



Published in final edited form as:

Magn Reson Med. 2020 December ; 84(6): 3054–3070. doi:10.1002/mrm.28338.

Advancing machine learning for MR image reconstruction with an open competition: Overview of the 2019 fastMRI challenge

Florian Knoll^{†,1}, Tullie Murrell^{†,2}, Anuroop Sriram^{†,2}, Nafissa Yakubova², Jure Zbontar², Michael Rabbat², Aaron Defazio², Matthew J. Muckley¹, Daniel K. Sodickson¹, C. Lawrence Zitnick², Michael P. Recht¹

¹Center for Advanced Imaging Innovation and Research (CAI²R), Department of Radiology, New York University Grossman School of Medicine, New York, NY, 10016 United States

²Facebook AI Research, Menlo Park, CA, 94025 United States

Abstract

Purpose: To advance research in the field of machine learning for MR image reconstruction with an open challenge.

Methods: We provided participants with a dataset of raw k-space data from 1,594 consecutive clinical exams of the knee. The goal of the challenge was to reconstruct images from these data. In order to strike a balance between realistic data and a shallow learning curve for those not already familiar with MR image reconstruction, we ran multiple tracks for multi-coil and single-coil data. We performed a two-stage evaluation based on quantitative image metrics followed by evaluation by a panel of radiologists. The challenge ran from June to December of 2019.

Results: We received a total of 33 challenge submissions. All participants chose to submit results from supervised machine learning approaches.

Conclusion: The challenge led to new developments in machine learning for image reconstruction, provided insight into the current state of the art in the field, and highlighted remaining hurdles for clinical adoption.

Keywords

Challenge; Image reconstruction; Parallel imaging; Machine Learning; Compressed Sensing; Fast Imaging; Optimization; Public Dataset

1 | INTRODUCTION

One of the fastest growing fields of research in medical imaging during the last several years is the use of machine learning methods for image reconstruction. Machine learning has been proposed for CT dose reduction [1, 2, 3, 4, 5, 6], attenuation correction for PET-MRI [7] and accelerated MR imaging [8, 9, 10, 11, 12, 13, 14, 15, 16]. Despite various methodological

Correspondence: Florian Knoll, Department of Radiology, NYU Grossman School of Medicine, New York, NY, 10016, United States, florian.knoll@nyumc.org.

[†]Indicates equal contributions.

advances, the methods developed in these studies were all trained and validated on small individual datasets collected by the authors, which in many cases were not shared with the research community. These limitations in data accessibility makes it challenging to reproduce different approaches, and to validate comparisons between them. The lack of broadly accessible data also restricts work on important medical image reconstruction problems to researchers associated with or cooperating with large university medical centers where imaging data is available. This restriction is a significant lost opportunity, given the substantial volume of ongoing research in basic science machine learning and data science.

Indeed, there is a striking contrast between specialized medical research and more general research in the field of machine learning, which has seen breakthrough improvements in diverse areas from image classification [17] with deep convolutional neural networks (CNNs) [18] to championship-level gaming [19]. The core technologies that led to these results had already been introduced around 1990 for applications like speech recognition [20] and written document parsing [21]. However, deep learning for computer vision did not expand beyond simple digit recognition tasks for the next 20 years. In retrospect, a single event is often identified as the key catalyst for the recent resurgence of machine learning technology [18]: The 2012 ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [22], in which a deep CNN achieved spectacular results for an image classification task [17]. Since then, every single winning entry in the competition was a form of deep CNN, with current winning entries even outperforming human performance. Winning the ILSVRC has become extremely prestigious and has attracted the interest of leading academic institutions and IT companies around the world. Performance on ILSVRC tasks has become a standard for the evaluation of new developments in computer vision. A similar event occurred in the field of medical image reconstruction with the 2016 Low Dose CT Grand Challenge organized by the Mayo Clinic [23]. Even after the conclusion of the challenge, the dataset provided by the organizers continues to be widely used by research groups for their own developments. It serves as a standard reference in the CT community for reconstruction advances. More recently, a dataset with the goal to facilitate deep learning for low-dose-CT was also made available. [24]

Our goal with the fastMRI challenge project was provide a similar stimulation to machine learning research in MR image reconstruction aimed at reducing MR examination times. In December of 2018, we released the first large-scale database of MRI scanner raw data from a clinical patient population [25, 26]. In the spirit of previous challenges organized by the ISMRM community [27], we then conducted a challenge to provide researchers in the field the opportunity to evaluate their methods in a large-scale, realistic setting with evaluation from clinical radiologists. We also aimed to spark interest in radiology and biomedical imaging within the large machine learning and computer vision research community. In this article, we describe the design and the results of the challenge as well as the lessons we learned from its organization.

2 | METHODS

2.1 | Challenge design principles

Our challenge was focused on accelerating MR image acquisitions. Two of the most influential developments in this arena during the last two decades have been parallel imaging [28, 29, 30] and compressed sensing [31]. Both of these approaches to rapid imaging are based on the principle of reducing the number of lines that are acquired in k-space, which reduces the scan time, and then exploiting redundancy in the measured data during the image reconstruction process. In parallel imaging, the redundancy arises from the simultaneous acquisition of MR signal with multiple receive coils; in compressed sensing, it derives from the observation that images are generally compressible. Machine learning approaches have generally adopted similar strategies for the acceleration of MRI, which set the main design criteria for our challenge.

We provided participants with sets of raw k-space data, and the goal of the challenge was to reconstruct images from these data. Since details about the dataset are reported in separate publications [25, 26], in this article we restrict our description of the dataset only to information that is relevant to the design of the challenge. We provided data for a total of 1,594 consecutive clinical proton-density-weighted MRI acquisitions of the knee in the coronal plane, both with (COR PD FS) and without (COR PD) frequency-selective fat saturation. In addition to their different image contrast, these two types of acquisition also vary in signal to noise ratio (SNR) by approximately a factor of 4 [32]. Data were acquired on three clinical 3T systems (Siemens Magnetom Skyra, Prisma, and Biograph-mMR) and one clinical 1.5T system (Siemens Magnetom Aera) using clinical multi-channel receive coils. Curation of the dataset was part of a study approved by our local institutional review board (IRB).

The selection of problems for the challenge was based on a three-way trade-off between a) providing a realistic scenario representative of actual clinical imaging exams, b) allowing fair and proper validation, and c) making the challenge practically and conceptually accessible for research groups outside the core field of MR image reconstruction. This led to the following design principles:

- To make the image reconstruction problem realistic, we provided actual raw (complex valued) k-space data obtained directly from our MRI scanners.
- To reduce the complexity of the challenge, we restricted ourselves to standard Cartesian 2D Turbo Spin Echo sequences that are part of the routine clinical protocol at our institution.
- In order to provide clear ground truth against which to compare image reconstructions, we did not provide prospectively undersampled data. Fully-sampled k-space data were acquired for all exams in the data set, and undersampling was performed retrospectively, so that no differences in conditions (e.g., in motion state or scanner calibration) between fully-sampled and undersampled acquisitions would complicate image comparisons.

- Since the goal of the challenge was to test reconstruction methods and not sampling trajectory design, we predefined the allowed undersampling patterns. We chose one-dimensional pseudo-random sampling in the phase encoding direction, with full sampling of a small central k-space region, as introduced in the original context of compressed sensing [31].
- For multi-coil acquisitions, our ground truth reference was the root-sum-of-squares combination of the fully-sampled multi-channel data after inverse Fourier transform. While this is not the optimal coil combination method in terms of SNR [33, 34], it does not bias the ground truth towards any particular approach to the estimation of coil sensitivities. We also removed readout-direction oversampling by cropping the reconstructed images to the central 320×320 pixel region. For the single-coil case, which is uncommon in clinical practice but was included to provide a low barrier to entry for those not familiar with multi-coil data acquisitions, we simulated a physically feasible ground truth using a linear combination of individual coil signals as described in [35].
- We used the structural similarity index (SSIM) [36] with respect to the fully sampled ground truth reference as an indicator of image quality. We calculated two other widely used quantitative metrics for our online leaderboard: pixelwise normalized mean square error (NMSE) and peak signal to noise ratio (PSNR). However, since all of these metrics provide limited insight into the diagnostic quality of medical images, we decided to use a ranking by an expert panel of musculoskeletal radiologists as the final metric to determine the winning entries in the challenge.
- We did not prohibit the use of additional non-fastMRI data in the development and training of the submissions. However, all participants who chose to use additional data were required to state this at the time of submission.

2.2 | Challenge tracks

One of our design goals was to test the submissions in different operating modes defined by the level of acceleration. We also wanted to make the challenge interesting for research groups with a focus on MR image reconstruction as well as for groups based in machine learning, computer vision and image processing. We therefore decided to organize the challenge into multiple submission tracks.

