

Video-Rate Stereo Matching using MRF TRW-S Inference on a Hybrid CPU+FPGA Computing Platform

Jungwook Choi and Rob A. Rutenbar {jchoi67, rutenbar}@illinois.edu ISFPGA 2013, Feb. 12 2013



Overview

- Overall goal: Explore probabilistic inference methods in custom hardware
- This work: Fast & high quality stereo matching via Markov random field (MRF) inference in FPGA
- Three key contributions:
 - CPU-FPGA partitioning of stereo matching procedure
 - Functional level pipelining across CPU and FPGAs
 - Frame level parallelization to process multi-frames in parallel
- Results: 15~40 frame/sec for QVGA stereo matching

Fast & High Quality Stereo Matching

Variety of real-world apps demand
Fast & High Quality stereo matching



Pedestrian detection for autonomous vehicle^{*} Laptop control using gestures^{**} (M. Gavrila, IJCV07, etc.) (TYZX G3 EVS)

* http://www.gavrila.net/Research/Pedestrian_Detection/pedestrian_detection.html ** http://www.tyzx.com/technology/gallery.html



MRF Inference for High Quality Stereo Matching





MRF Inference for High Quality Stereo Matching





Custom HW Accelerator for Fast Stereo Matching

- HW impl. for computation intensive MRF infer.
- Custom HW accel. is 2~4x faster than GPU impl.
 - Thanks to fully pipelined data path and streaming data access



]



Hybrid CPU+FPGA Platform

- Our platform: Convey HC-1
 - CPU-FPGA cache-coherent virtual memory system
 - Max memory BW: 1Kbit/cycle(~20GB/sec)/FPGA (runs @150MHz)
 - Non-blocking FPGA function call



 \mathbb{I}



CPU-FPGA Partitioning of Stereo Matching Procedure











CPU-FPGA Partitioning of Stereo Matching Procedure





Initialize	■ 3.2% <u>CPU</u>	+FPGA (58.5msec)
For (idx=0; idx <numframe; idx++)<="" td=""><td></td><td>RD_IMGS</td></numframe;>		RD_IMGS
Get frame images (RD IMGS)	18.6%	
Compute MRF costs (COM_COST)	38.1%	MRF_INFER
Run MRF inference (MRF_INFER)	32.3%	GET_DEPTH
Obtain depth result (GET_DEPTH) Produce depth map (WR_DEPTH)	52.570	WR_DEPTH
End	7.7	%



Video-Rate Stereo Matching CPU+FPGA System





FPGA

Function Level Pipelining

- Apply SW pipelining technique to overlap execution time of CPU and FPGA functions.
- While FPGA functions work on *current* frame, CPU functions process previous and next frames in Steady State

Prolog

Initialize For (idx=0; idx<NumFrame; idx++) Begin





Frame Level Parallelization



ECE ILLINOIS

Function Level Pipelining + Frame Level Parallelization



- If T_{CPU} < T_{FPGA} pipelining hides CPU execution time
- Frame level par. reduces per-frame execution time of FPGA functions
- But CPU functions are serialized

11/17

→No benefit!

Performance Result: Single Frame

Tsukuba (384x288,16)	Real-time BP [*] [Yang 2006]	Tile-based BP** [Liang 2011]	Fast BP ^{***} [Xiang 2012]	This work
GPU	NVIDIA GeForce 7900 GTX	NVIDIA GeForce 8800 GTS	NVIDIA GeForce GTX 260	N/A
# Iteration	(4 scales) = (5,5,10,2)	(B, T _I , T _O) = (12, 20, 5)	(3 scales) = (9,6,2)	T _O = 5
Time (msec)	80.8	97.3	61.4	26.10
Min. Energy	N/A	396,953	N/A	393,434

* Q. Yang, et al., "Real-time global stereo matching using hierarchical belief propagation," *BMVC*, 2006. ** Liang, et al., "Hardware-Efficient Belief Propagation," *IEEE Trans. Circ. Syst. Video Tech*, May 2011. *** X. Xiang, et al., "Real-time stereo matching based on fast belief propagation," *MACH VISION APPL*, 2012

I



Performance Result: Video Sample

- Tasks: Flower*(360x262,16), Ume*(320x240,16)
 - Require 80+ iterations to remove defect at the boundary region



Ume



I



Performance Result: Function Level Pipelining

- Function level pipelining hides CPU overhead
 - As frames processed, more CPU fn are overlapped by FPGA fn

Portion of CPU fn execution time (possible reduction)

Actual reduction in execution time by pipelining





Performance Result: Frame Level Parallelization

- Frame level parallelization speed-ups FPGA functions
 - CPU fn exec time is unchanged, but minor \rightarrow hidden by pipelin.



Number of frames processed in parallel

]



Performance Result: Video Demo

- Function level pipelining + frame level parallelization
 - 3-frame parallel stereo matching, 80 iterations for each inference
 - Speed: Flower 12.3 frame/sec, Ume 15.0 frame/sec

— Tenolor@courtshit TenutelTworkshote/citilTworkshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTheshotepathareaultscandinglicitieTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenuteTenute

Getting signature



Conclusion

- Video-rate of high quality stereo matching is done using MRF inference on a hybrid platform
- CPU-FPGA partitioned functions of stereo matching procedure are pipelined and frame-level parallelized
- Performance results:
 - Simple QVGA single frame: 40 frame/sec
 - Challenging QVGA video sample: 15 frame/sec
- Fastest ever implementation of TRW-S MRF inference

1



Thanks!

Any Questions?