

1 Deep learning for ultrasound beamforming

Ruud JG van Sloun, Jong Chul Ye and Yonina C Eldar

1.1 Introduction and relevance

Diagnostic imaging plays a critical role in healthcare, serving as a fundamental asset for timely diagnosis, disease staging and management as well as for treatment choice, planning, guidance, and follow-up. Among the diagnostic imaging options, ultrasound imaging (Szabo 2004) is uniquely positioned, being a highly cost-effective modality that offers the clinician an unmatched and invaluable level of interaction, enabled by its real-time nature. Its portability and cost-effectiveness permits point-of-care imaging at the bedside, in emergency settings, rural clinics, and developing countries. Ultrasonography is increasingly used across many medical specialties, spanning from obstetrics, cardiology and oncology to acute and intensive care, with a market share that is globally growing.

On the technological side, ultrasound probes are becoming increasingly compact and portable, with the market demand for low-cost ‘pocket-sized’ devices (i.e. “the stethoscope model”) expanding (Baran & Webster 2009). Transducers are miniaturized, allowing e.g. in-body imaging for interventional applications. At the same time, there is a strong trend towards 3D imaging (Provost, Papadacci, Arango, Imbault, Fink, Gennisson, Tanter & Pernot 2014) and the use of high-frame-rate imaging schemes (Tanter & Fink 2014); both accompanied by dramatically increasing data rates that pose a heavy burden on the probe-system communication and subsequent image reconstruction algorithms. Systems today offer a wealth of advanced applications and methods, including shear wave elasticity imaging (Bercoff, Tanter & Fink 2004), ultra-sensitive Doppler (Demené, Deffieux, Pernot, Osmanski, Biran, Gennisson, Sieu, Bergel, Franqui, Correas et al. 2015), and ultrasound localization microscopy for super-resolution microvascular imaging (Errico, Pierre, Pezet, Desailly, Lenkei, Couture & Tanter 2015), (Christensen-Jeffries, Couture, Dayton, Eldar, Hynynen, Kiessling, O’Reilly, Pinton, Schmitz, Tang et al. 2020).

With the demand for high-quality image reconstruction and signal extraction from less (e.g unfocused or parallel) transmissions that facilitate fast imaging, and a push towards compact probes, modern ultrasound imaging leans heavily on innovations in powerful digital receive channel processing. Beamforming, the process of mapping received ultrasound echoes to the spatial image domain, naturally lies at the heart of the ultrasound image formation chain. In this chapter, we discuss why and when deep learning methods can play a compelling role in the

digital beamforming pipeline, and then show how these data-driven systems can be leveraged for improved ultrasound image reconstruction (Van Sloun, Cohen & Eldar 2019).

This chapter is organized as follows. Sec. 1.2 briefly introduces various scanning modes in ultrasound. Then, in Sec. 1.3, we describe methods and rationale for digital receive beamforming. In Sec. 1.4 we elaborate on the opportunities of deep learning for ultrasound beamforming, and in Sec. 1.5 we review various deep network architectures. We then turn to typical approaches for training in Sec. 1.6. Finally, in Sec. 1.7 we discuss several future directions.

1.2 Ultrasound scanning in a nutshell

Ultrasound imaging is based on the pulse-echo principle. First, a radiofrequency (RF) pressure wave is transmitted into the medium of interest through a multi-element ultrasound transducer. These transducers are typically based on piezoelectric (preferably single-crystal) mechanisms or CMUT technology. After ionization, the acoustic wave backscatters due to inhomogeneities in the medium properties, such as density and speed of sound. The resulting reflections are recorded by the same transducer array and used to generate a so-called ‘brightness-mode’ (B-mode) image through a signal processing step termed beamforming. We will elaborate on this step in Sec. 1.3. The achievable resolution, contrast, and overall fidelity of B-mode imaging depends on the array aperture and geometry, element sensitivity and bandwidth. Transducer geometries include linear, curved or phased arrays. The latter is mainly used for extended field-of-views from limited acoustic windows (e.g. imaging of the heart between the ribs), enabling the use of angular beam steering due to a smaller pitch, namely, distance between elements. The elements effectively sample the aperture: using a pitch of half the wavelength (i.e. spatial Nyquist rate sampling) avoids grating lobes (spatial aliasing) in the array response. 2D ultrasound imaging is based on 1D arrays, while 3D imaging makes use of 2D matrix designs.

