

# A Study on Convolution Using Half-Precision Floating-Point Numbers on GPU for Radio Astronomy Deconvolution

**Mickaël Seznec**<sup>1,3</sup>, Nicolas Gac<sup>1</sup>, André Ferrari<sup>2</sup>, François Orioux<sup>1</sup>

<sup>1</sup> L2S (CentraleSupélec, Univ. Paris-Sud, CNRS) Gif-sur-Yvette, FRANCE

<sup>2</sup> Lab. J.-L. Lagrange (Univ. Nice Sophia Antipolis, CNRS, Observatoire Côte d'Azur) Nice, France

<sup>3</sup> Thales Research & Technology, Palaiseau, France

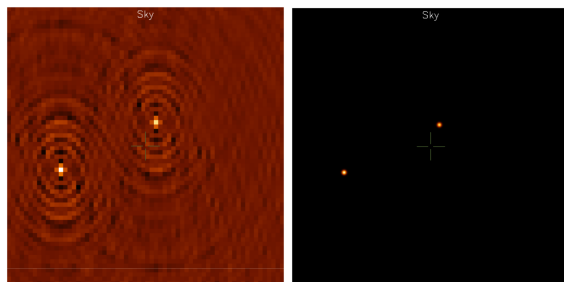
IEEE SIPS, Special Session on Computing for Radio-Astronomy  
2018/10/23

# Outline

- 1 Introduction to deconvolution
- 2 Half-floats
- 3 Application to radio astronomy
- 4 Conclusions

# Deconvolution

Trying to recover data from an acquisition...



*Source: I.M. Steward, Bielefeld University*

Helps in the analysis in many fields: from radio astronomy to microscope images.

# Modélisation

A physical measurement can be modeled as:

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \mathbf{n}$$

- $\mathbf{g} \in \mathbb{R}^M$  is the measured data
- $\mathbf{H} \in \mathbb{R}^{M \times N}$  is the linear observation model (here: a convolution)
- $\mathbf{f} \in \mathbb{R}^N$  is the ground truth, unknown
- $\mathbf{n} \in \mathbb{R}^M$  is added noise, usually Gaussian

# Regularization

$$J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \lambda\|\mathbf{D}\mathbf{f}\|^2$$

The criterion involves two terms: data fidelity and regularization to compensate the ill-posedness behavior.

In the following, we choose  $\mathbf{D} = \mathbf{I}^{M \times M}$

The solution is:

$$\hat{\mathbf{f}} = \arg \min_f J(\mathbf{f})$$

$$\hat{\mathbf{f}} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^t \mathbf{g}$$

The Hessian matrix's dimensions are too large to inverse it: we fall back on an iterative solver instead.

# The gradient descent algorithm

---

**Require:**  $H$ ,  $\lambda$ ,  $g$ ,  $\epsilon$ ,  $N$ ,  $c$

1: Set  $\mathbf{b} = H^t \mathbf{g}$  and  $\mathbf{Q} = H^t H + \lambda I$

2: Set  $\mathbf{f}^{(0)}$  and  $n \leftarrow 0$

3: **repeat**

4:      $\mathbf{k} \leftarrow H^t H \mathbf{f}^{(n)} - \mathbf{b} + 2\lambda \mathbf{f}^{(n)}$

5:      $\alpha \leftarrow \mathbf{k}^t \mathbf{k} / \mathbf{k}^t \mathbf{Q} \mathbf{k}$

▷ or  $\alpha \leftarrow c$  with  $c \leq \frac{2}{\|\mathbf{Q}\|}$

6:      $\mathbf{f}^{(n+1)} \leftarrow \mathbf{f}^{(n)} - \alpha \mathbf{k}$

7:      $n \leftarrow n + 1$

8: **until** Some criterion is met

▷ e.g.:  $n \geq N$

**return**  $\mathbf{f}^{(n)}$

---

# The gradient descent algorithm

---

**Require:**  $H$ ,  $\lambda$ ,  $g$ ,  $\epsilon$ ,  $N$ ,  $c$

1: Set  $\mathbf{b} = H^t \mathbf{g}$  and  $\mathbf{Q} = H^t H + \lambda I$

2: Set  $\mathbf{f}^{(0)}$  and  $n \leftarrow 0$

3: **repeat**

4:      $\mathbf{k} \leftarrow H^t H \mathbf{f}^{(n)} - \mathbf{b} + 2\lambda \mathbf{f}^{(n)}$