Regarding the different levels of acceleration, the goal of the first scenario was to operate in a mode where we expected the reconstruction to be challenging, but where reconstructed images that might be acceptable for clinical diagnosis were likely to be feasible. Based on our previous experience with similar data [10, 32], we chose an undersampling factor of $R=4$ for this scenario. The goal of the second scenario was to aim for a substantially higher acceleration than can be achieved with current reconstruction methods. We chose an undersampling factor of $R=8$ for this scenario. We did not expect to receive submissions with clinically acceptable image quality at this high level of acceleration. The goal for this scenario was to evaluate the performance of the submissions when they were pushed beyond reasonable limits, and to analyze failure modes.

In our experience, the steepest component of the learning curve for use of (Cartesian) MR data in image reconstruction relates to the proper handling of multi-channel raw k-space data of the sort required for parallel imaging. We therefore designed two additional tracks, which we termed the multi-coil and the single-coil track. For the multi-coil track, which was primarily aimed at research groups with a background in MR image reconstruction, we provided true multi-channel raw data from the MR scanners. Since most modern MR scans are performed using arrays of detector coils, this is a realistic scenario whose results are likely to be readily translatable to real-world imaging situations. For the single-coil track, we provided k-space data for which multi-channel information had been combined into a single channel that can be reconstructed with a simple inverse Fourier transform in the fully sampled case. However, instead of Fourier transforming previously-reconstructed images stored in DICOM format in an attempt to create k-space data (an approach sometimes observed in the literature, but not realistic or advisable for various reasons), we chose to retain the complex nature of the original data. A more detailed description of the channel combination we used is presented in [35]. The single-coil track was primarily aimed at research groups from machine learning, computer vision and image processing, who might be interested in applying their expertise in medical imaging applications. In particular, after a simple inverse Fourier transform of the data, the single-coil track enabled easy use of methods that are entirely based on image postprocessing. While the removal of multi-channel information decreases the complexity of working with the data, it also increases the difficulty of the reconstruction problem at any given acceleration factor, since the resulting single-channel data has reduced redundancy and more limited information content than the original multi-channel data. For the challenge phase, we therefore decided to limit ourselves to a single acceleration factor of $R=4$ in this track.

2.3 | Dataset split and leaderboard evaluation

We partitioned our dataset into six subsets for the individual tracks and the different phases of the challenge: training, validation, multi-coil test, single-coil test, multi-coil challenge, or single-coil challenge. Data from individual patient cases were randomly assigned to the individual subsets. Each patient case consists of an image volume of coronal images. The number of cases and the total number of slices are shown in Table 1. The dataset was made publicly available at <https://fastmri.med.nyu.edu/>.

We provided fully sampled k-space data and corresponding ground truth image reconstructions for the training and validation subsets, which could be used by the participants to develop and train their machine learning models and to determine any hyperparameters. In order to make most efficient use of our available data, we used the same training and validation cases for the multi-coil and single-coil tracks.

For the test set, we provided different subsets of undersampled k-space data for the single and the multi-channel tracks. Participants could upload their reconstruction results to our public leaderboard at <http://fastmri.org/>. We then calculated three commonly used quantitative image quality metrics: The first metric was the normalized mean square error (NMSE), defined as:

$$\text{NMSE}(\hat{u}, u) = \frac{\|\hat{u} - u\|_2^2}{\|u\|_2^2},$$

where \hat{u} is a reconstructed image volume and u is the fully sampled ground truth reference image volume. $\|\cdot\|_2^2$ is the squared Euclidean norm, and the subtraction is performed entry-wise. We computed NMSE values normalized over full image volumes rather than individual slices, since image-wise normalization can result in strong variations across a volume.

The second metric was peak signal-to-noise ratio (PSNR), which represents the ratio between the power of the maximum possible image intensity across a volume and the power of distorting noise and other errors:

$$\text{PSNR}(\hat{u}, u) = 10 \log_{10} \frac{\max(u)^2}{\text{MSE}(\hat{u}, u)}.$$

Here $\text{MSE}(\hat{u}, u)$ is the mean square error between \hat{u} and u defined as $\frac{1}{n} \|\hat{u} - u\|_2^2$ and n is the number of entries in the target volume u . Higher values of PSNR (as opposed to lower values of NMSE) indicate a better reconstruction.

The third metric was the structural similarity index (SSIM) [36], which measures the similarity between two images by exploiting the inter-dependencies among nearby pixels. SSIM is inherently able to evaluate structural properties of the objects in an image and is computed at different image locations by using a sliding window. The resulting similarity between two image patches \hat{m} and m is defined as

$$\text{SSIM}(\hat{m}, m) = \frac{(2\mu_{\hat{m}}\mu_m + c_1)(2\sigma_{\hat{m}m} + c_2)}{(\mu_{\hat{m}}^2 + \mu_m^2 + c_1)(\sigma_{\hat{m}}^2 + \sigma_m^2 + c_2)},$$

where $\mu_{\hat{m}}$ and μ_m are the average pixel intensities in \hat{m} and m , $\sigma_{\hat{m}}^2$ and σ_m^2 are their variances, $\sigma_{\hat{m}m}$ is the covariance between \hat{m} and m and c_1 and c_2 are two variables to stabilize the division; $c_1 = (k_1 L)^2$ and $c_2 = (k_2 L)^2$. For the evaluation of the challenge, we chose a window size of 7×7 , we set $k_1 = 0.01$, $k_2 = 0.03$, and defined L as the maximum value of the target volume, $L = \max(v)$.

We evaluated the performance on the complete test dataset as well as individual errors for the two image contrasts (with and without fat suppression. Participants could submit to each track leaderboard once a day before the submission deadline. Submissions were ranked by SSIM of R=8 undersampling. The challenge dataset was released on September 5th, 2019 (see below for a description of the timeline), and the submission window was then open for 14 days. The evaluation and the structure of the leaderboard for the challenge phase were identical to those for the test phase, but each team could only make one challenge submission, and challenge results were made available only after the submission window

was closed. A screenshot of the leaderboard for the multi-channel track of the completed challenge is shown in Figure 1.

2.4 | Challenge timeline and design of the evaluation

The challenge consisted of multiple phases according the following timeline:

- November 26, 2018: Release of the training and validation sections of the fastMRI dataset [25, 26].
- June 5, 2019: Official announcement of the challenge and release of the test set. The test set leaderboard was open for submission at this stage.
- September 5 to 19, 2019: Release of the challenge dataset and challenge submission window.
- September 19, 2019: Quantitative evaluation of the challenge submissions. We selected the top 4 submissions with highest SSIM on the challenge dataset from each track for the second phase of evaluation by a panel of radiologists. At this stage, we also asked all participants to provide an abstract for the 2019 Medical Imaging Meets NeurIPS workshop¹.
- September 20 to October 10, 2019: Radiologist evaluation phase. We sent 5 randomly selected cases (with and without fat suppression) from the top 4 submissions in each track plus the corresponding ground truth reconstructions to our panel of seven radiologists from multiple institutions, including NYU Langone Health, Cleveland Clinic, University of California San Diego, University of Wisconsin and Stanford University. Each radiologist looked at a total of 1840 images. We asked the panel to rank the submissions in terms of the overall image quality to select one winner in each track. We then averaged the rankings of the radiologists to determine the winners. In addition, we asked the radiologists to score each submission on a 4-point scale (1 is best and 4 is worst) for the following criteria: Presence of artifacts, image sharpness, perceived contrast-to-noise ratio and diagnostic confidence. This rating was performed to obtain additional meta-information about the readers' preferences, and to give the radiologists some suggestions on which to base their ranking. Subcriterion rankings were not directly used to determine the winners of the challenge. At the end of this stage, we notified the winners of the three tracks and shared their abstracts and identity with the organizers of the 2019 Medical Imaging Meets NeurIPS workshop.
- December 1, 2019: Publication of the challenge leaderboard with the results of the quantitative evaluation.
- December 14, 2019: Official announcement of the winners of the three tracks at the 2019 Medical Imaging Meets NeurIPS workshop, with oral presentation by the three winning teams.

3 | RESULTS

3.1 | Overview of submissions

We received a total of 33 challenge submissions. 8 groups submitted to the multi-coil track, each with submissions for both the R=4 and the R=8 tracks. At the time of this writing, four of the submitting groups have published manuscripts on their approach, in addition to their NeurIPS abstracts: *Sigma – Net* from team holyk-space [37], iRim from team AI Amsterdam [38], the pyramid convolutional RNN from team MSDC-RNN [39] and Adaptive CS-Net from Philips and LUMC [40, 41]. 17 groups submitted to the single coil R=4 track. 6 out of the 8 groups who submitted to the multi-coil track also submitted to the single-coil track. We did not require that the groups publicly disclose their names or affiliations at the submission stage. This was only required for the winners of each track. For the test-set leaderboard during the full duration of the challenge, we received more than 25 submissions for the multi-coil track and more than 70 submissions for the single coil track. All submissions used exclusively fastMRI data in the training of their approach for both challenge and public test submissions. We encouraged all participants to provide open source code together with their submissions, and 3 groups provided links to their open source code repositories.