Given the transducer’s physical constraints, getting the most out of the system requires careful optimization across its entire imaging chain. At the front-end, this starts with the design of appropriate transmit schemes for wave field generation. At this stage, crucial trade-offs are made, in which the frame rate, imaging depth, and attainable axial and lateral resolution are weighted carefully against each other: improved resolution can be achieved through the use of higher pulse modulation frequencies and bandwidths; yet, these shorter wavelengths suffer from increased absorption and thus lead to reduced penetration depth. Likewise, high frame rate can be reached by exploiting parallel transmission schemes based on e.g. planar or diverging waves. However, use of such unfocused transmissions comes at the cost of loss in lateral resolution compared to line-based scanning with tightly focused beams. As such, optimal transmit schemes depend on the

application. We will briefly elaborate on three common transmit schemes for ultrasound B-mode imaging, an illustration of which is given in Fig. 1.1.

1.2.1 Focused transmits / line scanning

In line scanning, a series of E transmit events is used to, with each transmit e , produce a single depth-wise line in the image by focusing the transmitted acoustic energy along that line. Such focused transmits are typically achieved using a subaperture of transducer elements $c \in \{e - L, \dots, e + L\}$, excited with time-delayed RF pulses. By choosing these transmit delays per channel appropriately, the beam can be focused towards a given depth and (for phased arrays) angle. Focused line scanning is the most common transmit design in commercial ultrasound systems, enjoying improved lateral resolution and image contrast compared to unfocused transmits. Line-by-line acquisition is however time consuming (every lateral line requires a distinct transmit event), upper bounding the frame rate of this transmit mode by the number of lines, imaging depth, and speed of sound. This constraint can be relaxed via multi-line parallel transmit approaches, at the expense of reduced image quality.

1.2.2 Synthetic aperture

Synthetic aperture transmit schemes allow for synthetic dynamic transmit focusing by acquiring echoes with the full array following near-spherical wave excitation using individual transducer elements c . By performing E such transmits (typically $E = C$, the number of array elements), the image reconstruction algorithm (i.e. the beamformer) has full access to all transmit-receive pairs, enabling retrospective focusing in both transmit and receive. Sequential transmissions with all elements is however time consuming, and acoustic energy delivered into the medium is limited (preventing e.g. harmonic imaging applications). Synthetic aperture imaging also finds application in phased array intravascular ultrasound (IVUS) imaging, where one can only transmit and receive using one element/channel at a time due to catheter-based constraints.

1.2.3 Plane wave / ultrafast

Today, an increasing amount of ultrasound applications rely on high frame-rate (dubbed *ultrafast*) parallel imaging based on plane waves or diverging waves. Among these are e.g. ultrasound localization microscopy, highly-sensitive Doppler, shear wave elastography and (blood) speckle tracking. Where the former two mostly exploit the incredible vastness of data to obtain accurate signal statistics, the later two leverage high-speed imaging to track ultrasound-induced shear waves or tissue motion to estimate elasticity, strain or flow parameters. In plane wave imaging, planar waves are transmitted to insonify the full region of interest

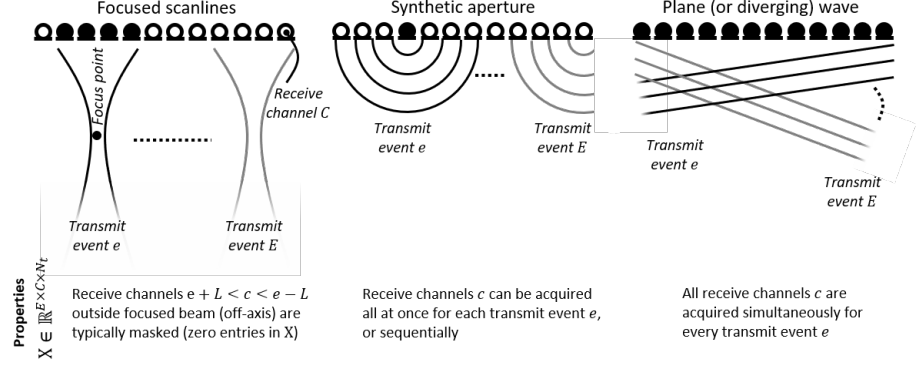


Figure 1.1 An illustration of three transmit types. Transmit events are denoted by e , receive channels as c and, for focused scanlines, $2L + 1$ is the size of the active aperture in terms of elements.

in a single transmit. Typically, several plane waves with different angles are compounded to improve image quality. For small-footprint phased arrays, diverging waves are used. These diverging waves are based on a set of virtual focus points that are placed behind the array, acting as virtual spherical sources. In that context, one can also interpret diverging wave imaging as a synthetic aperture technique.