5:      $\alpha \leftarrow \mathbf{k}^t \mathbf{k} / \mathbf{k}^t \mathbf{Q} \mathbf{k}$

▷ or  $\alpha \leftarrow c$  with  $c \leq \frac{2}{\|\mathbf{Q}\|}$

6:      $\mathbf{f}^{(n+1)} \leftarrow \mathbf{f}^{(n)} - \alpha \mathbf{k}$

7:      $n \leftarrow n + 1$

8: **until** Some criterion is met

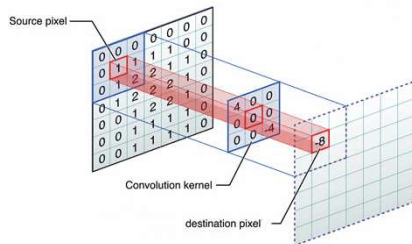
▷ e.g.:  $n \geq N$

**return**  $\mathbf{f}^{(n)}$

---

Up to 4 convolutions per loop!

## 2D convolution



Convoluting a  $N \times N$  image with a  $K \times K$  kernel can mainly be done in two ways:

- In direct space: complexity  $O(N^2 K^2)$
- In Fourier space: complexity  $O(N^3 \log(N))$

In special cases, optimization exists: separable kernel, overlap-add...



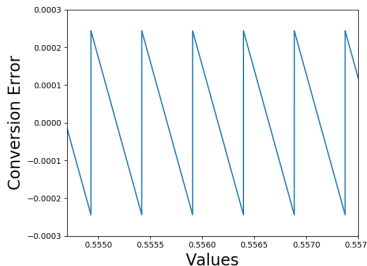
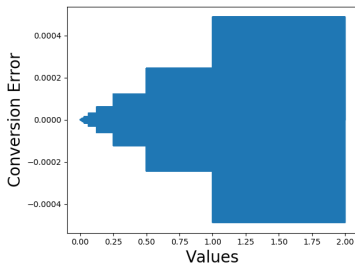
## 16-bit floats format



Figure: IEEE 754 binary16 format

- Supported on an increasing number of platforms
- Double the FLOPS on Nvidia GPUs via SIMD instructions
- Decreased memory and bandwidth usage
- Limited range (from  $10^{-8}$  to 65504) and accuracy

# From FP32 to FP16



- We suffer from quantization errors in the conversion
- The relative error stays the same

## Results on convolutions

Convolution of a  $512 \times 512$  image with a square kernel of variable size. Only the kernel execution time is measured.

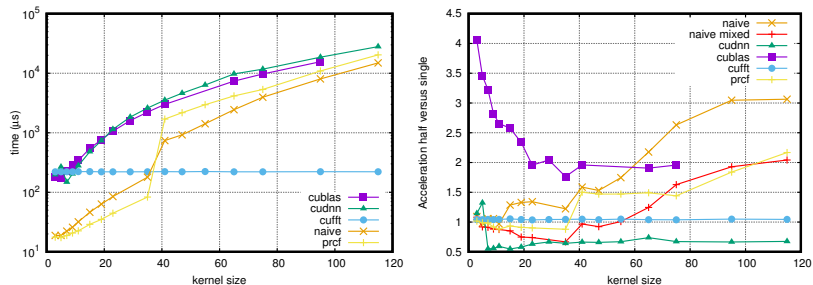


Figure: Left: results with FP32. Right: speed-up with FP16

- Naive and PRCF: own implementations.
- cuFFT looks constant because kernel is padded

# Accuracy

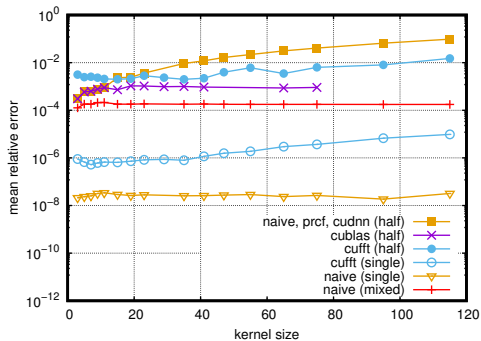
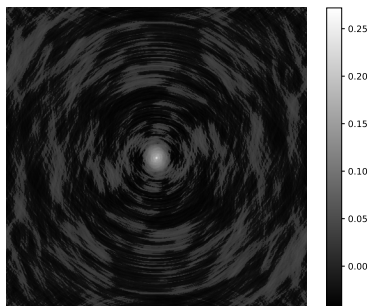