3.2 | Analysis of results

The SSIM scores of the challenge submissions are shown in Figure 1. As expected, there is a substantial difference in overall SSIM values between the multi-coil and the single-coil tracks. The average SSIM of all submissions was 0.924, 0.895 and 0.707 for the multi-coil R=4, multi-coil R=8 and single-coil R=4 tracks, respectively. Even the lowest-ranking multi-coil R=8 submission (SSIM=0.874) significantly outperformed the highest-ranking single-coil R=4 submission (SSIM=0.754).

Figure 2 shows selected results from the Multi-Coil R=4 track, for one particular slice with and one slice without fat suppression. Both of these cases were obtained from 1.5T systems (Siemens Magnetom Aera). Corresponding difference images to the fully sampled ground truth are shown in the supporting material. The top 4 submissions that were evaluated by the radiologists are ordered from left to right based on the radiologists rankings, next to the ground truth on the far left. The average rank of the 7 radiologists is displayed on top of each submission. The SSIM to the ground truth for each particular slice is shown in the bottom left of the plots. The case in the top row shows a subtle subchondral osteophyte, which was not visible in the accelerated reconstructions. In addition to the submissions, we are including a combined parallel imaging and compressed sensing (PI-CS) reconstruction using total generalized variation [42] for reference. All data processing, including estimation of the coil sensitivities, was done exactly as described in the referenced paper. The reference results can be reproduced with the corresponding software package that is available online². We optimized the regularization parameter of the reference reconstruction individually for each example such that SSIM was maximized. We have chosen this approach to ensure that the submissions are presented in context of the best possible reference reconstruction, and their superiority cannot simply be explained by a bad reference hyper-parameter setting. We would like to stress here that such an optimization requires the availability of the fully sampled ground truth, which was available to us but of course not to the participants. In

particular, for the results in Figure 2, we used a regularization parameter of 10^{-5} for the non fat-suppressed example, and 10^{-4} for the fat-suppressed example.

Figure 3 shows results for the Multi-Coil R=8 track. Corresponding difference images to the fully sampled ground truth are shown in the supporting material. The optimized regularization parameter for the PI-CS reference was 10^{-4} for both the non fat-suppressed and the fat-suppressed example. The case in the top row shows moderate artifact from a metal implant and was obtained on a 1.5 system (Siemens Magnetom Aera). None of the submissions was negatively affected by this irregularity. The case in the bottom row shows a meniscal tear. It was acquired on a 3T system (Siemens Magnetom Prisma). This pathology was not visible in the accelerated reconstructions.

Figure 4 shows results for the Single-Coil R=4 track. Both of these cases were obtained from 3T systems (Siemens Magnetom Skyra). For the single coil track, we performed a corresponding compressed sensing reconstruction without parallel imaging, again using a TGV regularizer. The optimized regularization parameters were 10^{-5} for the non fat-suppressed example, and 10^{-7} for the fat-suppressed example. Corresponding difference images to the fully sampled ground truth are shown in the supporting material.

One of the open questions in machine learning for image reconstruction is the assessment of failure modes of trained models, and whether there are certain situations where a particular model will outright fail [43]. To contribute to the investigation of this phenomenon, we identified the cases of our challenge where the individual submissions performed worst. Figure 5 shows the results of the cases with lowest SSIM over the whole 3D image volume of the challenge dataset from each track. The chosen regularization parameters for the PI-CS reference were 10^{-4} for Multi-Coil 4x and 8x, and 10^{-7} for Single-Coil 4x. Corresponding difference images to the fully sampled ground truth are shown in the supporting material. Notably, all methods performed worst on the same case within each track. The two cases for the two multi-coil tracks were obtained from 1.5T systems (Siemens Magnetom Aera), the case for the single-coil-track was from a 3T system (Siemens Magnetom Skyra). While the SSIM values are substantially lower for these cases, the results of the different methods are remarkably consistent and no obvious negative outlier in terms of image quality can be identified. In addition, even on the cases where they performed worst, all methods still outperformed the PI-CS reference.

Table 2 shows the average radiologist rankings as well as the overall SSIM, NMSE and PSNR values for the full challenge dataset for the top 4 submissions of each track. Figure 6 shows corresponding scatterplots after normalization of the scores (1 is best). In the case of multi-coil R=8, the highest ranked submission was also the one that had the highest SSIM, NMSE and PSNR values. For single-coil R=4, only SSIM showed a similar trend as the radiologists scores, while the other two metrics showed almost opposite trends. For the multi-coil R=4 track, the results are less conclusive. The top 4 submissions were very close together with all metrics. The differences in SSIM between the submissions were less than 1% in this track.

Additional insight into the radiologists' ratings is provided by Figure 7, which shows the individual rankings by the 7 radiologists for the top 4 submissions in all three submission tracks. For multi-coil R=8 and single-coil R=4 tracks, the radiologists had a strong preference for a single submission. The highest-rated submission in each of these tracks was ranked first by 5 radiologists, and ranked second by the remaining 2 radiologists. The results are substantially less consistent for the multi-coil R=4 track. The two highest-rated submissions each were also ranked worst by one reader. The lowest-rated submission was actually ranked best by 3 out of 7 radiologists, and ranked worst by the remaining 4.

Table 3 shows the average scores for the individual categories that the radiologists were asked to rate for the top 4 submissions in each track: Artifacts, sharpness, perceived contrast-to-noise ratio and diagnostic confidence, using a 4 point scale, where 1 is best and 4 is worst. These categories were intended as guidelines for radiologist ranking, not as strict criteria. However, by and large the radiologists chose to rank the submissions based on the sum of their scores in the different categories. For the multi-coil R=8 track, all radiologists ranked the submissions strictly based on their scores. 5 out of 7 radiologists for the multi-coil R=4 track and 6 out of 7 radiologists for the single-coil R=4 track ranked submissions strictly based on their scores. Even for the cases where the ranking deviated from the scores, the top-ranked submission always had the best scores as well.

4 | DISCUSSION

4.1 | Limits of the challenge design

The goal of our experimental design was that it is representative of clinical performance of the submissions. However, since this was the first time this challenge was held, we wanted to keep the evaluation as simple as possible at the same time. One of our most consequential decisions in terms of challenge design was to not generate any systematic differences between training, validation, test and challenge sets. All of these datasets were randomly selected from the same superset of data, and all consisted of coronal knee data from a limited set of MR scanners from a single vendor. This design substantially limits insight into robustness and generalization. It is possible to subsequently perform a more targeted analysis by, for example, only using a subset of the training data from one of the two contrasts or one field strength (1.5T or 3T) for training and validation, and the other set of data for testing. Given the importance of multi-coil data in the overall performance of the submissions, the challenge also didn't include substantial variations of receive coil geometries. All coil arrays were standard knee coil configurations from a single vendor, with the same number of receive channels (15).

Compressed sensing [31] relies on incoherence to remove aliasing artifacts. It is still an open question to what degree certain sampling patterns are beneficial for machine learning methods [44]. We chose one-dimensional pseudo-random sampling in the phase encoding direction for simplicity. We did not consider recent developments in sampling pattern design [45], patterns that are specifically designed to combine parallel imaging and compressed sensing [46], or approaches that use machine learning directly for the optimization of the sampling pattern [47, 48]. We also did not specifically consider how easily these particular patterns can be implemented prospectively in terms of scanner hardware and pulse-sequence

constraints. Therefore, our challenge does not allow to draw conclusions with respect to these questions.

In addition to maintaining homogeneity of the data on a technical level, we also decided not to perform curation of the data in terms of anatomical or pathological variations. It would be interesting to separate pathological from non-pathological cases and use only one of these individual subsets for training and validation, and the second subset for testing. Aside from pathology, similar experiments could be performed by grouping subsets of data based on age, height, weight, body mass index or gender. Experiments like these should be considered for future challenges in the area of machine learning for MR image reconstruction.

Another substantial limitation of our design was that in the radiologists evaluation, we only asked the readers to rate image submissions by image quality on a subjective level. Therefore, the evaluation is not unbiased regarding personal preferences in image quality. Ultimately, the goal of medical imaging is to provide diagnostic information to a referring physician. In the example of the meniscal tear that is shown in Figure 3, the relevant information content of the image is whether the radiologist correctly identifies this tear and makes the recommendation for the correct follow up procedure (for example, surgery vs no surgery). Therefore, a more objective evaluation would have been to ask our readers to provide a diagnostic reading in the same way they would do it in their clinical practice. Such an evaluation can then be done blinded for each individual submission as well as for the fully sampled ground truth reconstruction. The measure of success for a certain reconstruction method is then to what degree is the obtained diagnosis consistent with the diagnosis from the ground truth [49]. However, since we relied on voluntary radiologists with limited time, such an evaluation was outside the scope of our challenge.

