

Fast Bilateral Filter GPU implementation

Multi-Core Architectures and Programming

Gerhard Mlady, Rafael Bernardelli

Hardware/Software Co-Design, University of Erlangen-Nuremberg

July 21, 2016

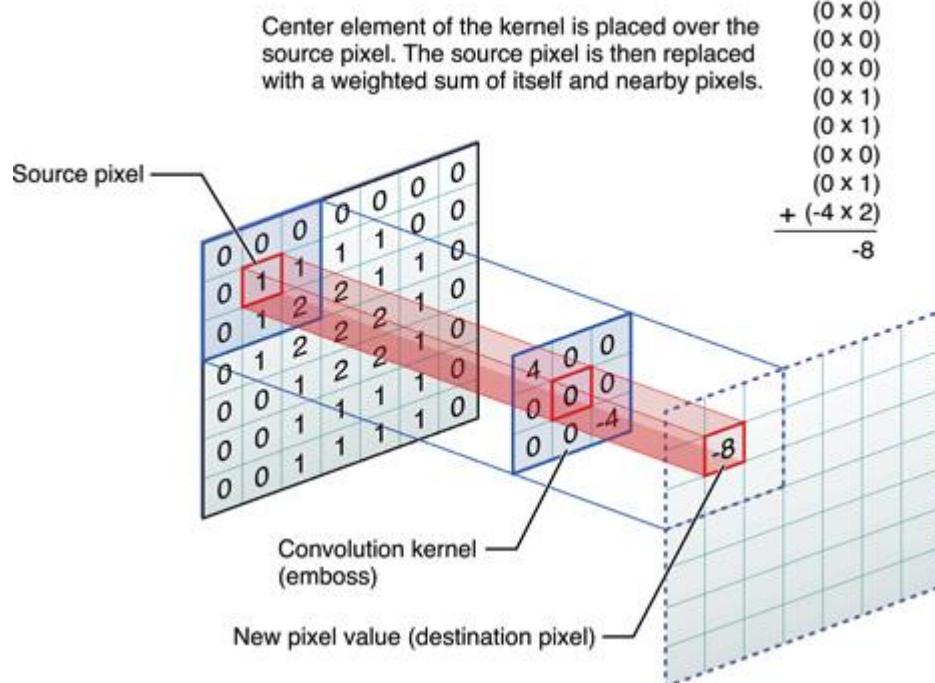


Overview

- **Fast Bilateral Filter**
- Implementation
- Benchmark

Introduction to Bilateral filter

First lets take a look on 2d convolution

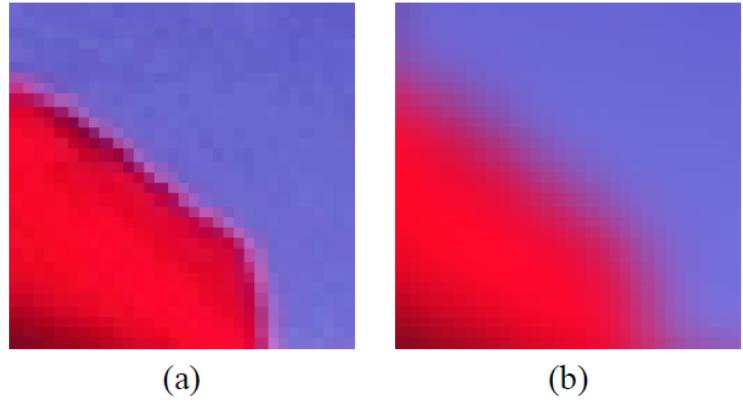


$$\begin{aligned}
 f(x) &= g(x) * k(x) \\
 f(x) &= \sum_{x' \in \mathcal{N}(n)} g(x') k(x - x') \\
 f(x) &= \mathcal{F}^{-1}\{\mathcal{F}\{g(x)\}\mathcal{F}\{k(x)\}\}
 \end{aligned}$$

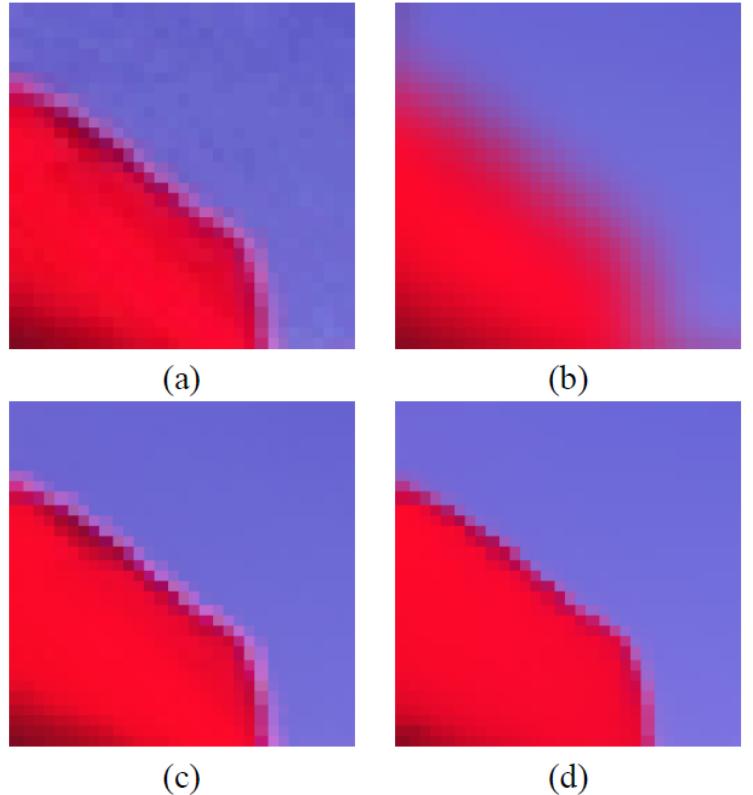
$\frac{1}{273}$

1	4	7	4	1
4	16	26	16	4
7	26	41	26	7
4	16	26	16	4
1	4	7	4	1

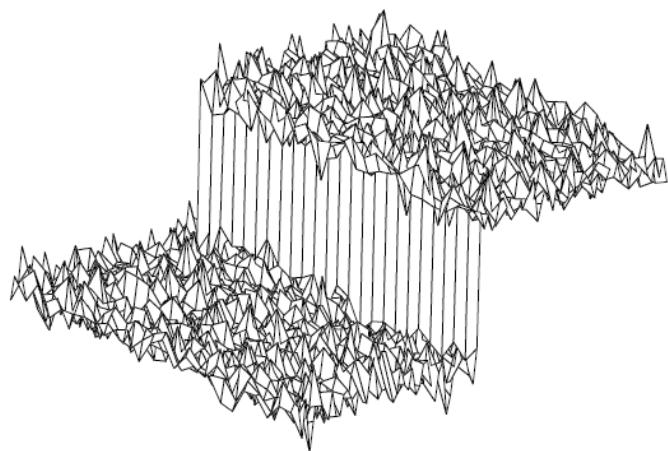
Gaussian kernel



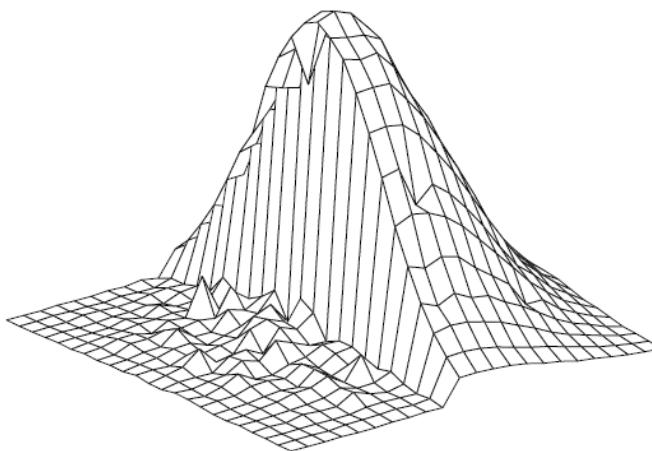
- (a) Original image
- (b) Gaussian filtered image



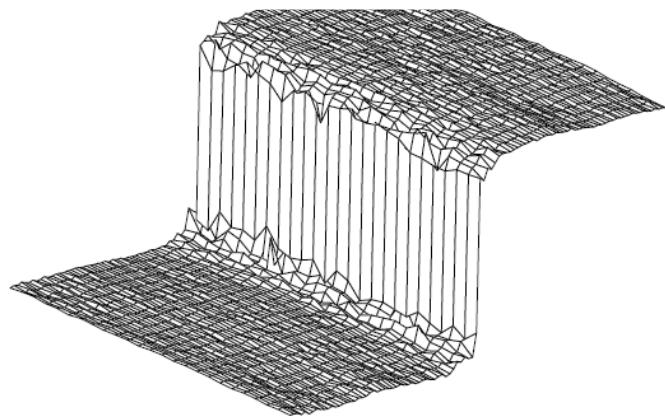
(a) Original image
(b) Gaussian filtered image
(c) Bilateral filtered image*
(d) Bilateral filtered image**



(a)



(b)



(c)

(a) Original image
(b) Bilateral kernel
(c) Result

Bilateral filter equations

$$f(\mathbf{x}) = k^{-1}(\mathbf{x}) \sum_{\mathbf{x}' \in \mathcal{N}} g(\mathbf{x}') c(\mathbf{x}, \mathbf{x}') s(g(\mathbf{x}), g(\mathbf{x}'))$$

$f(\mathbf{x})$: Output image

$g(\mathbf{x})$: Input image

Bilateral filter equations

$$f(\mathbf{x}) = k^{-1}(\mathbf{x}) \sum_{\mathbf{x}' \in \mathcal{N}} g(\mathbf{x}') c(\mathbf{x}, \mathbf{x}') s(g(\mathbf{x}), g(\mathbf{x}'))$$