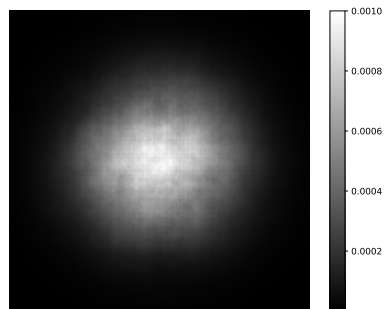
Figure: MRE compared to a reference implementation

$$MRE = \frac{1}{N} \sum \begin{cases} \left| \frac{x[i] - y[i]}{x[i]} \right|, & \text{if } x[i] \neq 0 \\ 0, & \text{otherwise} \end{cases}$$

# Application to radio astronomy



PSF (cubic root)



Sky (clipped to 0.001)

1000:1 ratio between star values and background field.

## Errors on “real” data

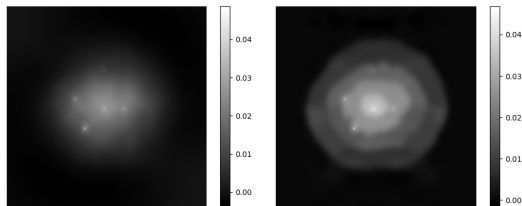
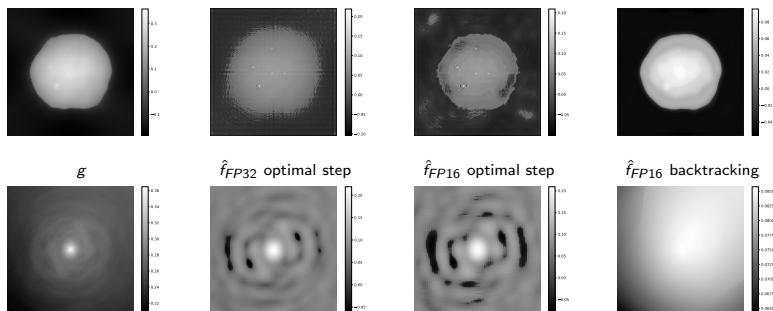


Figure: Left: convolution with FP32, right: with FP16

- Important loss of precision with a single convolution!

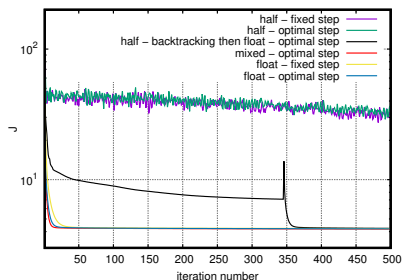
# Deconvolution Results



**Figure:** Reconstructions of  $f$  with different precisions (the cubic root is displayed for better contrast)

# Minimizing the criterion

$$J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \lambda\|\mathbf{f}\|^2$$



- Half-float only implementation does not have a satisfying behavior.
- Mixed strategy performs well.



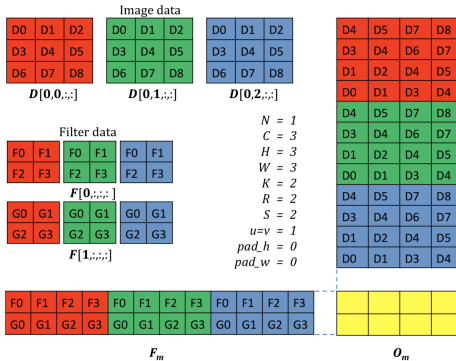
## Summary & Conclusions

- We developed a GPU convolution benchmark
- We made a CUDA/C++ implementation of a deconvolution algorithm

## Summary & Conclusions

- We developed a GPU convolution benchmark
- We made a CUDA/C++ implementation of a deconvolution algorithm
- Half floats can significantly speed-up convolutions (up to  $\times 2, \times 3$ )
- Loss of precision must be tightly controlled...
- ... but can be mitigated with the use FP32 compute / FP16 storage

Thank you!



*Sharan Chetlur et. al. 2015*