Daniel Hernández Juárez, Alejandro Chacón, Antonio Espinosa, David Vázquez, Juan Carlos Moure and Antonio M. López







Embedded real-time stereo estimation via Semi-Global Matching on the GPU

> Computer Architecture & Operating Systems Department (CAOS) **Autonomous University of Barcelona**



Centre de Visió per Computador



Motivation



- Driver Assistance: Needs distance of each object (to brake before collision)

• Embedded system: Real-time (>= 20 Frames Per Second) & low energy consumption



Stereo Vision: cameras



- Infer depth of each pixel using two cameras (stereo pair)
- Reconstruct 3D world







Why GPU?



- More GFLOPS, more Memory Bandwidth & low Energy consumption
- Needs Massive Parallelism: found in many Computer Vision algorithms

th & low Energy consumption any Computer Vision algorithms



- 2 images taken from cameras
- accelerators (GPUs) to achieve real-time and low energy consumption
- Propose and apply methodology for developing massively parallel algorithms; for efficient implementation

Identify and develop appropriate algorithms to reconstruct 3D world of objects from

Design and implement massively parallel schemes and data layouts for execution in

identify most common parallel execution patterns and generalize adequate strategies





Stereo matching & Semi-Global Matching (SGM)

Results

Index

Motivation

Massively-Parallel Algorithm Design

Conclusions & Future work



Stereo Vision Pipeline



Camera Calibration Focal length, image center, Rotation & translation, ...

Image Rectification Make images parallel & remove lens distortion Stereo Matching Compare pixels in both images to infer disparity **3D reconstruction** *Extract 3D points from disparity*



Stereo Vision: Overview



- **Disparity:** distance between the x position in the left and right image Near objects = higher disparity Far objects = lower disparity
- Epipolar geometry limits the search to 1 dimension



Stereo Matching: Find pixel along epipolar line

Left Image = Base Image



У

X

Right Image = Matching Image



Problem:

- Find pixel $p_{x,y}$ from Base Image (Left) into Matching Image (Right)
 - Only at the left along epipolar line



Stereo Matching: General Scheme

Left Image = Base Image



p(x): Preprocessing function Laplacian of Gaussian, Bilateral filtering, Rank transform, **Census transform**

- **Support window:** reduces ambiguity
- **Preprocessing + Cost function:** compensates for photometric distortions

Right Image = Matching Image



Stereo Matching: Computing Disparity

Left Image = Base Image



100	85	95	12	69	46	32	1
d=0							

Right Image = Matching Image



Stereo Matching: Local algorithm

Local method

disparity d is index with minimum cost







Consider global information

Decision depends on neighbors

99

. /

A better candidate?





Stereo Matching: Global algorithm

Global algorithm

Key assumption: Changes in distance are smooth

Define energy function over images and minimize it

- Local Cost $E(D) = E_{data}(D) + E_{smooth}(D)$ Consider Considers neighbor pixels

- **Problem:** Global minimization in 2D is NP-Complete
 - But, there is a faster option ...



Stereo Matching: Semi-Global Matching (SGM)



SGM:

- aggregate penalties from all directions and minimize
- similar accuracy of global algorithms, but lower computational complexity

8 Paths from all Directions r

• add smoothing penalties along several directions in 1D (dynamic programming method)





- Key idea: Penalize abrupt changes
- Penalty for small changes: considers slanted or curved surfaces
- Penalty for bigger changes: considers boundaries of real objects



Semi-global matching: Visual Example

Direction \rightarrow







Visual Example: effect of using different path directions

Direction \rightarrow



Direction ↓



Direction ←

Direction 1



Visual Example: Effect of Aggregation



Aggregation on all directions provides more robust results

→ + ↓ + ↑ + ←







Stereo matching & Semi-Global Matching (SGM)

Results

Index

Motivation

Massively-Parallel Algorithm Design

Conclusions & Future work



Massively Parallel Programming: Methodology

- 1. Identify massive parallelism & parallel patterns (Map, Reduce, Recurrence ...)
- 2. Propose work distribution into blocks of threads
- 3. Analyze global memory access
 - estimate arithmetic intensity
 - consider data reuse strategies
- 4. Design language-independent **pseudocode**
- 5. Translate to CUDA (or OpenCL, OpenACC ...)





Pipeline: Host and Device mapping



Complete stereo matching pipeline implemented on GPU (accelerator)



CSCT Transform



- Rationale: invariant to local intensity changes & tolerant to outliers
- Encode local neighborhood around each pixel into bit-vector 9x7 window (31-bit result)

es & tolerant to outliers cel into bit-vector







CSCT Transform: task distribution & performance analysis



- 1. 2D Stencil-pattern: parallelism in x & y axis
- 2. Naïve scheme: embarrasingly parallel (per-thread)



Parallelization Scheme	Naïve
Thread Parallelism (per image)	W×H
Compute Work per thread	62 ops
Total Global Data Read (Bytes)	62×W×H
Total Global Data Written (Bytes)	4×W×H

<u>Arithmetic Intensity</u>: 62 ops / (62+4) Bytes

xis per-threa







CSCT Transform: task distribution & performance analysis (2)



- 1. 2D Stencil-pattern: parallelism in x & y axis
- 2. 2D input data tile: cooperative read (per CTA) Data reuse from shared memory (per thread)

