

Therefore  $\rho_{i,j} = 1$  (from Definition 1. The proof of the cases 1B.2 and 1B.1B is analogical to the 1B.1).

2. Let  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$ . Then:
  - A. We want to show that  $\rho_{i,j} = 1$  becomes significant in case of previous bit plane, hence prior to  $\rho_{i,j}$  (Lemma 2).
  - B. We want to show that  $\rho_{i,j} = 1$  becomes significant in current bit plane, hence after  $\rho_{i,j}$  (Lemma 2).

Using induction assumption we expand each  $\rho_{i,j}^k = 0 \wedge \rho_{i,j} = 1$  as  $\rho_{i,j} = 1$  in case of 2B.1 or get

$$\rho_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 0 \right) \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 1 \right)$$

Each  $\rho_{i,j}$  can be written to a bitplane close two state variables has same value before last  $\rho_{i,j}$  is processed in current bitplane  $\rho$  (Lemma 2). Hence

$$\rho_{i,j} = 0 \wedge \rho_{i,j} = 1 \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 0 \right) \wedge \left( \bigvee_{k=1}^n \rho_{i,j}^k = 1 \right)$$

Henceforth  $\rho_{i,j} = 0 \wedge \rho_{i,j} = 1$  is right after last  $\rho_{i,j}$  is coded.  $\square$

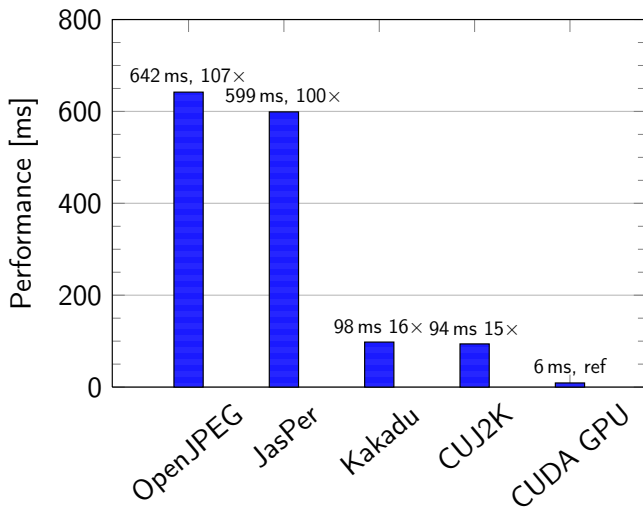
#### VI. REPRESENTATION BOUNDS

Mathematically we set up the benchmark only based on the SEB step. This step is processed in parallel. Performance of SEB step (GPU implementation) was compared with the sequential SPP implementation (OpenCV 2.2.3). In both cases, the number of coefficients (SEB) is 2.2.3 and one GPU implementation (SEB) is 1.1. All the GPU implementations were single-threaded.

Power of representation  
Theorem 1. The power of representation of the original SPP algorithm is based on the fact that the state variables  $\sigma$  and  $\rho$  are reproduced across bit planes and thus can be processed entirely on the first channel element.

# Modelování kontextu

Color image in 1080p resolution



# Aritmetické kódování – MQ-Coder

- Binární aritmetický kodér
  - LPS do  $[0, Q)$
  - MPS do  $[Q, A]$
  - $Q$  je pravděpodobnost výskytu LPS
  - výsledek  $C$
- Adaptivní
  - mění se význam LPS/MPS
- Kontextový
  - vstup má přiřazen kontext
  - kontextu určuje stav kodéru
    - ◆ mapování symbolů na MPS/LPS
    - ◆ pravděpodobnost LPS

MPS	$A = A(1 - Q)$	$C = C + AQ$
LPS	$A = AQ$	$C = C$



# Aritmetické kódování – MQ-Coder

- Problémy
  - příliš hrubý paralelismus
  - podmíněné větvené závislé na vstupních datech
  - nedá se odstěhovat na CPU – na GPU máme tou dobou 2× tolik dat ( $CX,D$ )

- Pozitiva
  - není tak náročné jako kontextové modelování

	720p	1080p	4K
Basic	38 ms	45 ms	93 ms

- Nenašli jsme způsob principiální paralelizace :-(
  - analýza různých optimalizačních technik
  - někteří výzkumníci vymýšlejí jiné kodéry – problém nekompatibility se standardem

# Aritmetické kódování – MQ-Coder

- Převod implementace do registrů
  - lokální datové struktury do registrů
  - průměrně: 240% zrychlení

⇒ registry jsou opravdu nejrychlejší (i vůči sdílené paměti)
- Rozvinování cyklů
  - navrženo jako optimalizační technika pro VLSI
  - zpracovává sérii MPS symbolů, pokud mají přiřazené stejné  $CX$
  - použití MAD instrukce:  $A = A - nQ_e$ ,  $C = C + nQ_e$
  - paralelní prohledávání – detekce sekvencí
  - průměrně: 31% zpomalení

⇒ zrychlení zpracování ani nevyváží režii předpočítání

⇒ zastoupení sérií MPS symbolů není dostatečně časté (17–25 % sekvencí 2–32)

# Aritmetické kódování – MQ-Coder

- Prefix sum
  - zpracovává sérii MPS symbolů nezávisle na  $CX$
  - nemůže používat MAD (různé  $CX$  znamená různé  $Q_e$ )
  - $Q_e$  jsou sčítány paralelně, součet použit pro úpravu  $A$  a  $C$