$$c(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2\sigma_d^2}}$$

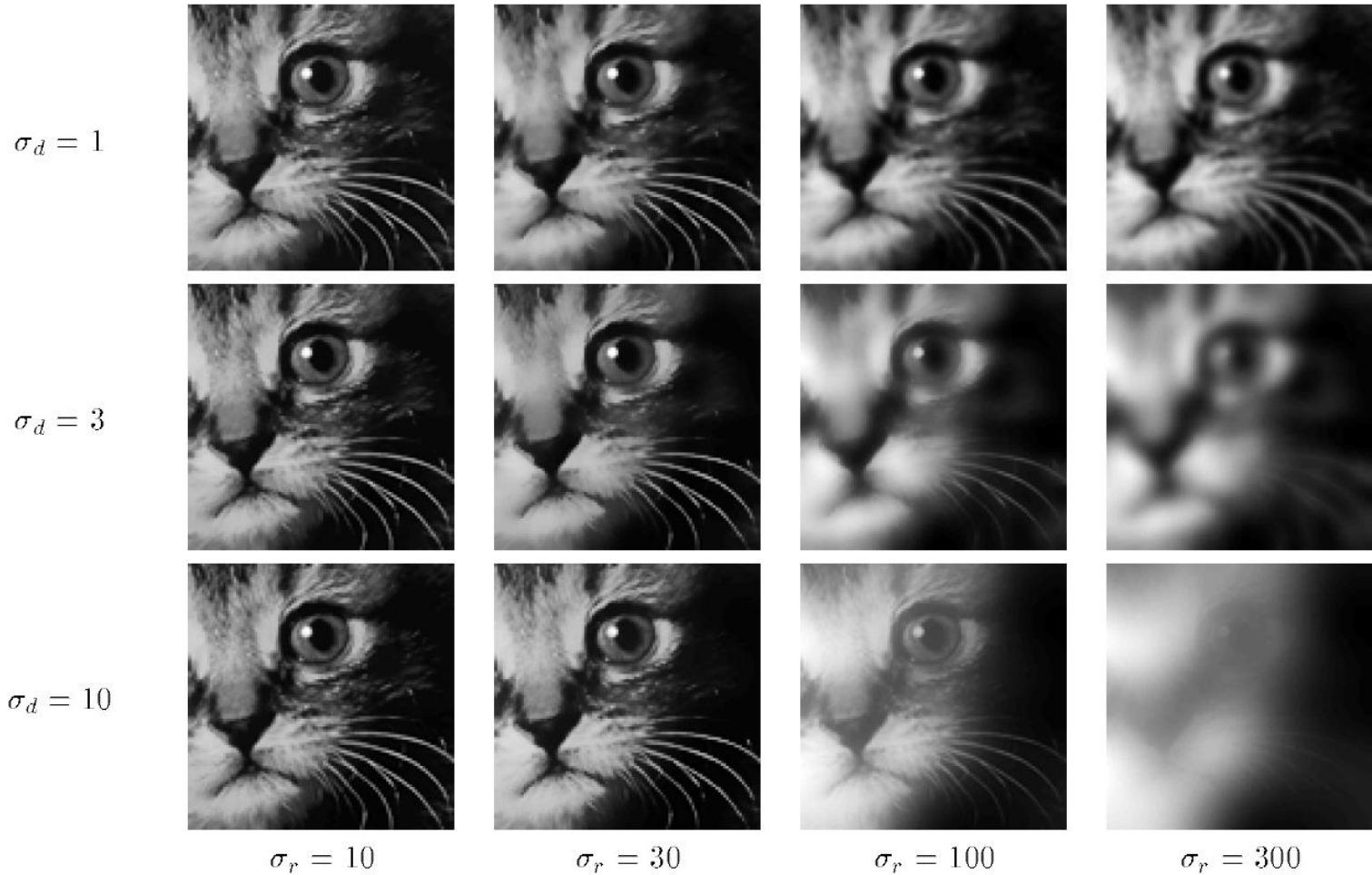
Spacial similarity function

$$s(g(\mathbf{x}), g(\mathbf{x}')) = e^{-\frac{(g(\mathbf{x})-g(\mathbf{x}'))^2}{2\sigma_r^2}}$$

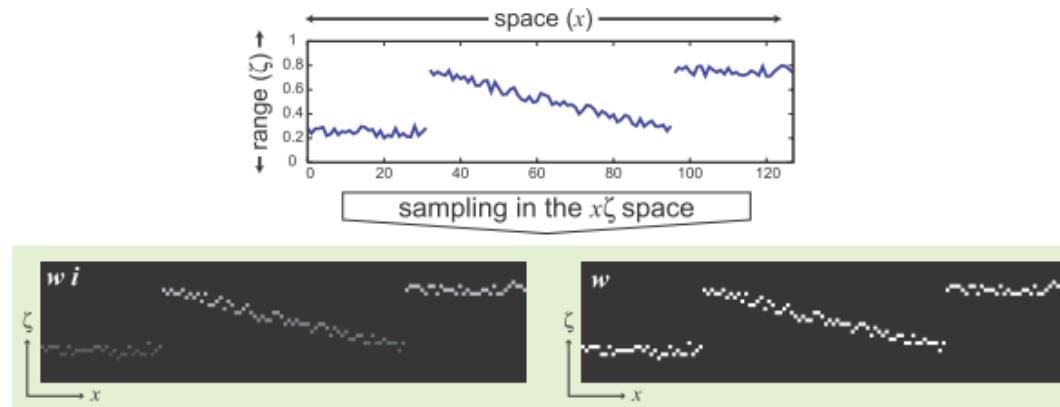
Range similarity function

$$k(\mathbf{x}) = \sum_{\mathbf{x}' \in \mathcal{N}} c(\mathbf{x}, \mathbf{x}') s(g(\mathbf{x}), g(\mathbf{x}'))$$

