Real-Time Frequency Selective Reconstruction through Register-Based Argmax Calculation

Andy Regensky, Simon Grosche, Jürgen Seiler, and André Kaup *Multimedia Communications and Signal Processing Friedrich-Alexander University Erlangen-Nürnberg (FAU)* Cauerstr. 7, 91058 Erlangen, Germany {andy.regensky, simon.grosche, juergen.seiler, andre.kaup}@fau.de

Abstract—Frequency Selective Reconstruction (FSR) is a stateof-the-art algorithm for solving diverse image reconstruction tasks, where a subset of pixel values in the image is missing. However, it entails a high computational complexity due to its iterative, blockwise procedure to reconstruct the missing pixel values. Although the complexity of FSR can be considerably decreased by performing its computations in the frequency domain, the reconstruction procedure still takes multiple seconds up to multiple minutes depending on the parameterization. However, FSR has the potential for a massive parallelization greatly improving its reconstruction time. In this paper, we introduce a novel highly parallelized formulation of FSR adapted to the capabilities of modern GPUs and propose a considerably accelerated calculation of the inherent argmax calculation. Altogether, we achieve a 100-fold speed-up, which enables the usage of FSR for real-time applications.

Index Terms-image reconstruction, parallelization, argmax

I. INTRODUCTION

Efficient image reconstruction algorithms are in high demand due to their broad applicability to diverse signal processing tasks, such as error concealment [1], image inpainting [2] or resolution enhancement [3]. Frequency Selective Reconstruction (FSR) [3] has been introduced as an image reconstruction algorithm in the context of image resolution enhancement and employs non-regular sampling to increase the quality of the reconstructed high resolution image. By tailoring the image acquisition process and the reconstruction algorithm to each other, accurate reconstructions can be achieved. In socalled quarter sampling [4], a higher resolution is achieved without increasing the number of pixels on the sensor by non-regularly covering 3/4 of each pixel and reconstructing the high resolution image using FSR. This yields a better quality of the high resolution image than the extrapolation of conventional image sensor data while reducing the cost compared to increasing the resolution of the image sensor.

The main problem of most modern reconstruction algorithms, be it in the context of image signals or in the related field of Compressed Sensing [5], [6] in general, is their high computational complexity and the resulting long execution times. For FSR, different approaches have been followed to reduce the overall reconstruction time on the CPU. In [7], FSR has been adapted to perform real-valued operations exclusively by reconstructing the image signal in the Hartley domain. In [8], constraints have been imposed on FSR in order to reduce the total number of operations that is required to reconstruct an image, yet, potentially impairing the overall reconstruction quality. In [9], extensive pre-computations are performed to speed up the iterative reconstruction procedure while increasing the memory cost of the algorithm.

On the other hand, many algorithms have been adapted to exploit the massive parallelization capabilities of modern GPUs. In [10], the Matching Pursuit algorithm [11] for sparse signal recovery has been sped up by performing the involved matrix-vector operations on the GPU. In [12], this idea was extended to Orthogonal Matching Pursuit (OMP) [13]. In [14], the idea is furthermore extended to the 2D-OMP [15], where a custom implementation of the required argmax operation on the GPU is employed, that exploits the increased speed of shared memory compared to global device memory.

FSR allows for a massive parallelization on the GPU as well. In this paper, we introduce a parallelization of the FSR algorithm on multiple levels and propose a highly effective, high speed argmax calculation. Other than [14], which uses shared memory for thread cooperation, we exploit the rapid register access between neighbouring threads on the GPU. By reducing global and shared memory access to a minimum, a considerable speed-up of the overall reconstruction procedure is achieved.

This paper is organized as follows: Section II shortly recaps the FSR algorithm forming the basis of this work. Section III describes the proposed parallelization of FSR for GPUs and explains the novel register-based argmax calculation in detail. Section IV compares the execution time of the proposed highly parallelized GPU implementation of FSR to its CPU counterpart, as well as a shared memory based GPU implementation, and evaluates the achievable framerates for different resolutions. Finally, Section V concludes this paper.