4.2 | Analysis of the submissions and results

The quantitative SSIM values (Figure 1) provide several interesting insights. First, the differences in SSIM values between submissions from the top teams are almost negligible. In each of the three tracks, the difference between the first and the fourth-ranked entry was less than 1%. For the multi-coil R=4 track, the radiologists' scoring showed a similar trend. Both the top two and the bottom two submissions were tied in the ranking. However, since all participants decided to use only NYU-provided training data, no conclusions can be drawn about potential improvements by using additional training data, either by expanding the dataset with additional knee data, or by using synthetic data and transfer learning. When it comes to the visual impression of the reconstructed images, we found it remarkable to what degree and how constantly all submissions outperformed a conventional parallel imaging and compressed sensing reconstruction (Figures 2 to 5).

It is often pointed out in the medical imaging community that quantitative metrics like NMSE, PSNR and SSIM are poor metrics to evaluate the quality of medical images. In our challenge, juxtaposition of the radiologists' scores with SSIM values (Figure 6) shows that for the two tracks where the radiologists picked a clear winner (multi-coil R=8 and single-coil R=4), the winner was also the submission with the highest SSIM value. NMSE and PSNR were aligned with this trend for multi-coil R=8, but for single-coil R=4, NMSE actually resulted in the opposite ranking order from that selected by the radiologists. While

none of the metrics in any track resulted in the same rank order as the radiologists, it is important to remember that the quantitative values were very closely spaced. It is interesting that for the track where the SSIM values of the top 4 submissions were essentially identical (multi-coil R=4), the individual radiologists also had substantial disagreement in their preference (Figure 7). The lowest-ranked submission was actually ranked best by 3 out of 7 radiologists, and ranked worst by the remaining 4. This indicates that for the multi-coil R=4 track, the submissions were most likely identical in terms of image quality, as correctly predicted by SSIM, and the ranking was determined by individual preferences for image quality by the radiologists. This means that in our challenge, SSIM actually did provide estimates of image quality that were consistent with the preferences of radiologists. Our results also suggest that radiologists' evaluations must be carried out at the level of diagnostic interpretation to allow their domain knowledge to provide substantial additional information.

While we knew that there would be a difference in performance between the single-coil and the multi-coil tracks, we were surprised by the degree of difference actually observed. From a linear algebraic point of view, the underlying problems in the different tracks are substantially different. The undersampled single-coil reconstruction problem is an undetermined system in which data acquisition violates the Shannon/Nyquist sampling theorem, and a solution can only be obtained by introducing prior knowledge and performing incoherent sampling, on the model of compressed sensing [31]. By contrast, for the multi-coil problem, even at R=8 acceleration, the number of receive channels (15) is still higher than the undersampling factor. The underlying problem involves an overdetermined system. However, the problem is ill-posed because the individual coil elements do not provide independent information, and prior knowledge is also needed to constrain the solution. While this may raise the question of whether fundamentally different approaches should be developed for the two scenarios, the results of the challenge indicate otherwise. Six out of eight participants in the multi-coil track also submitted to the single-coil track. Team holyk-space [37] and team AM used a dedicated approach for multi-coil data. Team AM explicitly estimated coil sensitivity maps using Espirit [50] and used a nullspace constraint on the fully-sampled center of k-space that is used to estimate the coil sensitivities. Team holyk-space learned the implicit weighting of the individual coils. In contrast, the remaining groups used essentially the same core method for both tracks, and only fine-tuned and re-trained for the different tracks. Also, the top three submissions in the single-coil track were from the same groups that submitted to the multi-coil track. As expected, the number of submissions for the single-coil track was substantially higher than for the multi-coil track, most likely due to the shallower learning curve and greater ease of use of the single-coil data. However, the results from the challenge show that in order to achieve the best possible image quality for accelerated MR scans, it is essential to take the multi-channel nature of MR acquisitions into account. Therefore, we plan to limit ourselves to multi-coil tracks for future iterations of our challenge.

The inability to correctly identify subtle pathology, even in the multi-coil R=4 results in Figure 2, must be considered in the light of clinical adoption. However, loss of low-contrast fine details is not necessarily a particular culprit of the supervised CNN learning-based reconstruction methods that formed the finalists in our challenge. The pathology in question

was also lost in the reference PI-CS reconstruction, and it is entirely possible that it would have been lost with any reconstruction approach at this level of acceleration. On the other hand, a common fear about machine learning reconstruction methods is that they react very unpredictably and unstably for cases that show severe abnormalities or deviations from normal anatomy. In our challenge, none of submissions showed any kind of deterioration for the case with the severe image artifact due to the metallic implant in the Multi-Coil R=8 track (Figure 3). This result is also reinforced with the worst-case results (Figure 5), which also didn't show severe outliers in terms of image quality, and all submissions still outperformed the parallel imaging compressed sensing reference. It is notable that the worst-case results were simply cases where the baseline ground truth SNR was lower and it can be expected that image quality is lower, and not cases that show a particular anatomical or pathological variation or a deviation in image contrast.. Separate dedicated studies will be required to investigate this effect, but this result is still encouraging from the point of view of robustness for clinical translation. Another common fear about machine learning reconstruction is that when the acceleration is pushed substantially, the CNNs create hallucinations of artificial structures that could be falsely identified as pathology. The results from our Multi-Coil R=8 track show that this was not the case for the submissions to our challenge. However, none of the finalists used an architecture that included a generative adversarial model [51], where the potential for such an error is generally considered highest.

All submissions used a supervised learning approach with deep Neural Networks. While it is tempting to conclude that a similar paradigm shift towards deep learning has occurred for MR image reconstruction as in the ImageNet challenge [22] for computer vision, in our opinion the results of our challenge do not allow us to draw that conclusion. First, the total number of submissions for our challenge was substantially smaller than for ILSVRC, and it was only the first time the challenge was held. Second, as described above, the design of the challenge essentially guaranteed that (supervised) machine learning methods would have strong performance on the challenge dataset. Third, in contrast to purely data-driven end-to-end learning in the true spirit of deep learning [18], the winners of all three tracks chose approaches that used a combination of a learned prior and a data-fidelity term that encodes information about the MR physics of the acquisition, in line with approaches that can be seen as neural network extensions of classic iterative image reconstruction methods [52, 10, 13, 53, 14, 54, 55]. Therefore, the submissions only covered a rather narrow part of the whole spectrum of MR image reconstruction methods, and a large number of alternative approaches both learning and non-learning based, exist. Finally, even though radiologist ratings were ultimately the deciding factor that determined the winners, and while they did have the fully sampled ground truth available as a reference, their ratings were essentially based on subjective impression of image quality and not on diagnostic equivalency. The translation of machine learning for reconstruction of accelerated MRI scans in routine clinical practice remains an open question for future research and development.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

We first would like to thank all participants of the challenge. We thank the radiologists who provided the scoring for the second evaluation phase: Drs. Christine Chung and Mini Pathria of UCSD, Dr. Michael Tuite of University of Wisconsin, Dr. Christopher Beaulieu of Stanford, Drs. Naveen Subhas and Hakan Iltaslan of the Cleveland Clinic, and Dr. David Rubin of NYU Langone Health. We thank our external advisors for the organization of the challenge: Dr. Daniel Rueckert of Imperial College London, Dr. Jonathan Tamir of University of Texas at Austin, Dr. Joseph Cheng of Apple AI research and Dr. Frank Ong of Stanford. We also thank our colleagues Mark Tygert, Michal Drozdal, Adriana Romero, Pascal Vincent, Erich Owens, Krzysztof Geras, Patricia Johnson, Mary Bruno, Jakob Asslaender, Yvonne Lui, Zhengnan Huang and Ruben Stern for their insights and feedback. We acknowledge grant support from the National Institutes of Health under grants NIH R01EB024532 and NIH P41EB017183.