# Aritmetické kódování – MQ-Coder

- Renormalizace
    - pokud  $A \approx 1$ , můžeme aproximovat:  $A(1 - Q) \approx A - Q$  a  $AQ = Q$
    - omezení počtu násobení
  - Vylepšená rozšířená renormalizace
    - renormalizace se provádí, pokud  $A$  klesne pod  $0x8000$
    - $n$ -krát se násobí  $A$  a  $C$ , až je  $A > 0x8000$
    - existující návrh urychlení:  $n$  určit vyhledávací tabulkou a použít shift
    - publikovaný algoritmus měl chybu – neřešil korektně přetečení
    - můžeme využít CLZ instrukci místo vyhledávací tabulky
    - průměrně: 39% zrychlení
- ⇒ shifty a CLZ jsou rychlé

# Aritmetické kódování – MQ-Coder

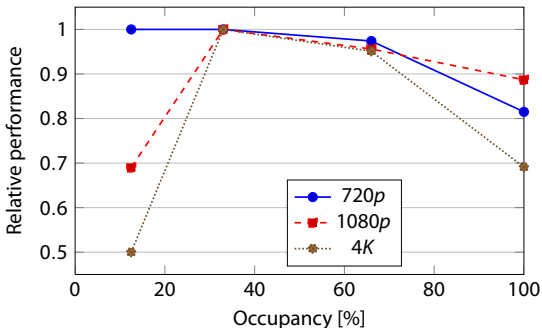
- Načítání dat po blocích
  - hrubé vláknění vede na nekoalescentní přístup do paměti – vlákna zpracovávají různé části bloku
  - pomůže načítat data po větších blocích<sup>1</sup>
  - využití **double** (8B) místo **int** (4B)
  - optimalizace počtu natahovaných **double** podle výkonu
  - průměrně: 33% zrychlení pro 16 natahovaných **double**

---

<sup>1</sup>Volkov, V.: Better Performance at Lower Occupancy. In: GPU Technology Conference 2010. (2010)

# Aritmetické kódování – MQ-Coder

- Optimalizace využití karty
  - díky hrubé granularitě každé vlákno potřebuje hodně zdrojů
  - vytváří tlak na registry
  - optimalizace využití karty na výsledný výkon



# Aritmetické kódování – MQ-Coder

- Shrnutí optimalizací

	<i>720p</i>		<i>1080p</i>		<i>4K</i>	
GPU Basic	38.0 ms	–	45.0 ms	–	92.9 ms	–
GPU R	16.2 ms	2.3×	18.9 ms	2.4×	48.5 ms	1.9×
GPU R+ERN	11.9 ms	3.2×	14.9 ms	3.0×	44.9 ms	2.1×
GPU R+LU	18.1 ms	2.1×	31.2 ms	1.4×	87.1 ms	1.1×
GPU R+PS	20.9 ms	1.8×	25.7 ms	1.8×	64.6 ms	1.4×
GPU R+CL	12.1 ms	3.1×	13.2 ms	3.4×	27.7 ms	3.4×
GPU R+ERN+CL	7.3 ms	5.2×	8.1 ms	5.6×	17.6 ms	5.3×

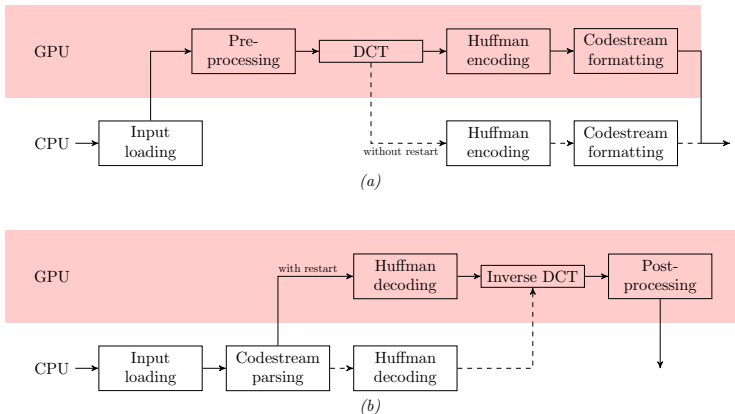
# Aritmetické kódování – MQ-Coder

- Výsledný výkon

	<i>720p</i>	<i>1080p</i>	<i>4K</i>
OpenJPEG 1.4	157 ms	316 ms	1081 ms
Jasper 1.900.1	89 ms	178 ms	594 ms
Kakadu 6.4 (4t)	41 ms	84 ms	284 ms
CUJ2K 1.1	≈25 ms	≈49 ms	≈166 ms
<b>CUDA GPU</b>	<b>7.3 ms</b>	<b>8.1 ms</b>	<b>17.6 ms</b>



# JPEG – schéma zpracování



# JPEG – Kódování entropie

- RLE + Huffman
  - 16 nul se kóduje jako speciální symbol 0xF0
  - zbytek nul se kóduje s nenulovými symboly Huffmanem
- Huffman: neadaptivní kodér
  - výběr symbolů se dá dělat plně paralelně
  - při znalosti velikosti symbolů je možné určit jejich umístění ve výsledku – prefix sum

# Paralení RLE