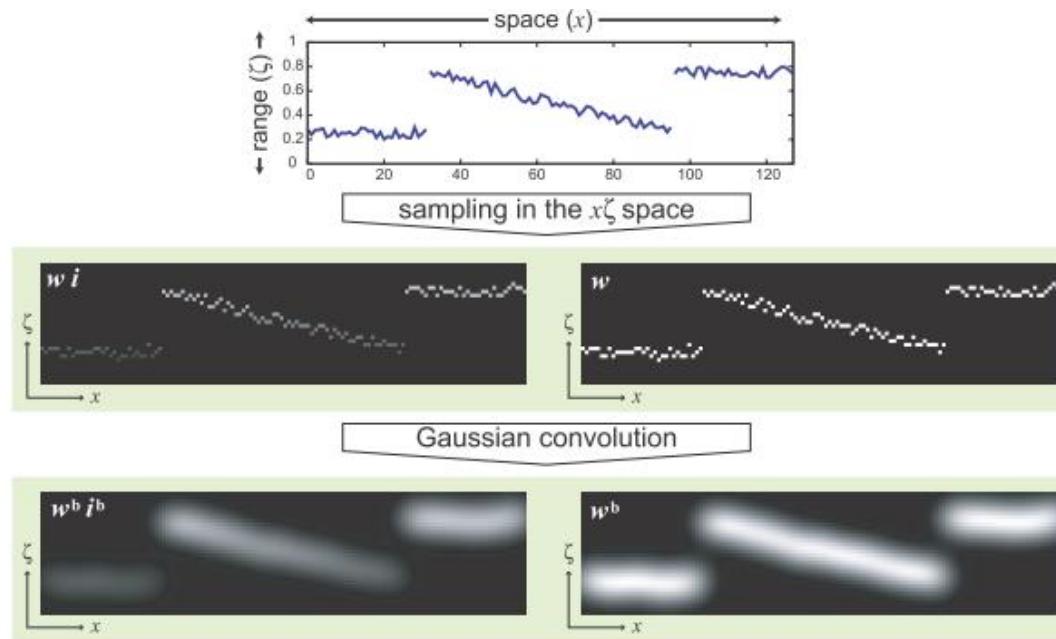
Normalizing factor



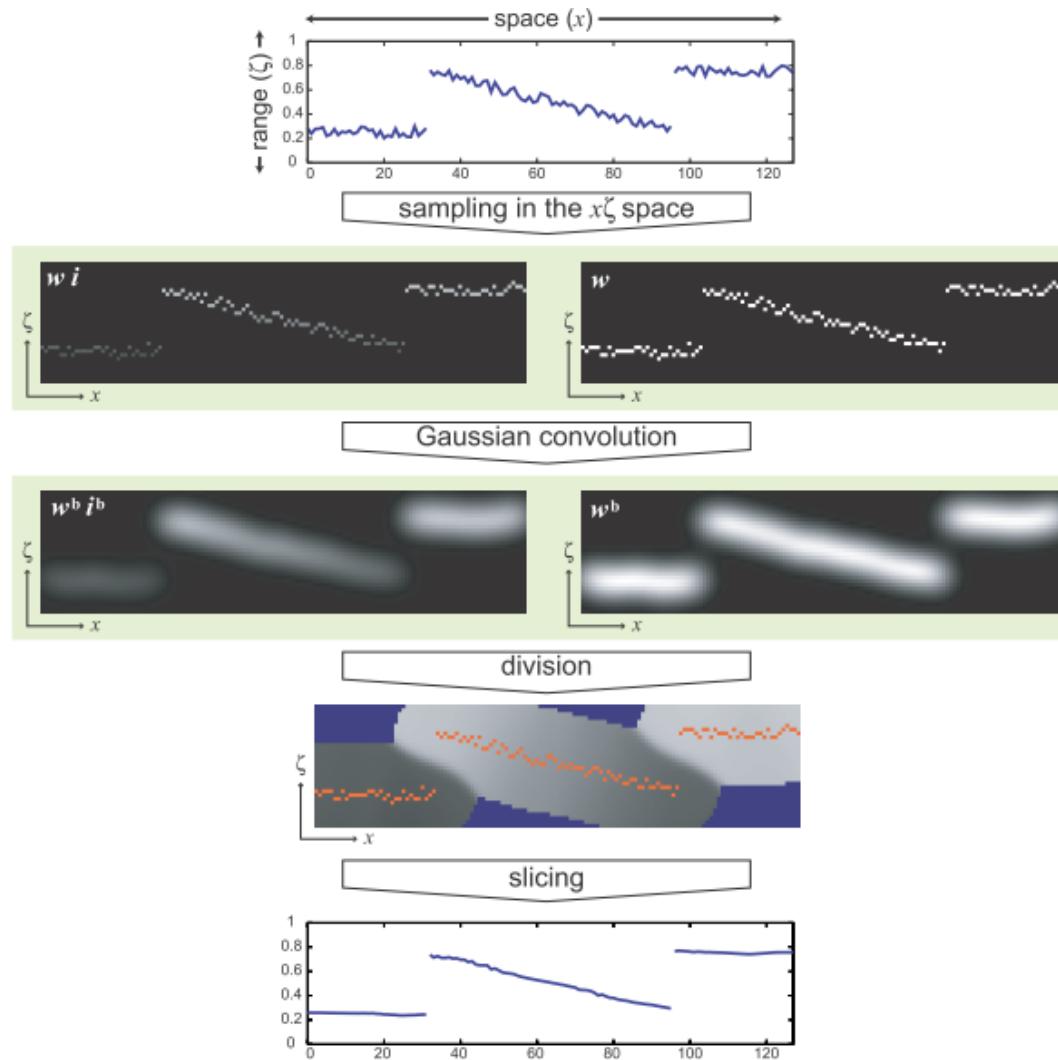
Fast Bilateral Filter



Fast Bilateral Filter



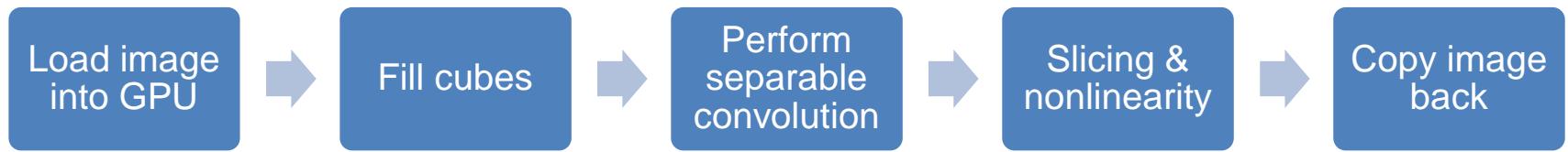
Fast Bilateral Filter



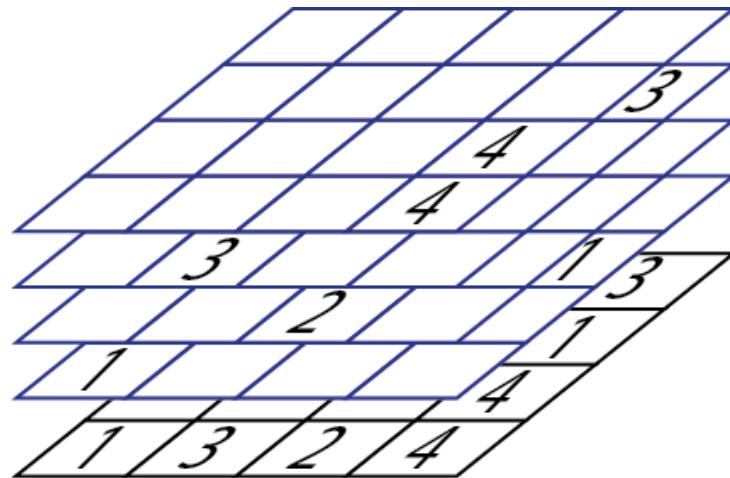
Overview

- **Fast Bilateral Filter**
- **Implementation**
- **Benchmark**

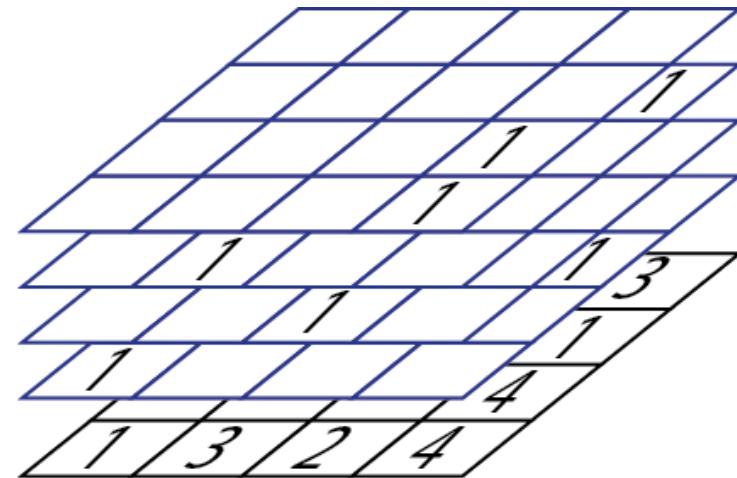
Implementation



Cube filling



WI



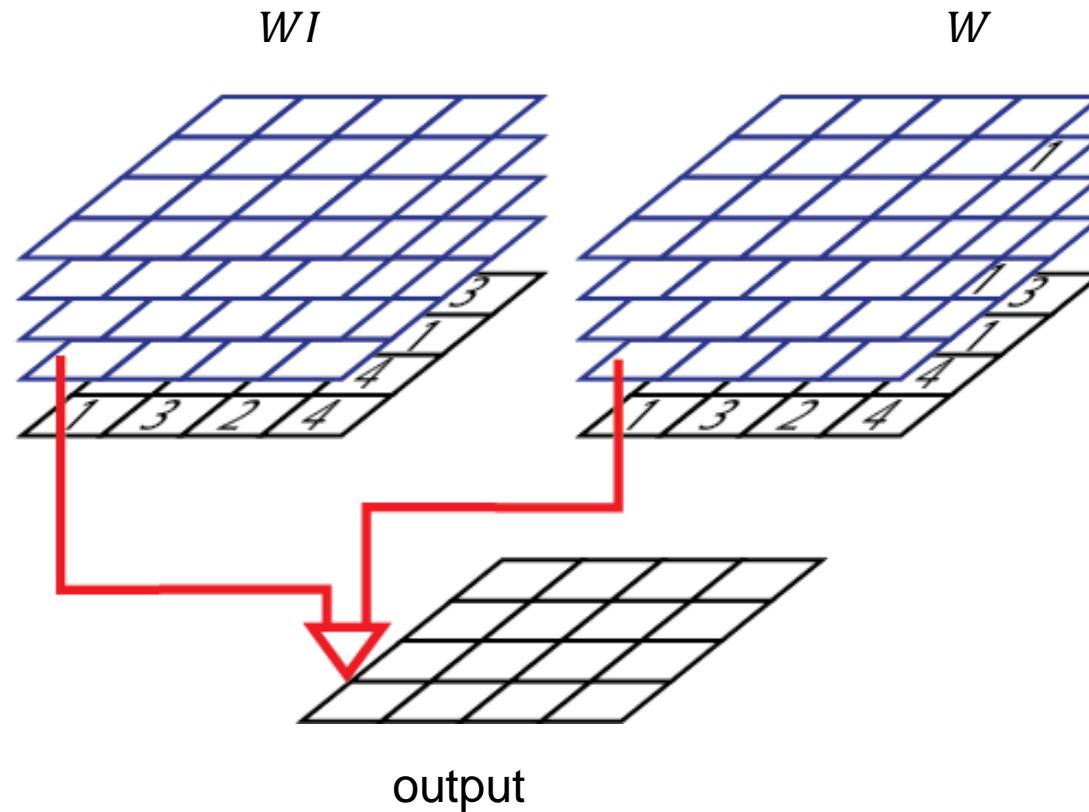
W

Perform separable convolution

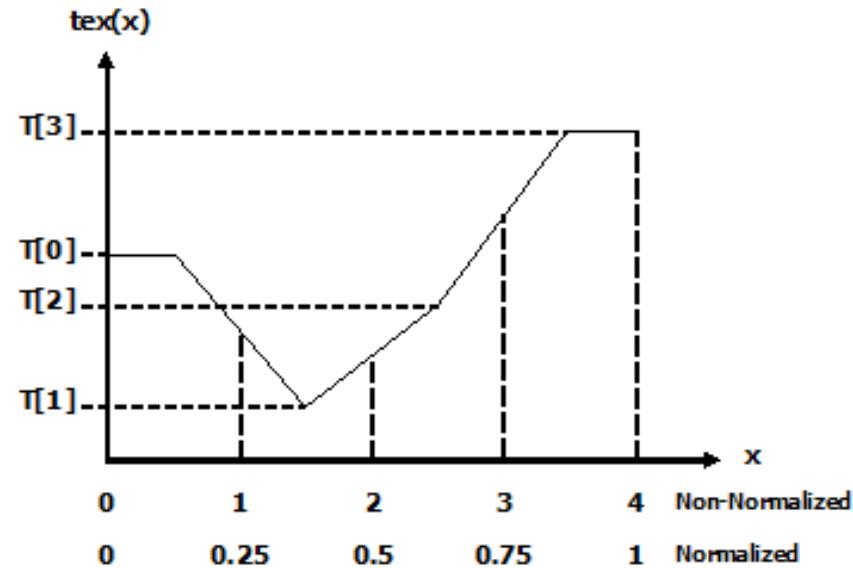
$$I * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} = I * \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} * [1 \quad 2 \quad 1]$$
$$O(NMk^2) \qquad \qquad O(2 * NMk)$$