Funding information

NIH: NIBIB, Awards NIH R01EB024532 and NIH P41EB017183.

references

- [1]. Kang E, Min J, Ye JC. A deep convolutional neural network using directional wavelets for low-dose X-ray CT reconstruction. *Medical Physics* 2017;44(10):e360–e375. [PubMed: 29027238]
- [2]. Chen H, Zhang Y, Kalra MK, Lin F, Chen Y, Liao P, et al. Low-Dose CT with a residual encoder-decoder convolutional neural network. *IEEE Transactions on Medical Imaging* 2017;.
- [3]. Jin KH, McCann MT, Froustey E, Unser M. Deep convolutional neural network for inverse problems in imaging. *IEEE Transactions on Image Processing* 2017;.
- [4]. Wolterink JM, Leiner T, Viergever MA, Isgum I. Generative Adversarial Networks for Noise Reduction in Low-Dose CT. *IEEE transactions on medical imaging* 2017;36(12):2536–2545. [PubMed: 28574346]
- [5]. Kobler E, Klatzer T, Hammernik K, Pock T. Variational Networks: Connecting Variational Methods and Deep Learning. In: *Proceedings of the German Conference on Pattern Recognition (GCPR)*; 2017 p. 281–293.
- [6]. Adler J, Öktem O. Learned Primal-Dual Reconstruction. *IEEE Transactions on Medical Imaging* 2018;37(6):1322–1332. <https://arxiv.org/pdf/1707.06474.pdf>. [PubMed: 29870362]
- [7]. Liu Y, Zhang Y. Low-dose CT restoration via stacked sparse denoising autoencoders. *Neurocomputing* 2018;284:80–89. 10.1016/j.neucom.2018.01.015.
- [8]. Hammernik K, Knoll F, Sodickson DK, Pock T. Learning a Variational Model for Compressed Sensing MRI Reconstruction. In: *Proceedings of the International Society of Magnetic Resonance in Medicine (ISMRM)*; 2016 p. 1088.
- [9]. Wang G A perspective on deep imaging. *IEEE Access* 2016;4:8914–8924.
- [10]. Hammernik K, Klatzer T, Kobler E, Recht MP, Sodickson DK, Pock T, et al. Learning a Variational Network for Reconstruction of Accelerated MRI Data. *Magnetic Resonance in Medicine* 2018;79(6):3055–3071. <http://arxiv.org/abs/1704.00447> [PubMed: 29115689]
- [11]. Aggarwal HK, Mani MP, Jacob M. MoDL: Model-Based Deep Learning Architecture for Inverse Problems. *IEEE Transactions on Medical Imaging* 2019;.
- [12]. Ye JC, Han Y, Cha E. Deep Convolutional Framelets: A General Deep Learning Framework for Inverse Problems. *SIAM Journal in Imaging Sciences* 2018;11(2):991–1048.
- [13]. Schlemper J, Caballero J, Hajnal JV, Price AN, Rueckert D. A Deep Cascade of Convolutional Neural Networks for Dynamic MR Image Reconstruction. *IEEE Transactions on Medical Imaging* 2018;.
- [14]. Qin C, Schlemper J, Caballero J, Price AN, Hajnal JV, Rueckert D. Convolutional recurrent neural networks for dynamic MR image reconstruction. *IEEE Transactions on Medical Imaging* 2019;.
- [15]. Zhu B, Liu JZ, Cauley SF, Rosen BR, Rosen MS. Image reconstruction by domain-transform manifold learning. *Nature* 2018;555(7697):487–492. <http://www.nature.com/doifinder/10.1038/nature25988> [PubMed: 29565357]

- [16]. Chen H, Zhang Y, Chen Y, Zhang J, Zhang W, Sun H, et al. LEARN: Learned Experts' Assessment-Based Reconstruction Network for Sparse-Data CT. *IEEE Transactions on Medical Imaging* 2018;37(6):1333–1347. [PubMed: 29870363]
- [17]. Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*; 2012 p. 1097–1105.
- [18]. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature* 2015 5;521(7553):436–444. 10.1038/nature14539. [PubMed: 26017442]
- [19]. Silver D, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G, et al. Mastering the game of Go with deep neural networks and tree search. *Nature* 2016 1;529(7587):484–489. 10.1038/nature14539. [PubMed: 26819042]
- [20]. Waibel A, Hanazawa T, Hinton G, Shikano K, Lang KJ. Phoneme Recognition Using Time-Delay Neural Networks. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 1989;.
- [21]. LeCun Y Handwritten Digit Recognition with a Back-Propagation Network. *Advances in Neural Information Processing Systems* 1989;.
- [22]. Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, et al. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* 2015;.
- [23]. McCollough CH, Bartley AC, Carter RE, Chen B, Drees TA, Edwards P, et al. Low-dose CT for the detection and classification of metastatic liver lesions: Results of the 2016 Low Dose CT Grand Challenge. *Medical Physics* 2017;.
- [24]. Leuschner J, Schmidt M, Baguer DO, Maaÿ P. The LoDoPaB-CT Dataset: A Benchmark Dataset for Low-Dose CT Reconstruction Methods 2019 10;<http://arxiv.org/abs/1910.01113>.
- [25]. Zbontar J, Knoll F, Sriram A, Muckley MJ, Bruno M, Defazio A, et al. {fastMRI: An} Open Dataset and Benchmarks for Accelerated {M}{R}{I}. arXiv:181108839 preprint 2018;.
- [26]. Knoll F, Zbontar J, Sriram A, Muckley MJ, Bruno M, Defazio A, et al. {fastMRI:} a publicly available raw k-space and {DICOM} dataset for accelerated {MR} image reconstruction using machine learning. *Radiology Artificial Intelligence* 2019;in press.
- [27]. Grissom WA, Setsompop K, Hurley SA, Tsao J, Velikina JV, Samsonov AA. Advancing RF pulse design using an open-competition format: Report from the 2015 ISMRM challenge. *Magnetic Resonance in Medicine* 2017;.
- [28]. Sodickson DK, Manning WJ. Simultaneous acquisition of spatial harmonics ({SMASH}): fast imaging with radiofrequency coil arrays. *Magn Reson Med* 1997 10;38(4):591–603. [PubMed: 9324327]
- [29]. Pruessmann KP, Weiger M, Scheidegger MB, Boesiger P. {SENSE}: sensitivity encoding for fast {MRI}. *Magn Reson Med* 1999 11;42(5):952–962. [PubMed: 10542355]
- [30]. Griswold MA, Blaimer M, Breuer F, Heidemann RM, Mueller M, Jakob PM. Parallel magnetic resonance imaging using the GRAPPA operator formalism. *Magn Reson Med* 2005 12;54(6):1553–1556. [PubMed: 16254956]
- [31]. Lustig M, Donoho D, Pauly JM. Sparse {MRI}: The application of compressed sensing for rapid {MR} imaging. *Magn Reson Med* 2007 12;58(6):1182–1195. [PubMed: 17969013]
- [32]. Knoll F, Hammernik K, Kobler E, Pock T, Recht MP, Sodickson DK. Assessment of the generalization of learned image reconstruction and the potential for transfer learning. *Magnetic Resonance in Medicine* 2019;.
- [33]. Roemer PB, Edelstein WA, Hayes CE, Souza SP, Mueller OM. The NMR phased array. *Magn Reson Med* 1990 11;16(2):192–225. [PubMed: 2266841]
- [34]. Walsh DO, Gmitro AF, Marcellin MW. Adaptive reconstruction of phased array {MR} imagery. *Magn Reson Med* 2000 5;43(5):682–690. [PubMed: 10800033]
- [35]. Tygert M, Zbontar J. Simulating single-coil MRI from the responses of multiple coils. arXiv 2018;.
- [36]. Wang Z, Bovik AC, Sheikh HR, Simoncelli EP. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing* 2004;13(4):600–612. [PubMed: 15376593]
- [37]. Schlemper J, Qin C, Duan J, Summers RM, Hammernik K, Sigma-net: Ensembled Iterative Deep Neural Networks for Accelerated Parallel MR Image Reconstruction. arXiv; 2019.