II. FREQUENCY SELECTIVE RECONSTRUCTION

FSR is a blockwise reconstruction algorithm, that subdivides the image into neighboring target blocks of size $B \times B$ pixels and reconstructs each block independently. During the iterative reconstruction procedure, it takes a neighborhood of L pixels to all sides of the target block into account to take advantage of additional information. The resulting block of size $S \times S = (B + 2L) \times (B + 2L)$ pixels is called the support block. FSR builds a model g[m, n]

^{values. Atthoug} decreased by p domain, the rec up to multiple However, FSR greatly improve introduce a now to the capabilit accelerated calo gether, we achied of FSR for real *Index Terms*Efficient imate mand due to t cessing tasks, ting [2] or rest Reconstruction hancement and quality of the rethe image acqui to each other, a called quarter without increa non-regularly of the high resol quality of the of conventiona compared to in The main p rithms, be it in field of Compi computational

^{© 2020} IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works. DOI: 10.1109/MMSP48831.2020.9287071.

for each support block as a weighted superposition of 2D Discrete Fourier Transform (DFT) basis images φ_{kl} [3], [16],

$$g[m,n] = \sum_{k=0}^{S-1} \sum_{l=0}^{S-1} c_{kl} \varphi_{kl}[m,n]$$
(1)

with pixel positions (m, n) and frequency components (k, l). It is capable of reconstructing images where a large subset of pixel values are missing. The sampled signal $\tilde{f}[m, n]$ for a support block is obtained from the support block signal f[m, n] as

$$\tilde{f}[m,n] = \begin{cases} f[m,n] & \text{for } (m,n) \in \mathcal{K}, \\ 0 & \text{otherwise,} \end{cases}$$
(2)

where the set \mathcal{K} subsumes the indices of all sampled pixels. It follows an iterative approach to generate the model g[m, n], where in each iteration ν , the coefficient c_{kl} of exactly one basis image $\varphi_{kl}[m, n]$ is changed. It is selected such that the weighted residual energy

$$E_w^{(\nu)} = \sum_{m,n} |r^{(\nu)}[m,n]|^2 w[m,n]$$
(3)

is minimized. Thereby, the residual $r^{(\nu)}[m,n]$ describes the difference between the sampled image $\tilde{f}[m,n]$ and the model g[m,n] in iteration ν as

$$r^{(\nu)}[m,n] = \tilde{f}[m,n] - g^{(\nu)}[m,n], \tag{4}$$

and w[m,n] describes the spatial weighting

$$w[m,n] = \begin{cases} \hat{\rho}^{\sqrt{\left(m - \frac{S-1}{2}\right)^2 + \left(n - \frac{S-1}{2}\right)^2}} & \text{for } (m,n) \in \mathcal{K}, \\ 0 & \text{otherwise,} \end{cases}$$
(5)

where the scalar decay factor $\hat{\rho}$ controls the speed of decay of the exponentially decreasing weight [3], [16]. Unlike [3], we refrain from taking into account already reconstructed pixel values in exchange for a more efficient parallelization of the overall algorithm. The spatial weighting is employed to ignore unknown samples $(m, n) \notin \mathcal{K}$ during the calculation of the weighted residual energy and to assign an exponentially decreasing weight to pixels further away from the block center.

As each basis image $\varphi_{kl}[m,n]$ corresponds to exactly one frequency component, the algorithm can be efficiently performed in the frequency domain. By describing the weighted residual $r_w^{(\nu)}[m,n]$ as

$$r_w^{(\nu)}[m,n] = r^{(\nu)}[m,n] \cdot w[m,n], \tag{6}$$

and formulating the model $g^{(\nu)}[m,n]$, the weighted residual $r_w^{(\nu)}[m,n]$ and the spatial weighting w[m,n] in the frequency domain, we get $G^{(\nu)}[k,l]$, $R_w^{(\nu)}[k,l]$ and W[k,l], respectively. The initial model $G^{(0)}[k,l]$ is set to 0. With \mathcal{F}_2 denoting the 2D DFT, the weighted residual $R_w[k,l]$ is therefore initialized to

$$R_w^{(0)}[k,l] = \mathcal{F}_2\{f[m,n] \cdot w[m,n]\}.$$
(7)

In each iteration, the basis image $\varphi_{kl}[m,n]$ is selected that reduces the weighted residual energy $E_w^{(\nu)}$ the most,

so that the frequency index $(u,v)^{(\nu)}$ of the selected basis image $\varphi_{uv}[m,n]$ in iteration ν is obtained as

$$(u,v)^{(\nu)} = \operatorname{argmax}_{(k,l)} \left(w_f[k,l] \cdot |R_w^{(\nu)}[k,l]|^2 \right).$$
(8)