With the expanding use of ultrafast transmit sequences in modern ultrasound imaging, a strong burden is placed on the subsequent receive channel processing. High data-rates not only raise substantial hardware complications related to data storage and data transfer, in addition, the corresponding unfocused transmissions require advanced receive beamforming to reach satisfactory image quality.

1.3 Digital ultrasound beamforming

We now describe how the received and digitized channel data is used to reconstruct an image in the digital domain, via a digital signal processing algorithm called beamforming.

1.3.1 Digital beamforming model and framework

Consider an array of C channels, and E transmit events, resulting in $E \times C$ measured and digitized RF signal vectors containing N_t samples. Denote $X \in \mathbb{R}^{E \times C \times N_t}$ as the resulting received RF data cube. Individual transmit events e can be tilted planar waves, diverging waves, focused scanlines through transmit beamforming, or any other desired pressure distribution in transmission. The goal of beamforming is to map this time-array-domain ‘channel’ data cube to

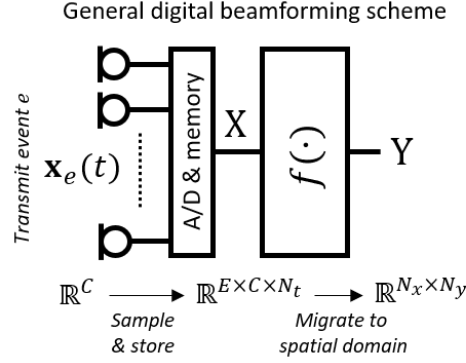


Figure 1.2 General beamforming scheme.

the spatial domain, through a processor $f(\cdot)$:

$$Y = f(X), \quad (1.1)$$

where $Y \in \mathbb{R}^{N_x \times N_y}$ denotes the beamformed spatial data, with N_x and N_y being the number of pixels in the axial and lateral direction, respectively. In principle, all beamforming architectures in ultrasound can be formulated according to (1.2) and the illustration in Fig. 1.2. The different approaches vary in their parameterization of $f(\cdot)$. Most are composed of a geometrical time-to-space migration of individual channels, and a subsequent combiner/processor. We will now go over some of the most common parameterizations for ultrasound beamforming. In Sec. 1.5, we will see that these conventional signal processing methods can also directly inspire parameterizations comprising deep neural networks.

1.3.2 Delay-and-sum

The industry standard beamforming algorithm is delay-and-sum beamforming (DAS). DAS is commonplace due to its low complexity, allowing for real-time image reconstruction at the expense of non-optimal image quality. Its processing can in general be written as (see Fig. 1.3 for an illustration):

$$Y = \sum_{E,L} W \odot D(X), \quad (1.2)$$

where $W \in \mathbb{R}^{E \times C \times N_x \times N_y}$ is an apodization weight tensor and \odot denotes the element-wise product. Note that apodization weights can be complex valued when $D(X)$ is IQ demodulated, allowing for phase shifting by W . In the remainder of this chapter, we will (without loss of generality) only consider real weights applied to RF data. Here, $D(\cdot)$ is a focussing function that migrates the RF time signal to space for each transmit event and channel, mapping X from $\mathbb{R}^{E \times C \times N_t}$ to $\mathbb{R}^{E \times C \times N_x \times N_y}$. This mapping is obtained by applying geometry-based time

delays to the RF signals with the aim of time-aligning the received echoes from the set of focal points.

DAS is typically employed in a so-called dynamic receive beamforming mode, in which the focal points change as a function of scan depth. In a specific variant of dynamic receive beamforming, *pixel-based beamforming*, each pixel is a focus point. Note that unlike its name suggests, dynamic does not mean that the beamformer is dynamically updating $D(\cdot)$ and W on the fly. Its name stems from the varying time-delays across fast time¹ (and therewith depth) to dynamically move the focal point deeper.

As said, channel delays are used to time-align the received echoes from a given position, and are determined by ray-based wave propagation, dictated by the array geometry, transmit design (plane wave, diverging wave, focused, SA), position of interest, and an estimate of the speed of sound. For each focal point $\{r_x, r_y\}$, channel c , and transmit event e , we can write the time-of-flight as:

$$\tau_{e,c,r_x,r_y} = \tau_{e,c,\mathbf{r}} = \frac{\|\mathbf{r}_e - \mathbf{r}\|_2 + \|\mathbf{r}_c - \mathbf{r}\|_2}{v}, \quad (1.3)$$

where $\tau_{e,c,\mathbf{r}}$ is the time-of-flight for an imaging point \mathbf{r} , \mathbf{r}_c is the position vector of the receiving element in the array, and v is the speed of sound in the medium. The vector \mathbf{r}_e depends on the transmit sequence: for focused transmits it is the position vector of the (sub)aperture center for transmit e ; for synthetic aperture it is the position vector of the e^{th} transmitting element; for diverging waves it is the position vector of the e^{th} virtual source located behind the array; for plane waves this is dictated by the transmit angle.