- [38]. Putzky P, Karkalousos D, Teuwen J, Miriakov N, Bakker B, Caan M, et al., i-RIM applied to the fastMRI challenge. arXiv; 2019.
- [39]. Wang P, Chen EZ, Chen T, Patel VM, Sun S. Pyramid Convolutional RNN for MRI Reconstruction 2019 12;<http://arxiv.org/abs/1912.00543>.
- [40]. Pezzotti N, de Weerd E, Yousefi S, Elmahdy MS, van Gemert J, Schülke C, et al. Adaptive-CS-Net: FastMRI with Adaptive Intelligence 2019 12;<https://arxiv.org/abs/1912012259>.
- [41]. Pezzotti N, Yousefi S, Elmahdy MS, van Gemert J, Schülke C, Doneva M, et al. An Adaptive Intelligence Algorithm for Undersampled Knee MRI Reconstruction: Application to the 2019 fastMRI Challenge 2020 4;<http://arxiv.org/abs/2004.07339>.
- [42]. Knoll F, Bredies K, Pock T, Stollberger R. Second order total generalized variation (TGV) for MRI. *Magnetic Resonance in Medicine* 2011;65(2):480–491. [PubMed: 21264937]
- [43]. Antun V, Renna F, Poon C, Adcock B, Hansen AC. On instabilities of deep learning in image reconstruction - Does AI come at a cost? 2019 2;<http://arxiv.org/abs/1902.05300>.
- [44]. Hammernik K, Knoll F, Sodickson D, Pock T. On the Influence of Sampling Pattern Design on Deep Learning-Based MRI Reconstruction. In: *Proceedings of the International Society of Magnetic Resonance in Medicine (ISMRM)*; 2017 p. 644.
- [45]. Adcock B, Hansen AC, Poon C, Roman B. Breaking the coherence barrier: A new theory for compressed sensing. *Forum of Mathematics, Sigma* 2017;5.
- [46]. Vasanawala SS, Murphy MJ, Alley MT, Lai P, Keutzer K, Pauly JM, et al. Practical parallel imaging compressed sensing MRI: Summary of two years of experience in accelerating body MRI of pediatric patients. In: *Proceedings - International Symposium on Biomedical Imaging*; 2011. .
- [47]. Seeger M, Nickisch H, Pohmann R, Schölkopf B. Optimization of k-space trajectories for compressed sensing by Bayesian experimental design. *Magn Reson Med* 2009 10;10.1002/mrm.22180.
- [48]. Zhu B, Liu JZ, Koonjoo N, Rosen BR, Rosen MS. AUTOMated pulse SEquence generation (AUTOSEQ) using Bayesian reinforcement learning in an MRI physics simulation environment. In: *ISMRM*; 2018.
- [49]. Obuchowski NA, Subhas N, Schoenhagen P. Testing for Interchangeability of Imaging Tests. *Academic Radiology* 2014;21(11):1483–1489. <http://www.sciencedirect.com/science/article/pii/S107663324002499>. [PubMed: 25300725]
- [50]. Uecker M, Lai P, Murphy MJ, Virtue P, Elad M, Pauly JM, et al. ESPIRiT an eigenvalue approach to autocalibrating parallel MRI: where SENSE meets GRAPPA In: *Magnetic Resonance in Medicine*, vol. 71 Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, California, USA; 2014 p. 990–1001. 10.1002/mrm.24751. [PubMed: 23649942]
- [51]. Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative Adversarial Nets. *Advances in Neural Information Processing Systems* 27 2014;p. 2672–2680. <http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf>.
- [52]. Yang Y, Sun J, Li H, Xu Z. ADMM-Net: A Deep Learning Approach for Compressive Sensing MRI. *Nips* 2017;(Nips):10–18. <http://arxiv.org/abs/1705.06869>.
- [53]. Cheng JY, Mardani M, Alley MT, Pauly JM, Vasanawala SS. Deep{SPIRiT}: Generalized Parallel Imaging using Deep Convolutional Neural Networks In: *Proceedings of the 26th Annual Meeting of the ISMRM, Paris, France*; 2018. .
- [54]. Aggarwal HK, Mani MP, Jacob M. MoDL: Model Based Deep Learning Architecture for Inverse Problems. *IEEE Transactions on Medical Imaging* 2018;p. Early view. <http://arxiv.org/abs/1712.02862>.
- [55]. Chun IY, Fessler JA. Deep BCD-Net Using Identical Encoding-Decoding CNN Structures for Iterative Image Recovery 2018;(1):1–5. <http://arxiv.org/abs/1802.07129>.

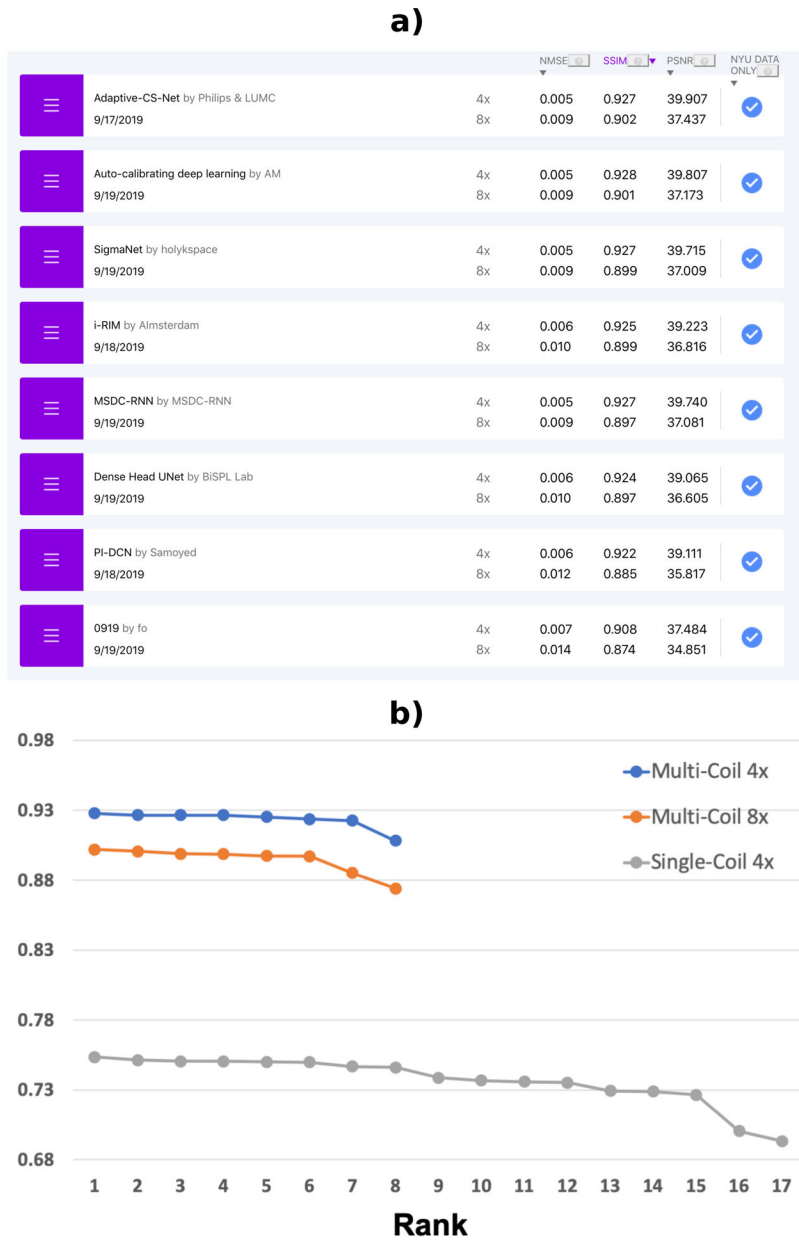


FIGURE 1.
 a) Online leaderboard at the completion of the challenge (December 2019). Three quantitative metrics are provided for R=4 and R=8 for both image contrasts on the webpage, along with a selection of reconstructed images. A short description of the reconstruction approach used, together with links to a corresponding paper or code repository, are also shown if the submitting groups provide this information. b) SSIM scores of the challenge submissions for each track. As expected, there is a substantial difference in overall SSIM values between the multi-coil and the single-coil tracks.

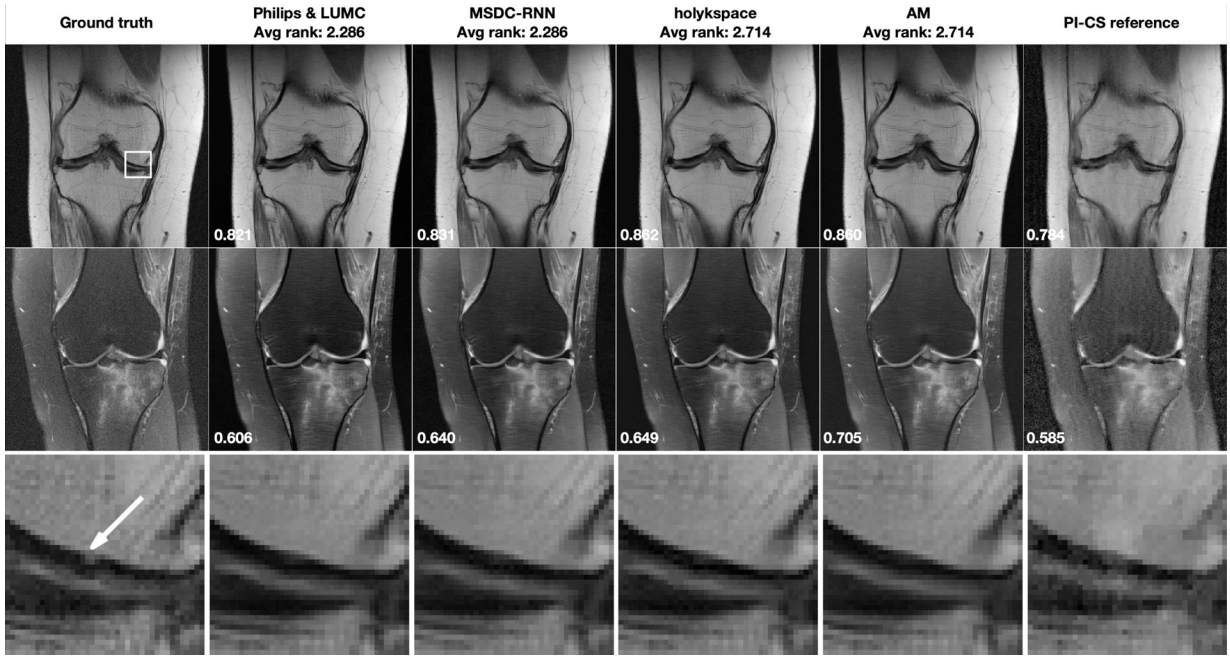


FIGURE 2. Multi-Coil R=4 track results: Selected results from the top 4 submissions in each track, for both image contrasts. The submissions are ordered from left to right based on the average of radiologists’ rankings. A combined parallel imaging and compressed sensing reconstruction using Total Generalized Variation (PI-CS) is shown for reference. SSIM to the ground truth for this particular slice is displayed in the bottom-left corner of each image. First row: Results for one slice from an acquisition without fat suppression. This case shows subtle pathology in the ROI indicated by a white rectangle in the ground truth image. Second row: One slice from an acquisition with fat suppression. Third row: Zoomed view of the ROI that shows a subchondral osteophyte (highlighted by a white arrow in the ground truth reconstruction). This pathology is not visible in any of the accelerated reconstructions.