Thereby, the frequency weighting $w_f[k, l]$ is incorporated to favor lower frequency basis images over higher ones [3]. It is defined as

$$w_f[k,l] = \left(1 - \sqrt{2}\sqrt{\frac{\tilde{k}^2}{S^2} + \frac{\tilde{l}^2}{S^2}}\right)^2 \tag{9}$$

with $\tilde{k} = \frac{S}{2} - |k - \frac{S}{2}|$ and $\tilde{l} = \frac{S}{2} - |l - \frac{S}{2}|$. The projection coefficient $p_{uv}^{(\nu)}$ is chosen such that the selected basis image $\varphi_{uv}[m,n]$ maximally reduces the weighted residual energy $E_w^{(\nu)}$ resulting in

$$p_{uv}^{(\nu)} = \frac{R_w^{(\nu)}[u,v]}{W[0,0]}.$$
(10)

Once the frequency component $(u, v)^{(\nu)}$ to be updated has been determined, and its projection coefficient $p_{uv}^{(\nu)}$ has been calculated, the model G[k, l] and the weighted residual $R_w[k, l]$ are updated accordingly

$$G^{(\nu+1)}[u,v] = G^{(\nu)}[u,v] + \gamma p_{uv}^{(\nu)} S^2,$$
(11)

$$R_w^{(\nu+1)}[k,l] = R_w^{(\nu)}[k,l] - \gamma p_{uv}^{(\nu)} W[k-u,l-v] \ \forall \ (k,l).$$
(12)

Thereby, the model G[k, l] needs to be updated for the selected frequency component $(u, v)^{(\nu)}$ only, whereas the weighted residual $R_w[k, l]$ needs to be updated for all of the S^2 frequency components (k, l). The scalar orthogonality deficiency compensation factor $0 < \gamma \leq 1$ is introduced to reduce the interference between the different basis images [3]. This procedure of selecting a basis function according to (8), computing the projection coefficient according to (10), and updating the model and the residual according to (11) and (12), respectively, is then repeated for a fixed number of iterations.

Eventually, the final model G[k, l] is transformed back into the spatial domain to obtain g[m, n]. The final support block reconstruction $\hat{f}[m, n]$ is obtained by taking over the sampled support block signal $\tilde{f}[m, n]$ for known samples, and taking over the final model g[m, n] for unknown samples. To end the reconstruction of the current block, the reconstructed target block signal is extracted from the reconstructed support block signal $\hat{f}[m, n]$ and placed at the corresponding target block in the global reconstruction. Please refer to [3] for a more detailed explanation of FSR and its parameters.

III. FREQUENCY SELECTIVE RECONSTRUCTION ON THE GPU

FSR in the frequency domain shows great potential for large speed improvements through a massive parallelization of the internal processing procedure. Especially on GPUs, which are specifically designed for running many lightweight tasks simultaneously, a careful design of the parallelized FSR algorithm promises a huge gain in processing speed. In this section, we describe the applied GPU thread model and explain the highly parallelized FSR algorithm that is specifically adapted to it. Furthermore, we take a deeper look at the novel register-based argmax calculation on the GPU, which serves as a major accelerator of the overall reconstruction procedure.

A. Parallelization of FSR

Parallelizing FSR requires a certain understanding of the thread model that the applied GPU uses. The thread model describes how data-parallel and task-parallel workloads are subdivided on the device, and hence, understanding its basics is essential to be able to design an efficient GPU algorithm. As Nvidia[®]'s CUDA[®] [17], [18] programming language enjoys high popularity in the scientific community and due to the broad availability of compatible GPU devices, we employ the CUDA thread model in our work. Note though, that competitor GPU manufacturers commonly use similar models with slightly varying names.

Fig. 1 shows the CUDA thread model [17] where a grid consists of multiple blocks and each block consists of multiple threads. Thereby, all blocks in a grid can run independent from each other, whereas all threads in a block are executing the same program on multiple data elements. All threads in a block run in a shared memory space, i.e., on the same Streaming Multiprocessor (SM) [17], and can interact with each other. Threads in different blocks can not natively interact with each other and data transfer among them requires the involvement of comparably slow global memory. With this basic understanding of the thread model, an efficient way to parallelize FSR can be derived.