N, M : Image dimensions
 k : convolution kernel length

Perform Slicing & Nonlinearity



Texture fetching

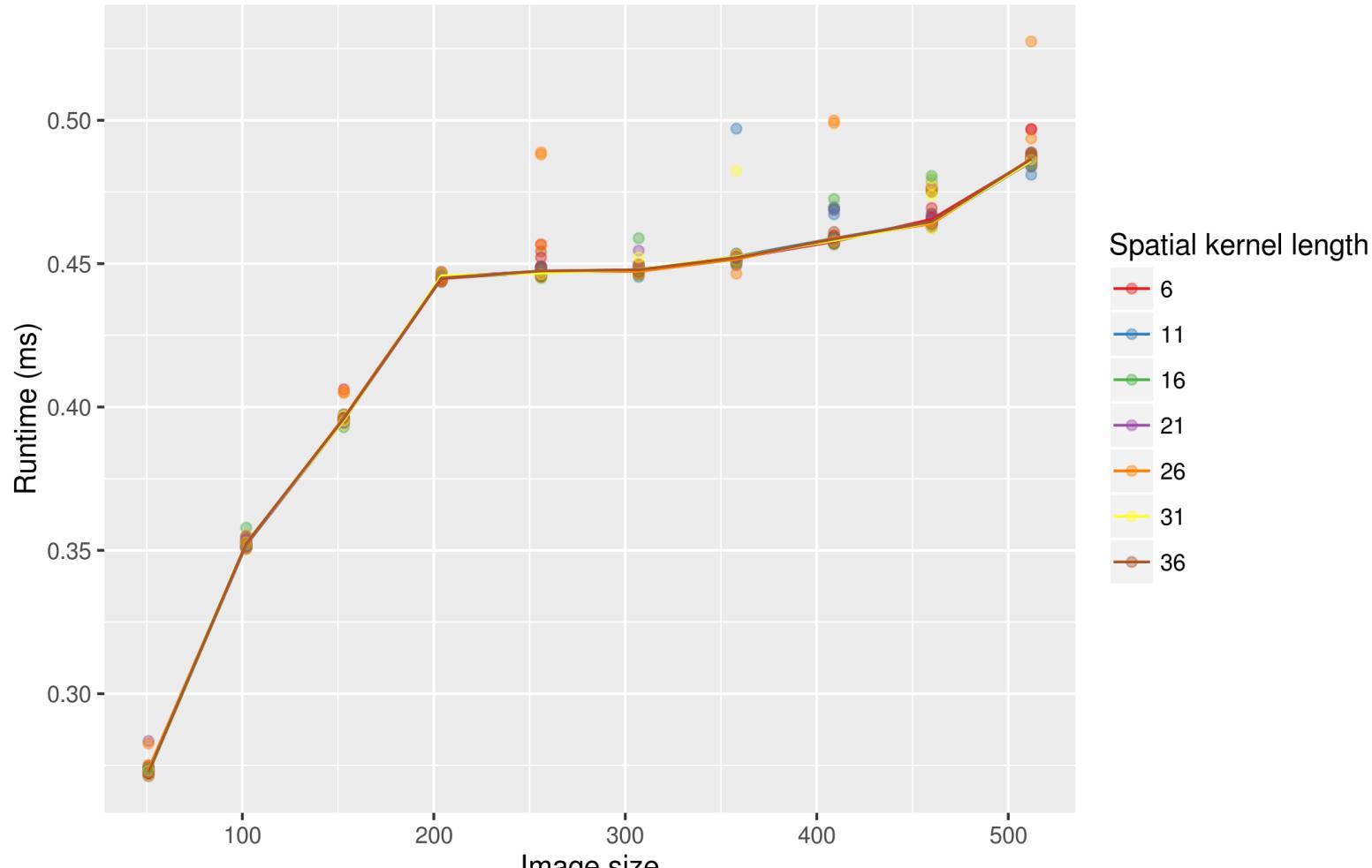


Overview

- **Fast Bilateral Filter**
- **Implementation**
- **Benchmark**

Benchmark

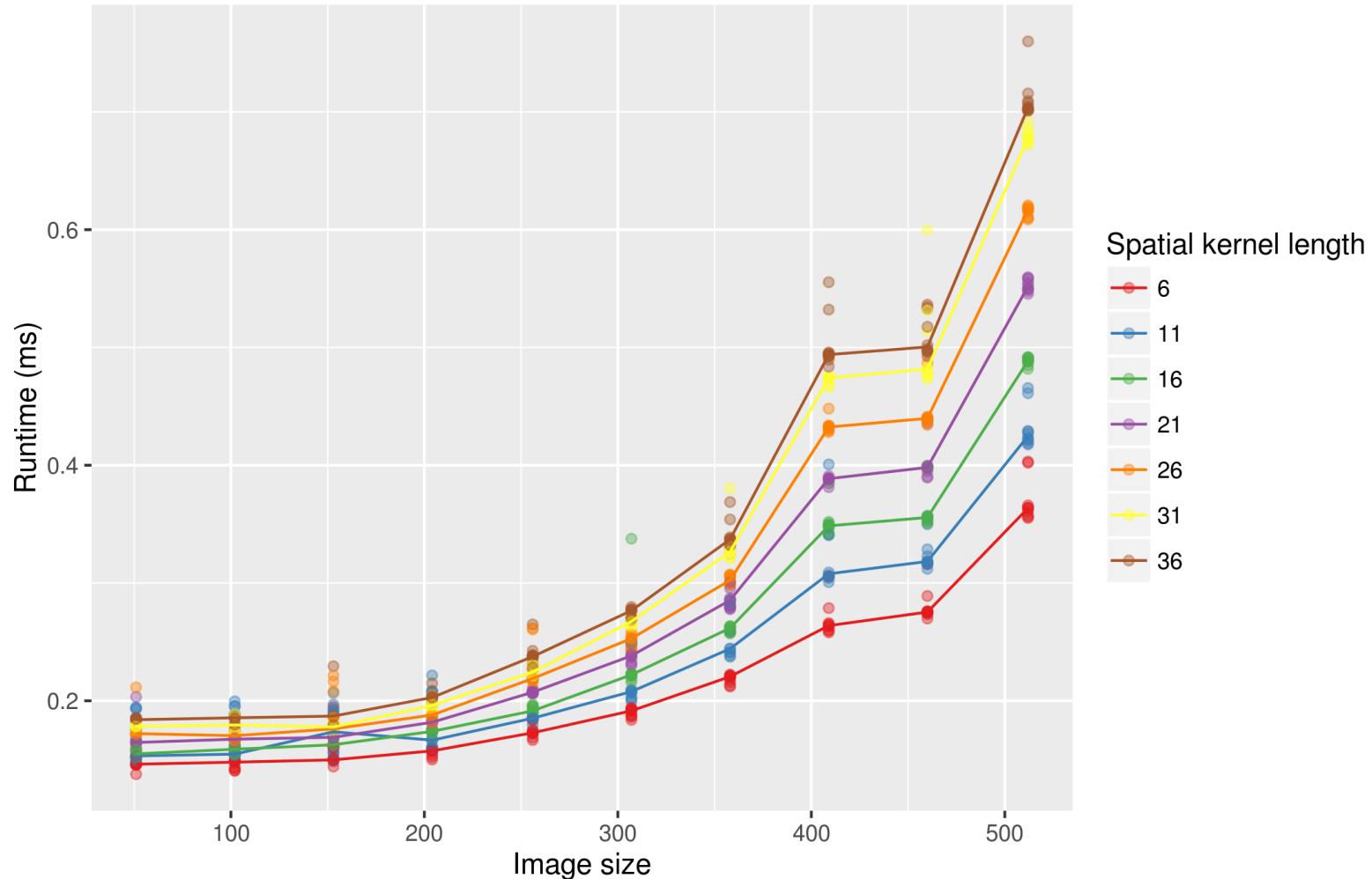
Filling on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

Benchmark

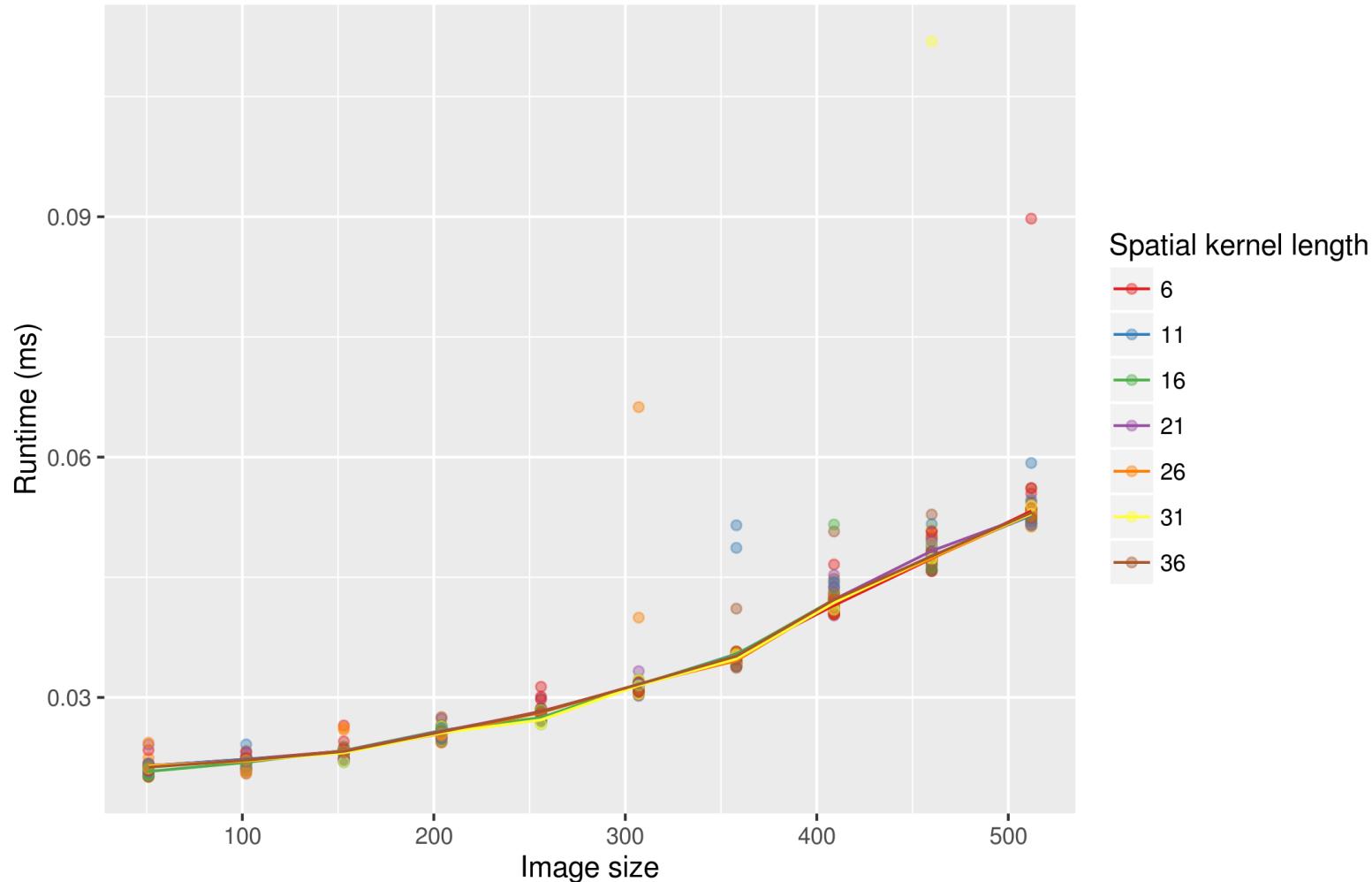
Convolution on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

Benchmark

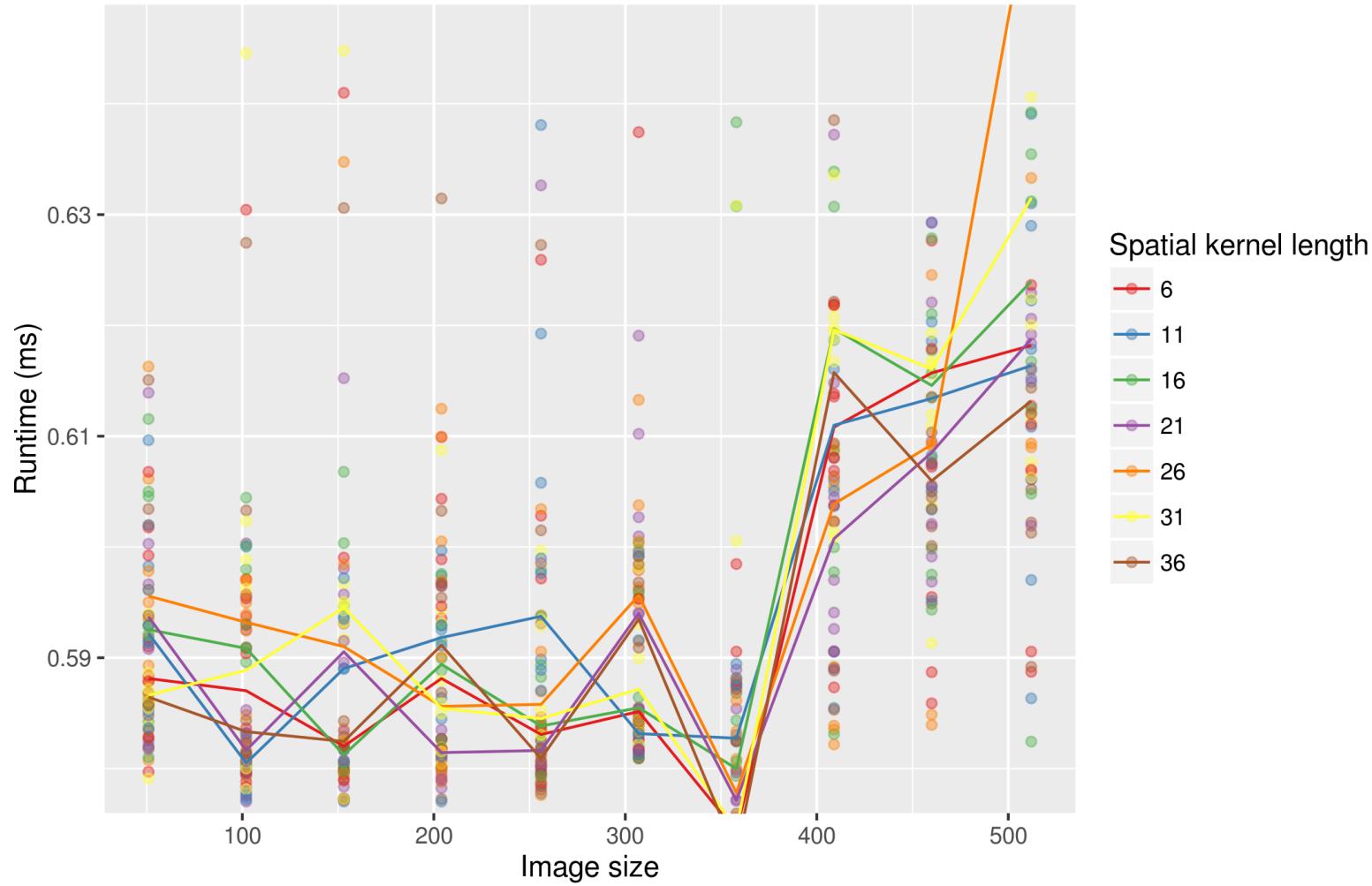
Slicing on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