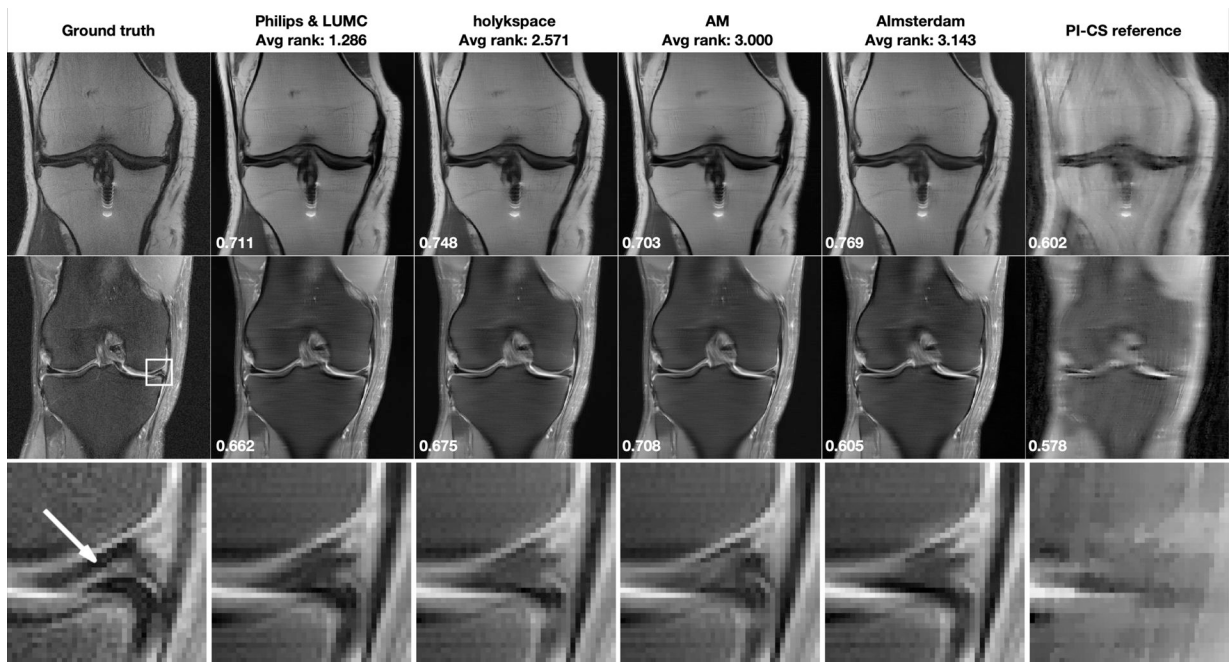


FIGURE 3.

Multi-Coil R=8 track results: Selected results from the top 4 submissions in each track. The submissions are ordered from left to right based on the average of radiologists' rankings. A combined parallel imaging and compressed sensing reconstruction using Total Generalized Variation (PI-CS) is shown for reference. SSIM to the ground truth for this particular slice is displayed in the bottom-left corner of each image. First row: Results for one slice from an acquisition without fat suppression. This case shows moderate artifact from a metal implant. Second row: One slice from an acquisition with fat suppression. This case shows a meniscal tear in the ROI indicated by a white rectangle in the ground truth image. Third row: Zoomed view of the ROI that shows a meniscal tear (highlighted by a white arrow in the ground truth reconstruction). This pathology is not well seen in any of the accelerated reconstructions.

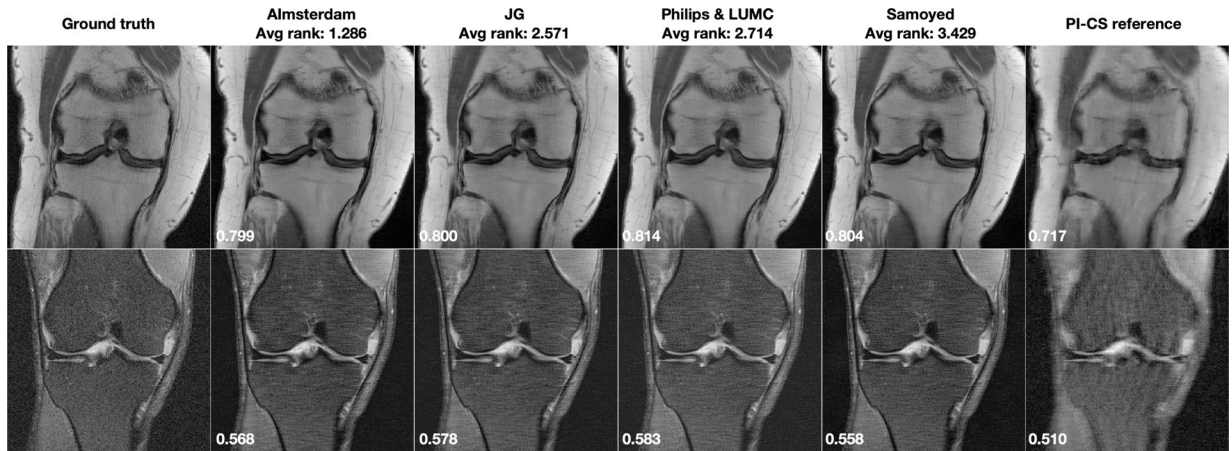


FIGURE 4.

Single-Coil R=4 track results: Selected results from the top 4 submissions for each track. SSIM to the ground truth for this particular slice is displayed in the bottom-left corner of each image. A combined parallel imaging and compressed sensing reconstruction using Total Generalized Variation (PI-CS) is shown for reference.

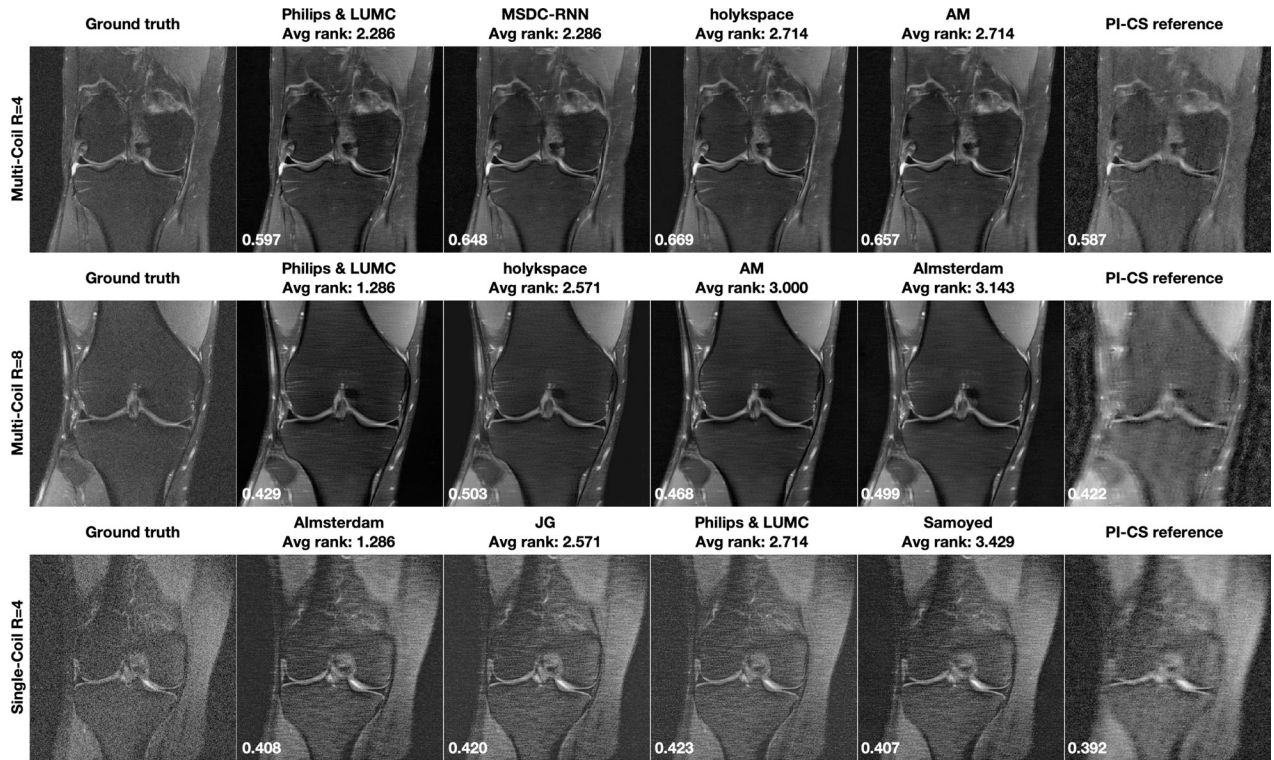


FIGURE 5.