The parallelization of FSR happens in two stages: a taskparallel stage and a data-parallel stage. The task-parallel stage is thereby established similarly to how one would parallelize the reconstruction process on a multi-core CPU. As FSR is a blockwise procedure and all blocks can be reconstructed independent from each other, parallelizing the reconstruction of all blocks is an obvious solution. This means that the reconstruction of each support block is assigned to a single thread block on the GPU, each. For an image of size $Y \times X$, a number of

$$N = \left| \frac{Y}{B} \right| \cdot \left| \frac{X}{B} \right| \tag{13}$$

support blocks need to be reconstructed, which results in N thread blocks being employed for the reconstruction of the image. Note, that one can create more thread blocks than there are SMs on the GPU, as they are not required to run simultaneously and one SM can hold multiple thread blocks. The data-parallel stage follows a more sophisticated layout for the reconstruction of one support block using FSR. In general, the parallelization can be interpreted as each thread being responsible for one pixel of the support block. This means, that each thread is dedicated to a pixel position (m, n) in the spatial domain, or a frequency component (k, l) in the frequency domain, and a support block is reconstructed

 \mathbf{Grid}

Block	Block	Block	Block
Threads	Threads	Threads	Threads
•••••	•••••	•••••	•••••
••••	•••••	•••••	•••••
			• • • • • • • • • • • •
Block	Block	Block	Block
Threads	Threads	Threads	Threads
••••	••••	••••	•••••

Fig. 1. Schematic of the employed thread model.

by S^2 threads in parallel. Starting from the regarded sampled support block signal f[m, n], the residual (7) is initialized by first multiplying each pixel of the sampled signal $\tilde{f}[m, n]$ with the spatial weighting w[m, n] simultaneously, where each thread performs a single multiplication, and then executing a 2D FFT using all threads to obtain the residual in the frequency domain. The FFT algorithm is also executed on the GPU but not investigated further since the optimized and tested cuFFT library [19] is used for all FFT and IFFT invocations on the GPU. The model G[k, l] is initialized with all threads setting the value of their respective frequency component to 0. From here, the iterative procedure begins. In each iteration, a basis image is selected, its optimal projection coefficient is calculated, and the model and the residual are updated accordingly. For the selection of the basis image (8), each thread calculates the objective value for the frequency component (k, l) it is assigned to. Then, all threads in the regarded thread block are cooperating to find the index of the maximum frequency component in $\lceil 2 \log_2(S) \rceil$ steps. The thread cooperative argmax calculation is explained in detail in Section III-B. Once the basis image to be added has been found, its projection coefficient (10) has to be computed. For this computation, only one thread is involved as the projection coefficient is only required for the selected basis image. The same is true for the update of the model (11), so that this computation is performed by a single thread, as well. The update of the residual (12) can be done by all threads simultaneously, where a thread dedicated to frequency component (k, l)accesses the corresponding value W[k-u, l-v] = W[i, j]from the Fourier transformed spatial weighting with

$$i = \begin{cases} k - u & k - u \ge 0, \\ S + k - u & \text{otherwise,} \end{cases}$$
(14)

$$j = \begin{cases} l - v & l - v \ge 0, \\ S + l - v & \text{otherwise.} \end{cases}$$
(15)

After the desired number of iterations has been executed, the final model G[k, l] is transformed back into the spatial domain

using the IFFT provided by the cuFFT library. Eventually, each thread writes its assigned pixel value from the reconstructed support block signal g[m, n] to the corresponding pixel position in the final reconstructed image.

B. Register-based argmax calculation

On NVIDIA GPUs, threads are scheduled to run in groups called warps [17]. Thereby, a warp consists of 32 threads on all current Nvidia GPUs. Within each warp, all threads perform the same instructions on different data elements simultaneously. The CUDA programming model provides functionalities for threads within the same warp to read out each others registers. Direct register access is extremely fast and allows for a high-performance implementation of the argmax calculation that occurs in the iterations of FSR.

As a first step, all threads of the current block compute the objective value for their corresponding frequency component (k, l)

$$obj[k, l] = w_f[k, l] \cdot |R_w[k, l]|^2,$$
 (16)

which is stored in a register objective in each thread. Furthermore, each thread stores its frequency component (k, l) in registers index_k and index_l.