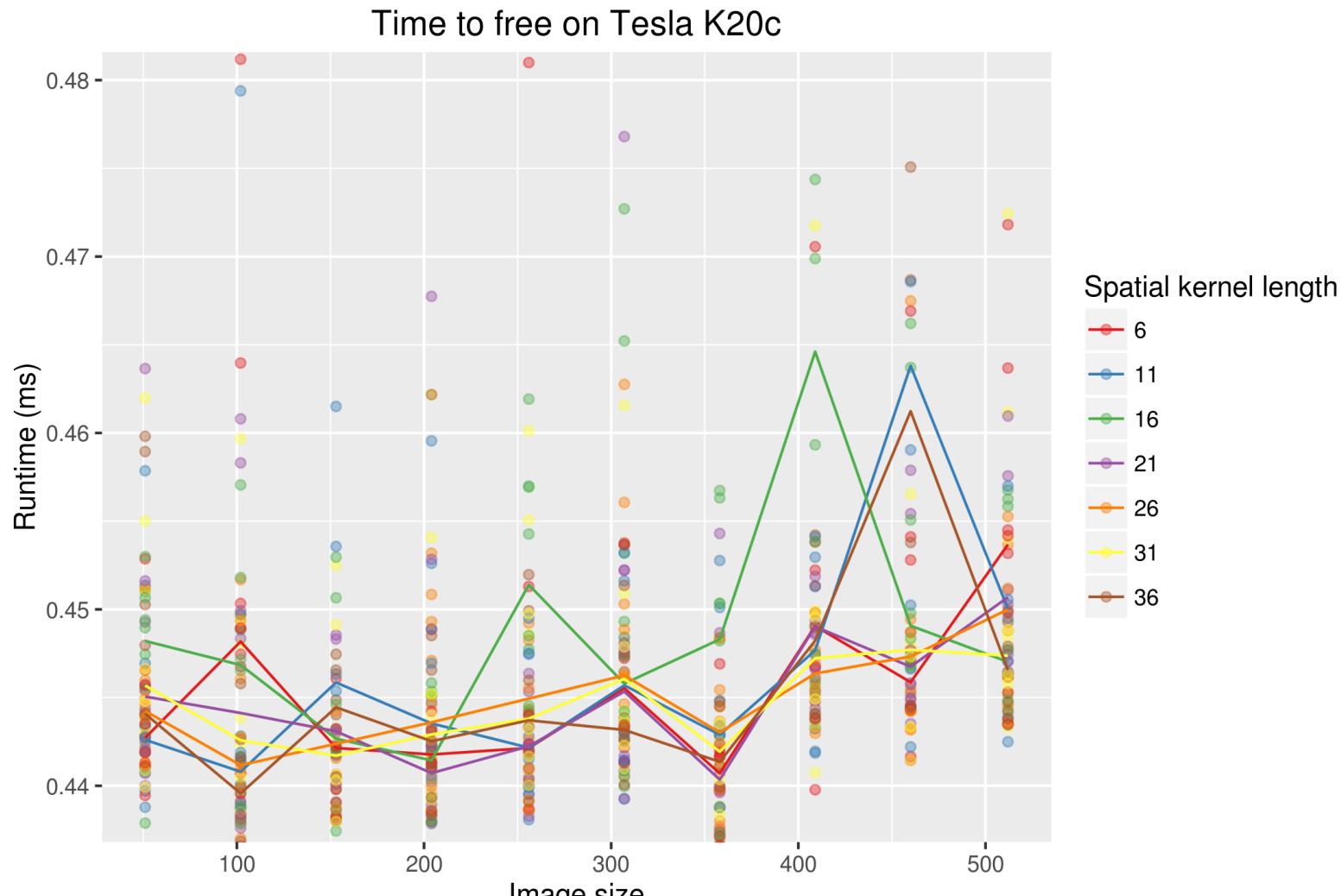
Benchmark

Allocation on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

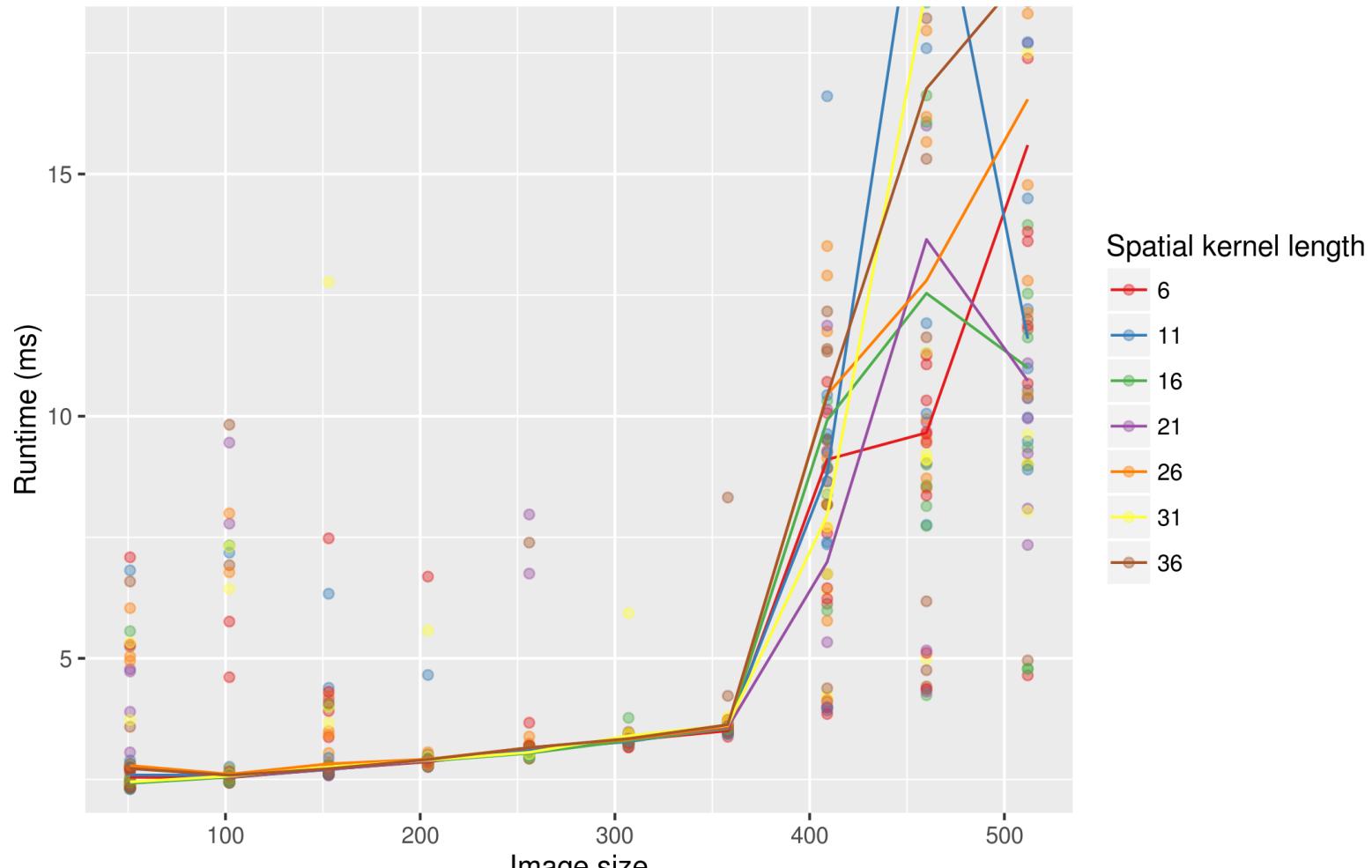
Benchmark



for intensity kernel length, intensity and spatial scaling 11

Benchmark

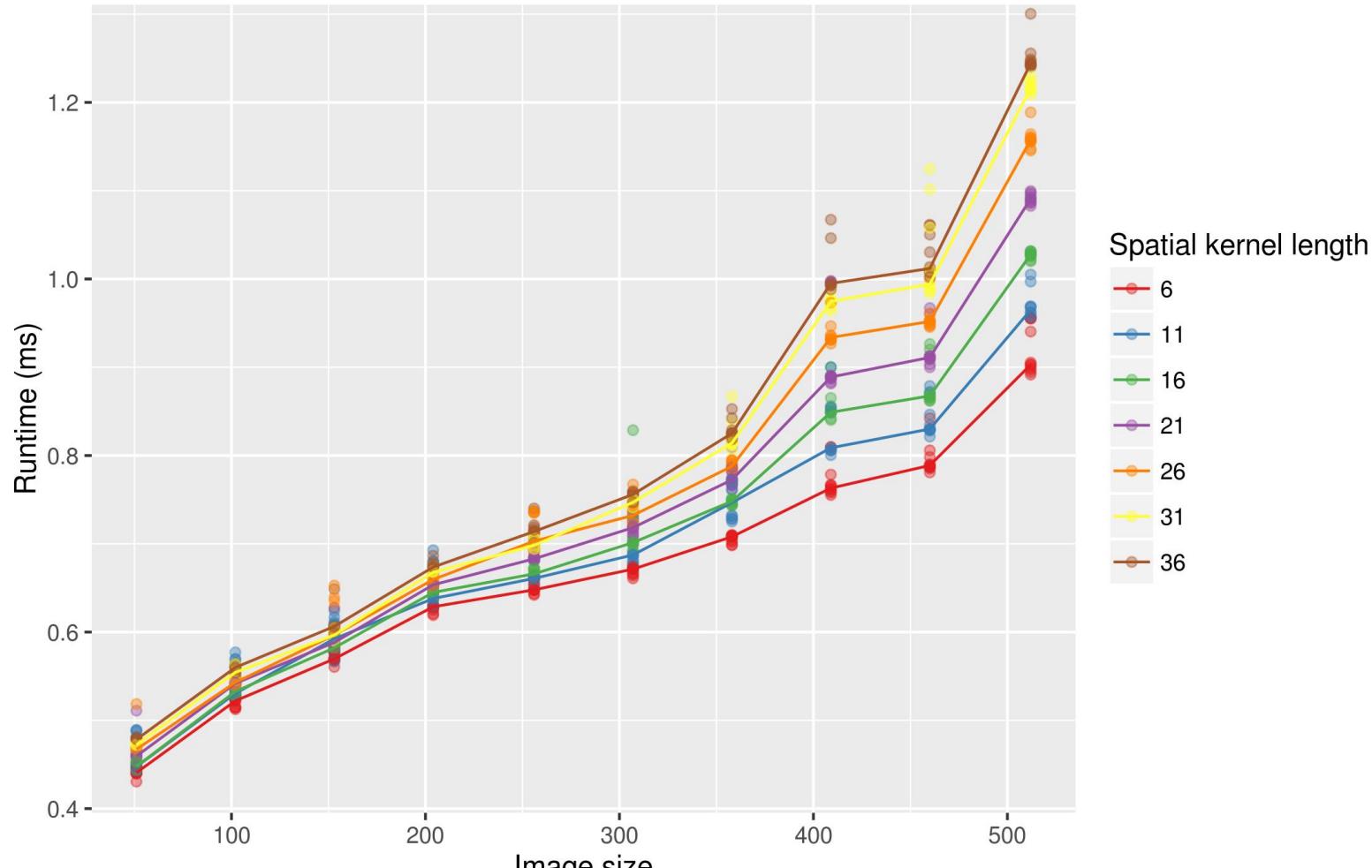
Full run on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