Parallelization Scheme	Naïve	2D-tiled
Thread Parallelism (per image)	W×H	W×H
Compute Work per thread	62 ops	62 ops
Total Global Data Read (Bytes)	62×W×H	≈ 1.5 ×₩×ŀ
Total Global Data Written (Bytes)	4×W×H	4×₩×H

Arithmetic Intensity: 62 ops / 5,5 Bytes







CSCT Transform: pseudocode

Algorithm 1: CSCT: 2D-tiled, read-cooperative parallel scheme input : I[H][W], H, Woutput: CSCT[H][W]1 parallel for y=0 to H step WarpSize do parallel for x=0 to W step WarpSize do $\mathbf{2}$ CTA parallel for yCTA, xCTA = (0,0) to (WarpSize, WarpSize) do $\mathbf{3}$ copy $(WarpSize + 8) \times (WarpSize + 6)$ tile of I[[] into SharedI[][]; 4 CTA Barrier Synchronization; 5 $CSCT[y+yCTA][x+xCTA] = CSCT_{9.7}(SharedI, xCTA, yCTA);$ 6











1. 2D to 3D pattern: parallelism in x, y & d axis 2. Naïve scheme: embarrasingly parallel (per-thread)

Naïve
W×H×D
2 ops
8×W×H×D
W×H×D

Arithmetic Intensity: 1 op / (8+1) Bytes









- 1. 2D to 3D pattern: parallelism in x, y & d axis
- 2. Increase work per thread to D output results (reduce parallelism) 1D input data tile: cooperative read (per CTA) Data reuse from shared memory (per thread)

Parallelization Scheme	Naïve	1D-t
Thread Parallelism	W×H×D	W>
Compute Work per thread	2 ops	2D
Total Global Data Read (Bytes)	8×W×H×D	12×\
Total Global Data Written (Bytes)	W×H×D	W×ŀ

Arithmetic Intensity: 2D ops / (12+D) Bytes









SGM Cost Aggregation: path direction 4



- 1. Parallelism in x axis + Recurrent pattern in y axis + 3-1 stencil and minimum reduction in d axis
- 2. 1D CTA arrangement: stencil communication (SharedM) & cooperative minimum reduction

Parallelization Scheme	Recurrent
Thread Parallelism	W×D
Compute Work per thread	k×H ops
Total Global Data Read (Bytes)	H×W×D
Total Global Data Written (Bytes)	H×W×D











SGM aggregation and Disparity Computation



1. 3D addition to 3D: parallelism in x, y and d axis 3D to 2D (d axis) reduction pattern (minimum)









SGM aggregation and Disparity Computation (2)



- ND: Number of path Directions
- 1. 3D to 2D (d axis) reduction pattern: parallelism in x & y axis
- 2. 1D CTA arrangement: cooperative minimum reduction

Parallelization Scheme	reduction
Thread Parallelism	W×H×D
Compute Work per thread	ND ops
Total Global Data Read (Bytes)	ND×H×W×D
Total Global Data Written (Bytes)	H×W











- Fused: Last path + Disparity Computation (1.35x speed up)
- Fused: Matching Cost + horizontal path directions (1.13x speed up)
- **CTA-to-Warp** conversion (Kepler architecture and newer):
 - avoids shared memory by using fast register-to-register communication (shuffle)
 - reduces instruction count and increases instruction parallelism
- Vectorize Cost Aggregation inner loop: compute 4 costs at the same time (3× speed up)

Additional Optimizations









Stereo matching & Semi-Global Matching (SGM)

Results

Index

Motivation

Massively-Parallel Algorithm Design

Conclusions & Future work



Experimental Methodology



NVIDIA Tegra X1 Cores: 256 Peak Bandwidth: 25.6 GB/s **TDP 10 W**

- Metrics:
 - Performance= Frames Per Second (fps)
 - Energy Efficiency= fps / Watt
 - Detailed performance per kernel: Instruction Count, Memory Bandwidth & IPC
- Images stored in memory



NVIDIA Titan X

Cores: 3072

Peak Bandwidth: 336.5 GB/s TDP 250 W



Results: Accuracy & Performance



 Accuracy reduced very slightly with 4 SGM path directions with respect to 8 SGM path directions • Real-time on Tegra X1 is achieved for 4 SGM path directions on images of 1280x480 pixels







- Titan X GPU provides more than 10x the performance of Tegra X1
- Tegra X1 offers more than 2x





Stereo matching & Semi-Global Matching (SGM)

Results

Index

Motivation

Massively-Parallel Algorithm Design

Conclusions & Future work



Conclusions & Future work

- real-time stereo
- We have proposed baseline parallel schemes and data layouts that follow general optimization rules
- Evaluate performance on the new embedded Pascal GPUs with larger images and a higher number of disparity levels
- Include and evaluate well-known post-filtering algorithms:
 - Left-Right Consistency Check
 - Subpixel Estimation
 - Adaptive P2

Low-consumption embedded GPU systems (Tegra X1) are capable of attaining



Thank you

Daniel Hernández Juárez, Alejandro Chacón, Antonio Espinosa, David Vázquez, Juan Carlos Moure and Antonio M. López

More information:

www.cvc.uab.es/people/dhernandez

Autonomous University of Barcelona