The argmax calculation is then performed in two phases. In the first phase, the argmax is searched within each warp of the current block independently. Thereby, the threads within the current block are assigned to warps in a consecutive fashion. In the second phase, the results from the warp-internal first phase are further processed by a single warp to find the argmax result of the overall thread block. Both, the first phase and the second phase perform the same underlying iterative, threadcooperative procedure to rapidly find the argmax result by optimally utilizing the available GPU resources. This procedure consists of three steps that are executed simultaneously by all threads in a warp:

- A thread at index i reads the register value objective of the offset thread at index i + offset and stores it in a register objective_offset. Preceeding the first iteration, the offset is initialized to half the number of threads per warp.
- 2) If the objective value of the offset thread is larger than the objective value of the current thread (objective_offset > objective), the current thread replaces its register values objective, index_k and index_l with the corresponding register values from the offset thread.
- 3) While the offset is larger than 1 (offset > 1), the offset is halved and the procedure continues with Step 1. Otherwise, the warp-internal argmax calculation finishes and the argmax result is guaranteed to reside within the first thread of the warp.

The described procedure is depicted in Fig. 2, where a simplified argmax calculation with 8 threads per warp is performed. Note, that if a thread requests an out-of-bound thread register, its own register value is returned instead.



Fig. 2. Schematic of the argmax calculation within a single warp. Simplified representation with 8 threads per warp. Cell values correspond to the frequency component. Colors represent the objective value, where blue represents low values and red represents high values. In each iteration, each thread compares its own objective value (vertical arrow) with the offset thread's objective value (angled arrow) and takes over the register values from the thread with the higher objective value (solid black arrow). After the last iteration, the result of the argmax calculation lies in the first thread of the warp (black box).

Thereby, a thread register is considered out-of-bounds if it is not a member of the current warp.

Once the first phase completes, each warp's first thread contains the warp-internal argmax result (maximum objective value and corresponding frequency component) and stores it in shared memory for later access. In the second phase, the argmax of the overall block is calculated by evaluating the argmax among all warp-internal results from the first phase. As CUDA allows a maximum number of 1024 threads per block, there are at most 32 warps per block. Hence, there exist at most 32 warp-internal argmax results in shared memory after finishing the first phase. Consequently, in the second phase, a single warp can perform the argmax calculation using the described highly efficient thread-cooperative procedure. This is achieved by each thread reading the results of one warpinternal argmax calculation from shared memory and storing the corresponding values in its local registers objective, index_k and index_1. The iterative, thread-cooperative procedure is then performed as described above. Finally, independent of the chosen blocksize ($S^2 \leq 1024$), the final argmax result will always reside in the first thread of the warp that executed the second phase of the described two-phase procedure. From there, the frequency component (k, l) corresponding to the maximum objective value can be extracted and the FSR algorithm can continue as described in Section III-A. Due to the minimal access of shared memory (global memory is not involved), the register-based argmax calculation is extremely fast and yields a considerable speed-up of the overall



Fig. 3. Comparison of the execution time of the FSR on the CPU and the GPU with a shared memory and our register-based argmax approach.

reconstruction procedure.

IV. SIMULATIONS

To evaluate the performance of the proposed parallelization of FSR for GPUs including the novel register-based argmax calculation, the reconstruction procedure on the GPU is compared to the reconstruction procedure on the CPU in terms of quality and speed. To realize a fair comparison, a highly optimized CPU implementation in C++ is employed. All CPU computations are performed on one core of an Intel[®] Xeon® E5 2690 @ 2.6 GHz. Simulations that use the GPU implementation run on an Nvidia GeForce® RTX 2080 Ti with 11 GB of global device memory featurng 4608 cores. Thereby, code that runs on the GPU is written in Nvidias CUDA programming language and all parts that require CPU interaction (memory transfer to/from the GPU, starting the reconstruction procedure) are written in $C++^1$. The cuFFT library [19] is used for the necessary FFT calculations on the GPU. Execution time and Peak Signal to Noise Ratio (PSNR) are averaged over 100 grayscale images from the TECNICK image dataset [20] for a given parameterization. If not stated otherwise, the images are processed at their original resolution of 1200×1200 pixels. The target blocksize is set to B = 4 pixels, the spatial decay factor to $\hat{\rho} = 0.7$ and the orthogonality deficiency compensation factor to $\gamma = 0.5$. The sampled image f is obtained through a quarter sampling sensor [3], which is emulated by sampling a random pixel in each neighboring 2×2 block of the original image f.