Benchmark

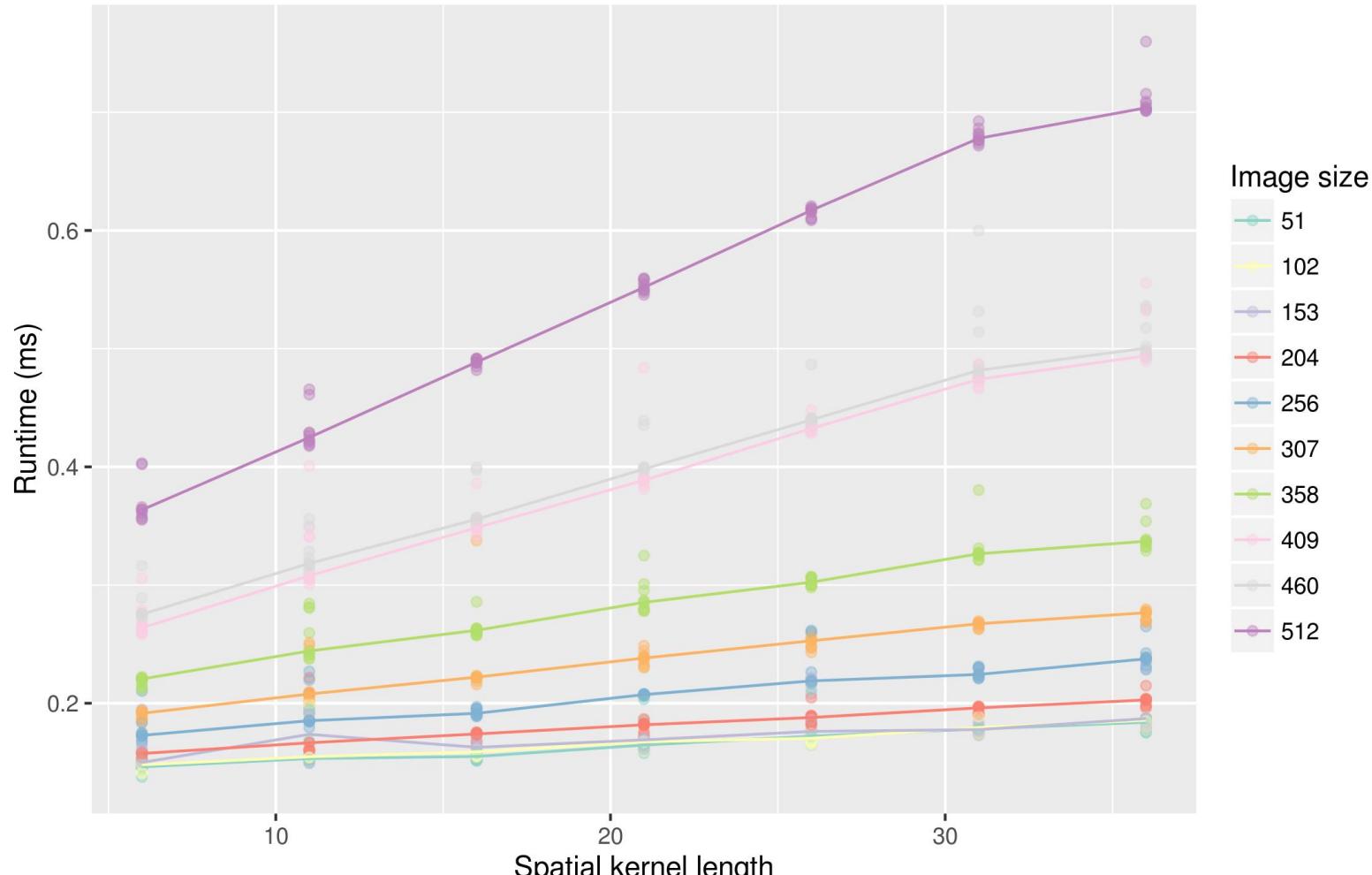
Filling + Convolution + Slicing on Tesla K20c



for intensity kernel length, intensity and spatial scaling 11

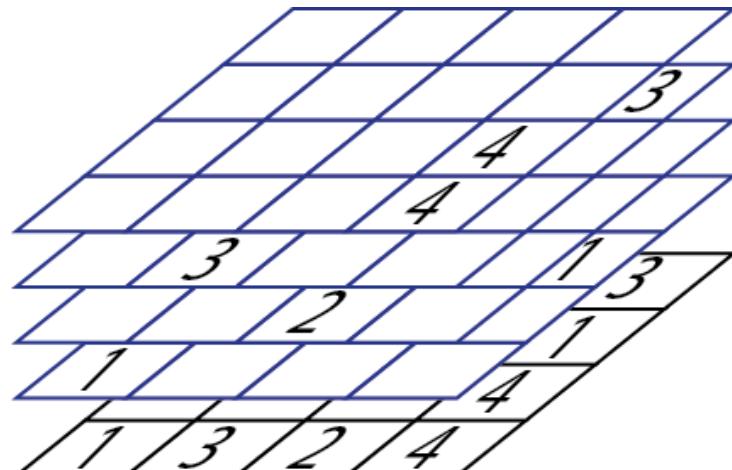
Benchmark

Convolution on Tesla K20c

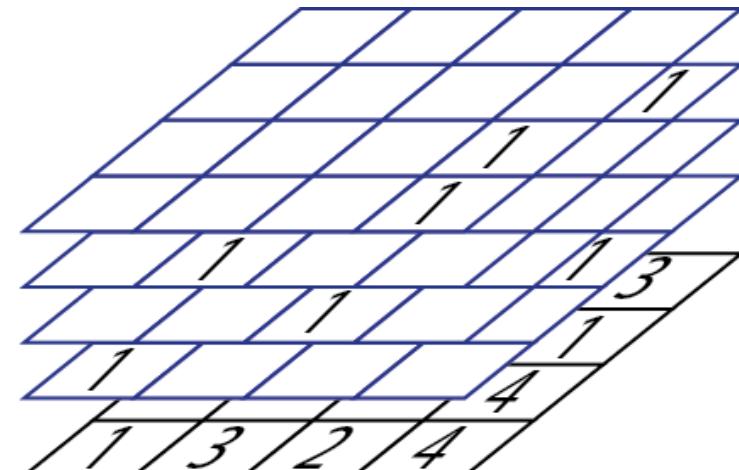


for intensity kernel length, intensity and spatial scaling 11

Benchmark – Best method for filling



WI



W

```

__global__ void cubefilling_loop(const float* image, float *dev_cube_wi, float *dev_cube_w, const dim3 image_size, int scale_xy, int scale_eps,
dim3 dimensions_down)
{
    unsigned int i = blockIdx.x * blockDim.x + threadIdx.x;
    unsigned int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < dimensions_down.x && j < dimensions_down.y) {

        size_t cube_idx_1 = i + dimensions_down.x*j;
        #pragma unroll
        for (int ii = 0; ii < scale_xy; ii++)
        {
            #pragma unroll
            for (int jj = 0; jj < scale_xy; jj++)
            {
                size_t i_idx = scale_xy*i + ii;
                size_t j_idx = scale_xy*j + jj;
                if (i_idx < image_size.x && j_idx < image_size.y)
                {

                    float k = image[i_idx + image_size.x*j_idx];
                    size_t cube_idx_2 = cube_idx_1 +
                    dev_cube_wi[cube_idx_2] += k;
                    dev_cube_w[cube_idx_2] += 1.0f;
                }
            }
        }
    }
}

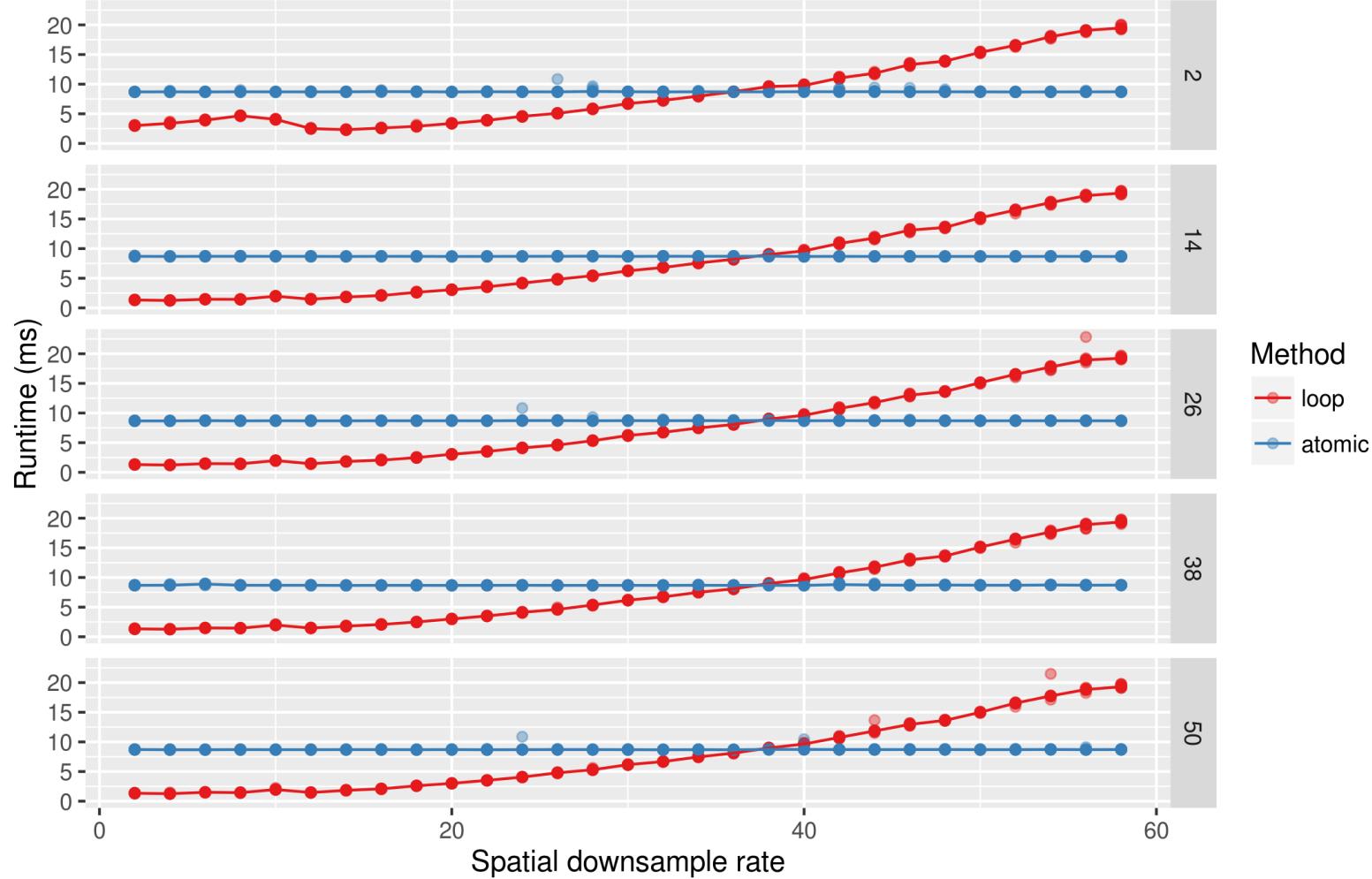
```

```
__global__ void cubefilling_atomic(const float* image, float *dev_cube_wi, float *dev_cube_w, const dim3 image_size, int scale_xy, int
scale_eps, dim3 dimensions_down)
{
    const size_t i = blockIdx.x * blockDim.x + threadIdx.x;
    const size_t j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < image_size.x && j < image_size.y) {
        const float k = image[i + image_size.x*j];
        const size_t cube_idx = (i / scale_xy) + dimensions_down.x*(j / scale_xy) +
dimensions_down.x*dimensions_down.y*((int)k / scale_eps);

        atomicAdd(&dev_cube_wi[cube_idx], k);
        atomicAdd(&dev_cube_w[cube_idx], 1.0f);
    }
}
```

