

IMPLEMENTING ADAPTIVE BEAMFORMING ON THE GPU FOR REAL TIME ULTRASOUND IMAGING

J. P. Åsen¹, J. I. Buskenes², C.-I. C. Nilsen², A. Austeng² and S. Holm^{1,2}

¹Mi Lab, Norwegian University of Science and Technology, Trondheim, Norway

²Department of Informatics, University of Oslo, Oslo, Norway

ABSTRACT

1. INTRODUCTION

Introduction to Capon adaptive beamforming, ultrasound imaging and GPU/CUDA.

Delay and sum beamformer (DAS):

$$z[n] = \sum_{m=0}^{M-1} \mathbf{w}_m^* \mathbf{x}_m[n - \Delta_m[n]] = \mathbf{w}^H \mathbf{x}[n], \quad (1)$$

where M is the number of elements.

Capon weights are found by solving the following minimization problem:

$$\min_{\mathbf{w}} E\{|z[n]|^2\} = \mathbf{w}^H \mathbf{R} \mathbf{w} \quad (2)$$

$$\text{subjected to } \mathbf{w}^H \mathbf{1} = 1, \quad (3)$$

where $\mathbf{R} = E\{\mathbf{x}[n]\mathbf{x}[n]^H\}$.

The solution to this minimization problem is

$$\mathbf{w}[n] = \frac{\mathbf{R}[n]^{-1} \mathbf{a}}{\mathbf{a}^H \mathbf{R}[n]^{-1} \mathbf{a}}, \quad (4)$$

where $\mathbf{a} = \mathbf{1}$ when \mathbf{x} are pre-delayed and $\mathbf{w} \in \mathbb{C}^M$.

Estimation of $\mathbf{R}[n]$: $\hat{\mathbf{R}}[n] = \mathbf{x}[n]\mathbf{x}[n]^H$.

In order to get a well conditioned $\hat{\mathbf{R}}$, avoid signal cancellation and to get DAS-like speckle, $\hat{\mathbf{R}}$ has to be averaged over $L \leq M/2$ subarrays and N_{avg} time samples.

$$\hat{\mathbf{R}} = \frac{1}{(2N_{avg} + 1)K} \sum_{n'=n-N_{avg}}^{n+N_{avg}} \sum_{l=0}^K \mathbf{x}_l[n'] \mathbf{x}_l[n']^H, \quad (5)$$

where $K = M - L + 1$, $Y_{avg} \sim \tau/T_s$ (τ is the imaging pulse length and T_s is the sampling period.) and \mathbf{x}_l is the l^{th} subarray $[x_l[n], \dots, x_{l+L}[n]]$.

Building and calculating the inverse of $\hat{\mathbf{R}}$ is computational demanding, limiting the methods accessibility. Following the innovation in GPU computing, this has now started to change.

2. METHOD

2.1. Building the Sample Covariance Matrix

Building $\hat{\mathbf{R}}$ using a sliding window approach across K and N_{avg} . We have a limited amount of fast, near-core memory on the GPU. On NVIDIA architecture this memory is known as shared memory, and is restricted to 48KB per compute block. The maximum block size of 32x32, further restrict how large arrays we can handle inside a single compute block. We therefore restrict L to a maxima of 32 elements. Since $L \leq M/2$, $M \leq 64$. We can not afford to hold the full $\hat{\mathbf{R}}$ ($M = L$) when $\hat{\mathbf{R}} > 32$, because of limited amount of memory. ... Explain in detail how `buildR` is implemented in a kernel. ... Give more details on how much shared memory we have, and how it can be divided between compute blocks occupying one stream multiprocessor (SM).

2.2. Solving Multiple Small Linear Systems

Most focus on large matrices in the literature. We have used an unreleased GPU implementation of Gauss Jordan (GJ) elimination (by NVIDIA) to solve $\mathbf{R}\mathbf{b} = \mathbf{1}$, where $\mathbf{b} = \mathbf{R}^{-1}\mathbf{1}$. Can include details on GPU implementation of $\mathbf{U}^H \mathbf{D} \mathbf{U}$, but this has not proved to be faster than NVIDIA's GJ implementation. However, the complexity for solving with $\mathbf{U}^H \mathbf{D} \mathbf{U}$ decomposition is supposed to be 1/3 of GJ.

3. RESULTS

Present images before and after Capon weights has been applied. Comment on the effect of selecting different parameters. Present graph showing running times for Cython-Capon v.s. CUDA-Capon for different choices of parameters.

(A prerequisite for this is to get the GPU-processing up and running on simulated data before Easter!! Maybe also a Matlab plug-in.)

4. DISCUSSION

Discuss results presented in Section 3.

5. CONCLUSION

...