As the proposed parallelization of FSR follows the same mathematical procedure as the conventional FSR on the CPU, both implementations produce identical results for a given parameterization. However, due to the massively parallel execution of FSR on the GPU, a considerable gain of speed with respect to the single-threaded CPU implementation can be observed as visible in Fig. 3. Thereby, results for a shared memory GPU implementation are included, where the argmax calculation is performed using a shared memory



Fig. 4. Extensive parameter exploration regarding quality and execution time for the CPU and the GPU implementations of the FSR. Each data point represents a single parameterization averaged over all test images. In total, 640 parameterizations were tested. For details, see text.

approach as described in [14] instead of the novel registerbased approach. Other than that, the shared memory GPU implementation is identical to the introduced parallelization of FSR. It is clearly visible, that the register-based argmax calculation yields notable speed improvements over the shared memory approach. With the register-based argmax calculation, speed improvements of more than $100 \times$ are possible with respect to the single-threaded CPU implementation. While the CPU performance could be greatly improved by higher clock frequencies and performing the reconstruction of different support blocks in parallel on multiple CPU cores, note, that even for a twice as fast clock frequency, more than 50 CPU cores are required to compete with the performance of the highly parallelized FSR on a single GPU.

Fig. 4 verifies that the significant gain of speed can be generalized to a wide range of parameters. Each data point represents the PSNR and execution time for a given parameterization averaged over all test images. The number of iterations I has been varied in the range [100; 400] using a step size of 100, the support blocksize S has been varied in the range [8; 32] using a step size of 8, the spatial decay factor $\hat{\rho}$ has been varied in the range [0.68; 0.82] using a step size of 0.02, and the orthogonality deficiency compensation factor γ has been varied in the range [0.2; 0.6] using a step size of 0.1. Therefore, 640 parameterizations are tested for the GPU and the CPU implementation, each. A notable gain of speed can be observed for all parameterizations by employing the highly parallelized GPU implementation compared to its CPU counterpart. For all tested parameterizations, one can expect the GPU implementation to be roughly $100 \times$ faster than the single-threaded CPU implementation, validating our findings from above.

With many parameterizations of the parallelized FSR fin-

¹The source code for all evaluated implementations of FSR is publicly available at https://gitlab.lms.tf.fau.de/lms/gpu-fsr



Fig. 5. Average framerate of the novel parallelized FSR for various resolutions on different GPUs compared to a multithreaded CPU implementation on an Intel CoreTM i9-9900 @ 3.1 GHz.

ishing considerably below one second while still producing decent results in terms of quality, a look at the achievable framerates measured in frames per second (fps) is of interest. Fig. 5 compares the average framerate of the parallelized FSR for various common video resolutions on different GPUs to a multithreaded implementation of FSR on the CPU. Thereby, an Intel Core[™] i9-9900 @ 3.1 GHz is employed for the CPU computations, where the 8 cores of the CPU perform the reconstruction of different support blocks simultaneously. The investigated video resolutions are 640×480 (VGA), 800×600 (SVGA), 1920×1080 (FHD) and 3840×2160 (UHD). The support blocksize is set to S = 16. Considering VGA, more than 30 fps can be achieved which makes the proposed highly parallelized FSR capable of real-time applications. Furthermore, it is visible that the algorithm scales reliably with the computing power of different GPUs. For the larger FHD resolution, more than 1 fps is achieved on all applied GPUs, while the powerful RTX 2080 Ti reaches almost 10 fps. Even for the state-of-theart UHD resolution, most of the GPUs among the test field are still able to steadily reconstruct the original image data in less than one second. The multithreaded CPU implementation reaches more than 1 fps only for the low resolutions VGA and SVGA, and, with a loss of about 70% in terms of framerate, is well behind even the slow GTX 1060 for all investigated resolutions.

V. CONCLUSION

In this paper, a highly parallellized implementation of FSR carefully designed for the execution on GPUs has been introduced. In combination with the novel highly effective, high-speed argmax calculation based on direct register access between neighbouring threads, the execution time of the reconstruction procedure can be considerably reduced. The registerbased argmax calculation proved to be a major accelerator of the overall reconstruction procedure leading to notable speed improvements over the shared memory based approach. Depending on the applied hardware, speed-ups of more than $100 \times$ are possible compared to previous approaches without losing image quality. These speed-ups are reliably obtained for a wide range of parameterizations. All in all, the proposed highly parallelized approach is able to reconstruct multiple frames per second on a wide range of hardware and is capable of real-time applications with more than 30 frames per second.