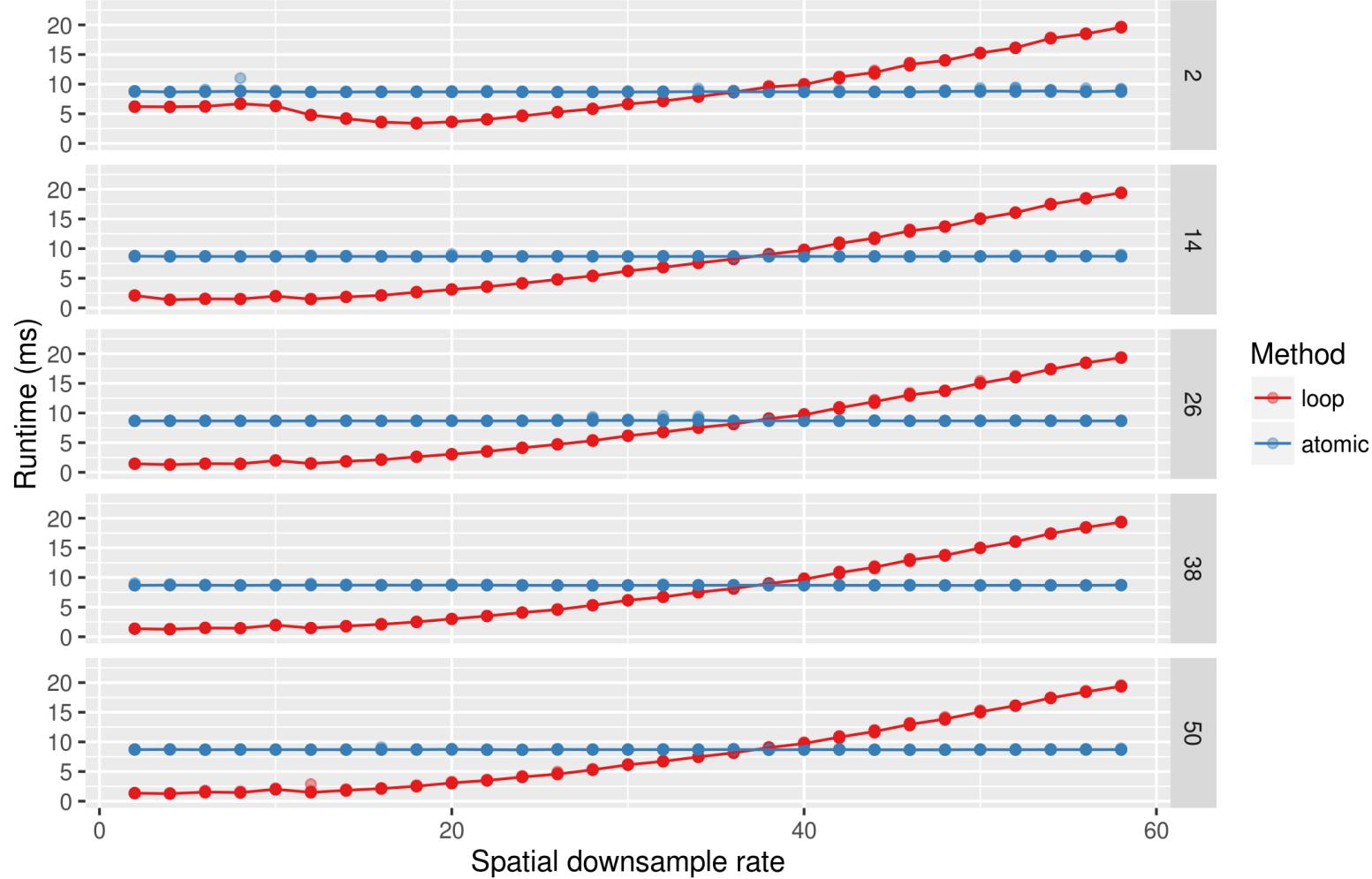
Benchmark – Best method for filling

Filling on GeForce 940M - Lena



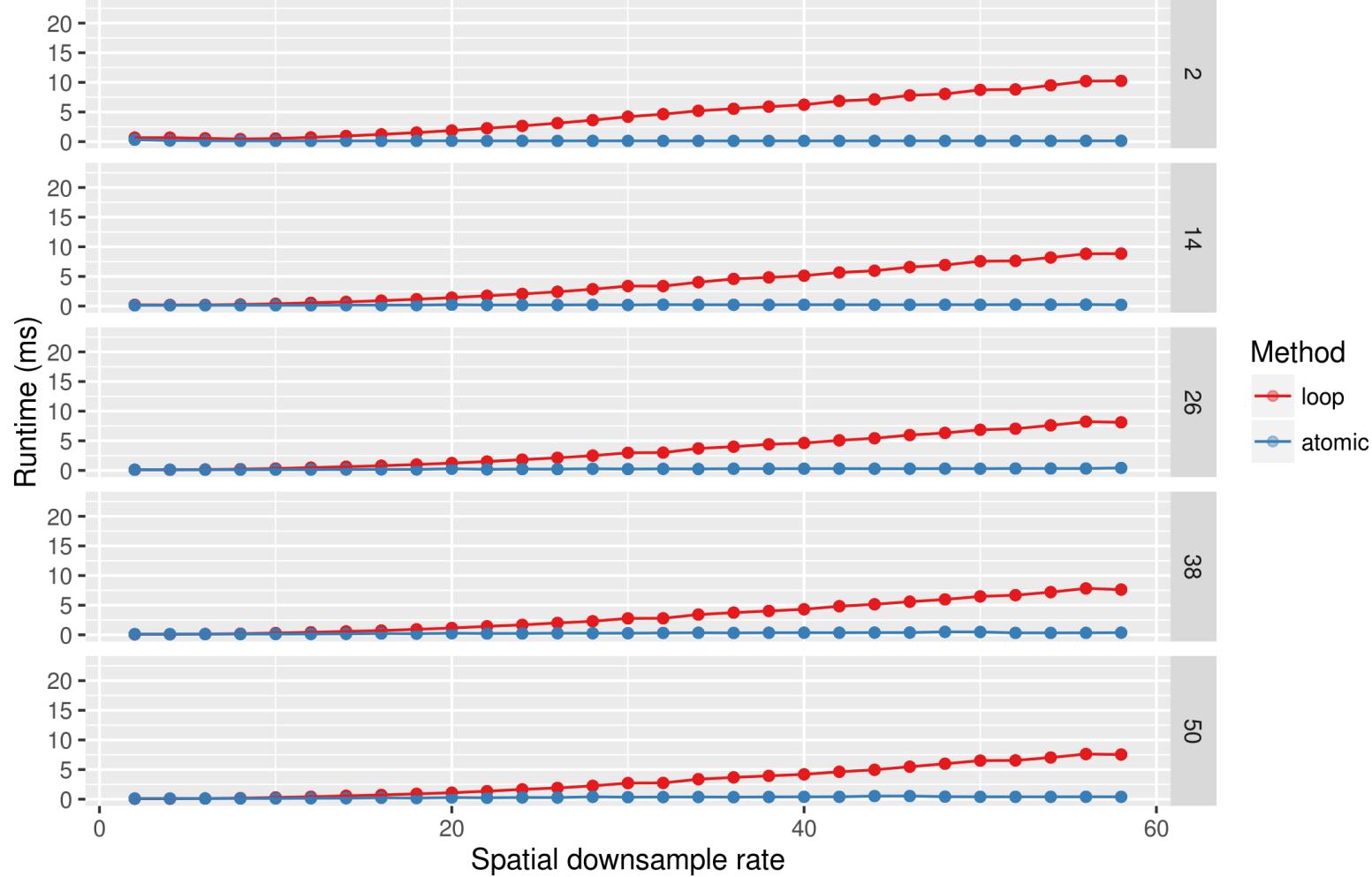
Benchmark – Best method for filling

Filling on GeForce 940M - Random



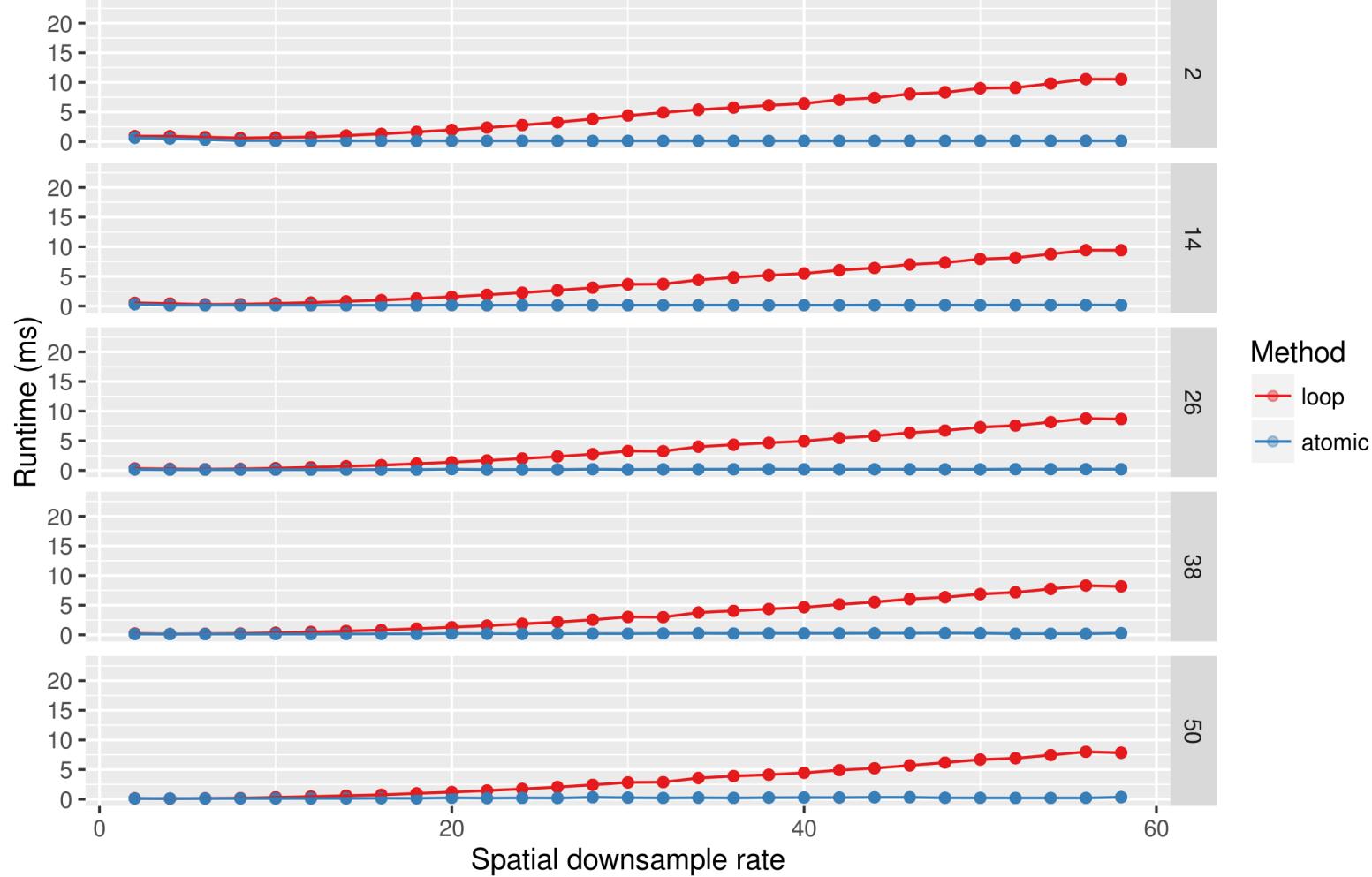
Benchmark – Best method for filling

Filling on Tesla K20c - Lena

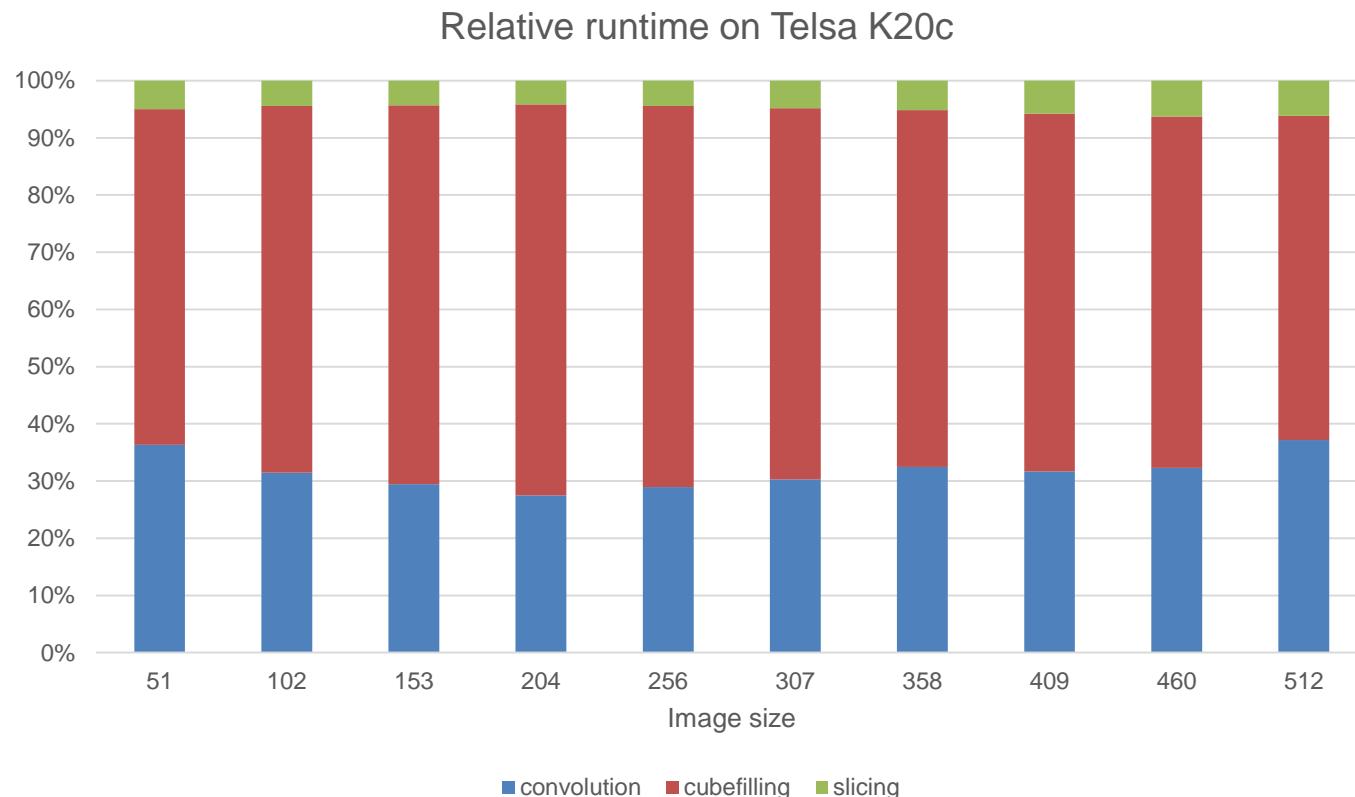


Benchmark – Best method for filling

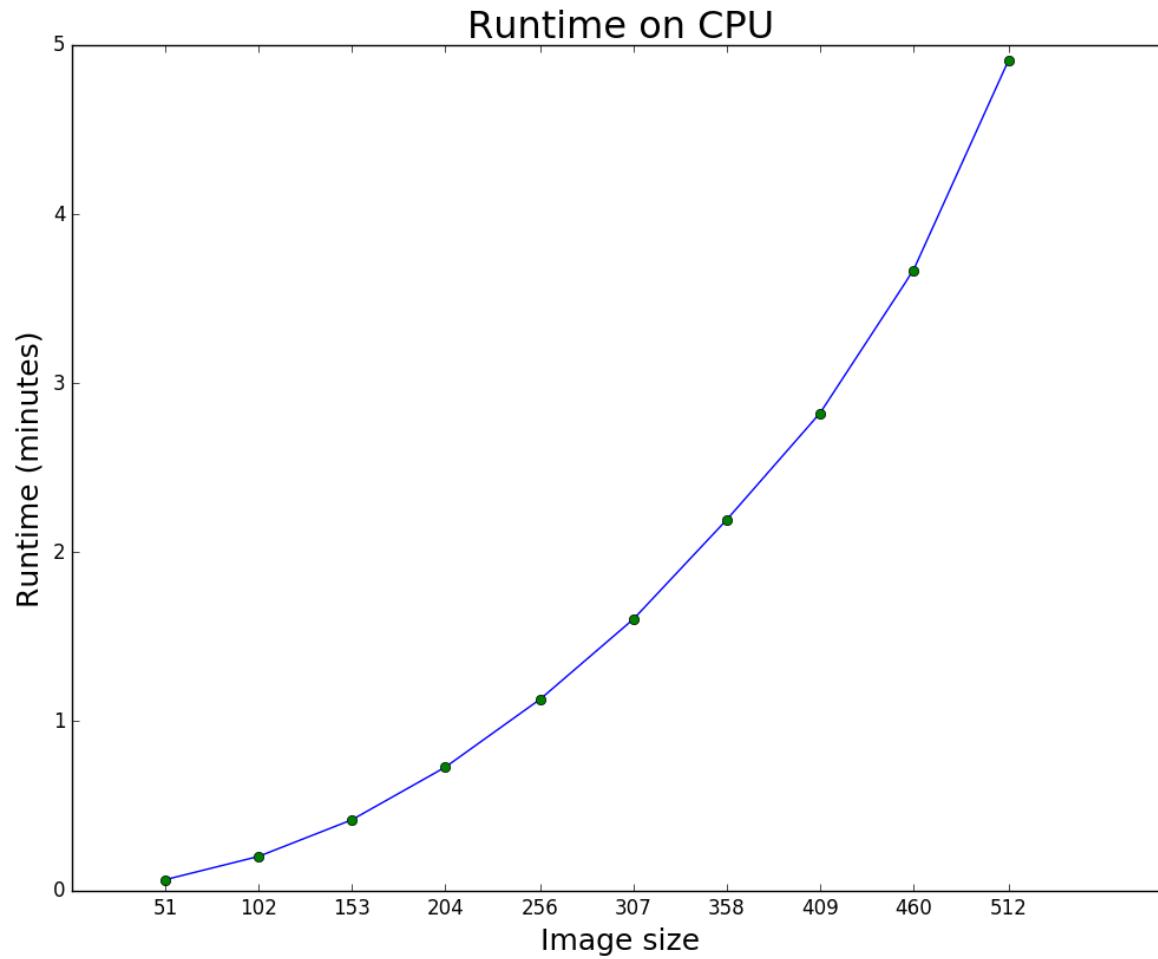
Filling on Tesla K20c - Random



Benchmark



Benchmark - CPU



for intensity kernel length 21 and spatial k.l. 25

Issues on the Implementation

- On windows, a register key must be set to allow the GPU to run a kernel more than 2 seconds
(HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Control\GraphicsDrivers > TdrDelay=10)
- On (almost) every CUDA API call, retrieve cudaStatus and check for errors. It doesn't raise any exceptions.
- On linux, add some lines to `.bashrc`

```
export PKG_CONFIG_PATH=/scratch-local/usr/lib64/pkgconfig/
export PATH=$PATH:/opt/cuda/bin
export LD_LIBRARY_PATH=/scratch-local/usr/lib/:$(LD_LIBRARY_PATH)
```

If you are interested on the code and more plots...

- <https://github.com/bernardelli/Bilateral-Filter-GPU>

Thank You For Your Attention!

References

- <https://developer.apple.com/library/prerelease/content/documentation/Performance/Conceptual/vImage/ConvolutionOperations/ConvolutionOperations.html>
- <http://homepages.inf.ed.ac.uk/rbf/HIPR2/gsmooth.htm>
- C. Tomasi and R. Manduchi **Bilateral Filtering for gray and color images.** In *Proceedings of IEEE International Conference on Computer Vision*, pages 839-846
- http://people.csail.mit.edu/sparis/publi/2006/eccv/Paris_06_Fast_Approximation.pdf
- http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#linear-filtering_linear-filtering-of-1-d-texture-of-4-texels