Results from the case with the lowest averaged SSIM over the whole image volume from the top 4 submissions for each track. A combined parallel imaging and compressed sensing reconstruction using Total Generalized Variation (PI-CS) is shown for reference. Notably, all methods performed worst on the same case within each track, and all methods outperformed the PI-CS reference. SSIM to the ground truth for the shown slice is displayed in the bottom-left corner of each image.

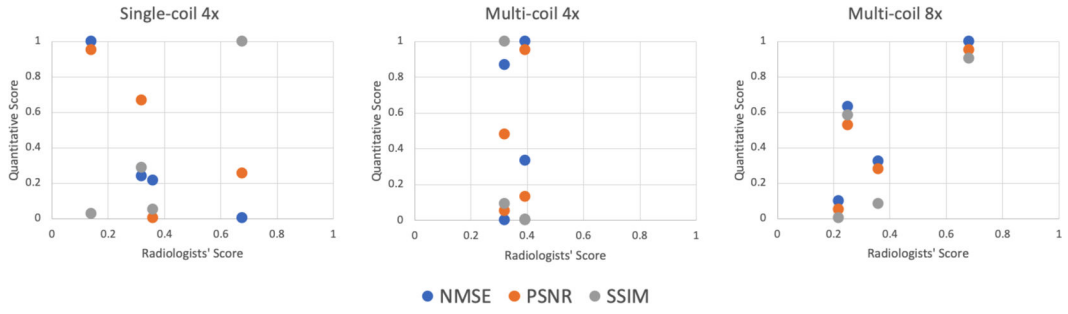


FIGURE 6. Scatterplots of the quantitative scores vs the average radiologists’ score based on ranking (1 is best) for the top 4 submissions in all three submission tracks. NMSE, PSNR, and SSIM scores are normalized so that the best score corresponds to a value of 1 for convenient visualization. For multi-coil R=8, the highest ranked submission was also the one that had the highest SSIM, NMSE and PSNR values. For single-coil R=4, only SSIM showed a similar trend as the radiologists scores, while the other two metrics showed almost opposite trends. For the multi-coil R=4 track, the top 4 submissions were very close together with all metrics. Therefore, the results are less conclusive.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

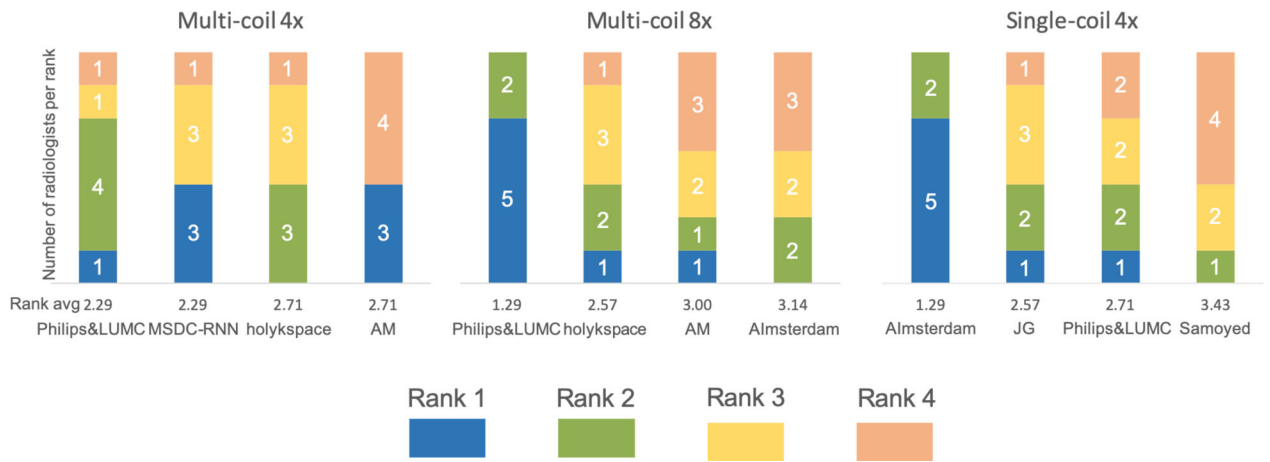


FIGURE 7.

Individual rankings by the 7 radiologists for the top 4 submissions in all three submission tracks. The ranks from 1 to 4 are color coded, and the number of radiologists who scored each submission at a certain rank is shown in the plot. For multi-coil R=8 and single-coil R=4 tracks, the radiologists had a strong preference for a single submission. The results are less consistent for the multi-coil R=4 track.

TABLE 1

Overview of the dataset that was provided for the fastMRI challenge.

	Cases		Slices	
	Multi-coil	Single-coil	Multi-coil	Single-coil
training	973	973	34,742	34,742
validation	199	199	7,135	7,135
test	118	108	4,092	3,903
challenge	104	92	3,810	3,305

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

TABLE 2

Average radiologist rankings and corresponding SSIM, NMSE and PSNR scores for the full challenge dataset for the top 4 submissions of the three tracks.

Multi-Coil R=4					
Team name	Rank	Avg radiologist rank	SSIM	NMSE	PSNR
Philips & LUMC	1(tie)	2.285	0.927	0.005	39.907
MSDC-RNN	1(tie)	2.285	0.927	0.005	39.740
holyspace	3(tie)	2.714	0.927	0.005	39.715
AM	3(tie)	2.714	0.928	0.005	39.807
Multi-Coil R=8					
Team name	Rank	Avg radiologist rank	SSIM	NMSE	PSNR
Philips & LUMC	1	1.286	0.901	0.0086	37.437
holyspace	2	2.571	0.899	0.0092	37.009
AM	3	3.000	0.901	0.0089	37.173
AMsterdam	4	3.143	0.898	0.0096	36.816
Single-Coil R=4					
Team name	Rank	Avg radiologist rank	SSIM	NMSE	PSNR
AMsterdam	1	1.286	0.754	0.031	32.549
JG	2	2.571	0.750	0.031	32.476
Philips & LUMC	3	2.714	0.751	0.030	32.666
Samoyed	4	3.428	0.751	0.029	32.761

TABLE 3

Average scores for the individual categories that the 7 radiologists were asked to rate, for the top 4 submissions from each track. Ratings followed a 4 point scale, where 1 is best and 4 is worst. In the Multi-Coil R=4 track, 5 out of 7 radiologists based their rankings strictly on their scores for this track. For the remaining 2 radiologists, the top-ranked submission also had the best overall score. In the Multi-Coil R=8 track, all radiologists based their rankings strictly on their scores. In the Single-Coil R=4 track, 6 out of 7 radiologists based their rankings strictly on their scores. For the remaining radiologist, the top-ranked submission also had the best overall score.

Multi-Coil R=4				
Team name	Artifacts	Sharpness	Contrast to noise	Diagnostic Confidence
Philips & LUMC	2.714	2.286	2.286	2.000
MSDC-RNN	2.571	2.286	2.429	1.857
holyspace	2.000	3.000	2.714	1.857
AM	2.000	3.000	2.000	2.000
Multi-Coil R=8				
Team name	Artifacts	Sharpness	Contrast to noise	Diagnostic Confidence
Philips & LUMC	1.714	2.286	2.286	2.286
holyspace	2.143	3.143	2.286	2.857
AM	1.857	3.286	2.429	3.143
Almsterdam	2.714	2.857	3.000	3.143
Single-Coil R=4				
Team name	Artifacts	Sharpness	Contrast to noise	Diagnostic Confidence
Almsterdam	2.429	2.286	2.143	2.286
JG	3.000	2.714	2.429	2.571
Philips & LUMC	2.714	3.000	2.714	2.857
Samoyed	3.143	3.286	3.000	3.286