REFERENCES

- G. Zhai, X. Yang, W. Lin, and W. Zhang, "Bayesian Error Concealment With DCT Pyramid for Images," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 20, no. 9, pp. 1224–1232, Sep 2010.
- [2] C. Guillemot and O. Le Meur, "Image Inpainting : Overview and Recent Advances," *IEEE Signal Process. Mag.*, vol. 31, no. 1, pp. 127–144, Jan 2014.
- [3] J. Seiler, M. Jonscher, M. Schöberl, and A. Kaup, "Resampling Images to a Regular Grid From a Non-Regular Subset of Pixel Positions Using Frequency Selective Reconstruction," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 4540–4555, Nov 2015.
- [4] M. Schöberl, J. Seiler, S. Foessel, and A. Kaup, "Increasing imaging resolution by covering your sensor," in *Proc. 18th IEEE Int. Conf. Image Process.*, Sep 2011, pp. 1897–1900.
- [5] D. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr 2006.
- [6] E. Candes, J. Romberg, and T. Tao, "Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb 2006.
- [7] N. Genser, S. Grosche, J. Seiler, and A. Kaup, "Sparse Hartley Modeling for Fast Image Extrapolation," in *IEEE 20th Int. Work. Multimed. Signal Process.*, Aug 2018, pp. 1–6.
- [8] N. Genser, J. Seiler, and A. Kaup, "Spectral Constrained Frequency Selective Extrapolation for Rapid Image Error Concealment," in *Proc.* 25th Int. Conf. Syst. Signals Image Process., Jun 2018, pp. 1–5.
- [9] J. Seiler and A. Kaup, "A Fast Algorithm for Selective Signal Extrapolation with Arbitrary Basis Functions," *EURASIP J. Adv. Signal Process.*, vol. 2011, no. 1, p. 495394, Dec 2011.
- [10] M. Andrecut, "Fast GPU Implementation of Sparse Signal Recovery from Random Projections," Sep 2008.
- [11] S. Mallat and Zhifeng Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Trans. Signal Process.*, vol. 41, no. 12, pp. 3397– 3415, Dec 1993.
- [12] Y. Fang, L. Chen, J. Wu, and B. Huang, "GPU Implementation of Orthogonal Matching Pursuit for Compressive Sensing," in *IEEE 17th Int. Conf. Parallel Distrib. Syst.*, Dec 2011, pp. 1044–1047.
- [13] Y. Pati, R. Rezaiifar, and P. Krishnaprasad, "Orthogonal matching pursuit: recursive function approximation with applications to wavelet decomposition," in *Proc. 27th Asilomar Conf. Signals, Syst. Comput.*, Nov 1993, pp. 40–44.
- [14] Y. Dai, D. He, Y. Fang, and L. Yang, "Accelerating 2D orthogonal matching pursuit algorithm on GPU," J. Supercomput., vol. 69, no. 3, pp. 1363–1381, Sep 2014.
- [15] Y. Fang, J. Wu, and B. Huang, "2D sparse signal recovery via 2D orthogonal matching pursuit," *Sci. China Inf. Sci.*, vol. 55, no. 4, pp. 889–897, Apr 2012.
- [16] A. Kaup, K. Meisinger, and T. Aach, "Frequency selective signal extrapolation with applications to error concealment in image communication," *AEU - Int. J. Electron. Commun.*, vol. 59, no. 3, pp. 147–156, Jun 2005.
- [17] NVIDIA Corporation, "CUDA C Programming Guide," Apr 2018. [Online]. Available: https://docs.nvidia.com/pdf/CUDA_C_ Programming_Guide.pdf
- [18] J. Nickolls, I. Buck, M. Garland, and K. Skadron, "Scalable parallel programming with CUDA," *Queue*, vol. 6, no. 2, p. 40, Mar 2008.
- [19] NVIDIA Corporation, "cuFFT Library User's Guide," Aug 2019.
 [Online]. Available: https://docs.nvidia.com/cuda/pdf/CUFFT_Library.
 pdf
- [20] N. Asuni and A. Giachetti, "TESTIMAGES: A Large Data Archive For Display and Algorithm Testing," J. Graph. Tools, vol. 17, no. 4, pp. 113–125, Oct 2013.