For any focus point $\{r_x, r_y\}$, the response at channel c for a given transmit event e is thus given by:

$$\mathbf{z}_{e,c}[r_x, r_y] = D_{e,c}(\mathbf{x}_{e,c}; \tau_{e,c,r_x,r_y}), \quad (1.4)$$

where $\mathbf{x}_{e,c}$ denotes the received signal for the e^{th} transmit event and c^{th} channel, and $D_{e,c}(\mathbf{x}; \tau)$ migrates $\mathbf{x}_{e,c}$ from time to space based on the geometry-derived delay τ . To achieve high-resolution delays in the discrete domain, (1.4) is typically implemented using interpolation or fractional delays with polyphase filters. Alternatively, delays can be implemented in the Fourier domain, which, as we shall discuss later, has practical advantages for e.g. compressed sensing applications (Chernyakova & Eldar 2014).

After migrating X to the spatial domain the apodization tensor W is applied, and the result is (coherently) summed across the e (transmit event) and c (channel) dimensions. Design of the apodization tensor W inherently poses a compromise between main lobe width and side lobe intensity, or equivalently, resolution and contrast/clutter. This can be intuitively understood from the far-field Fourier

¹ In ultrasound imaging a distinction is made between slow-time and fast-time: slow-time refers to a sequence of snapshots (i.e., across multiple transmit/receive events), at the pulse repetition rate, whereas fast-time refers to the time axis of the received RF signal for a given transmit event.

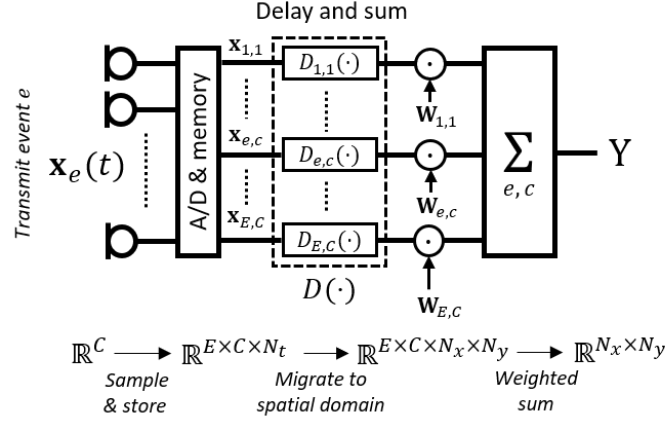


Figure 1.3 Delay-and-sum beamforming.

relationship between the beampattern and array aperture: analogous to filtering in frequency, the properties of a beamformer (spatial filter) are dictated by the sampling (aperture pitch), filter length (aperture size), and coefficients (apodization weights). Typical choices include Hamming-style apodizations, suppressing sidelobes at the expense of a wider main lobe. In commercial systems, the full (depth/position-dependent) weight tensor is carefully engineered and fine-tuned based on the transducer design (e.g. pitch, size, center frequency, near/far-field zones) and imaging application (e.g. cardiac, obstetrics, general imaging, or even intravascular).

1.3.3 Adaptive beamforming

Adaptive beamforming aims to overcome the inherent tradeoff between sidelobe levels and resolution of static DAS beamforming by making its apodization weight tensor \mathbf{W} fully data-adaptive, i.e., $\mathbf{W} \triangleq \mathbf{W}(\mathbf{X})$. Adaptation of \mathbf{W} is based on estimates of the array signal statistics, which are typically calculated instantaneously on a spatiotemporal block of data, or through recursive updates. Note that in general, the performance of these methods is bounded by the bias and variance of these statistics estimators. The latter can be reduced by using a larger block of samples (assuming some degree of spatio-temporal stationarity), which comes at the cost of reduced spatiotemporal adaptivity, or techniques such as sub-array averaging.

We will now briefly review some typical adaptive beamforming structures. In general, we distinguish between beamformers that act on individual channels and those that use the channel statistics to compute a single weighting factor across all channels, so called postfilters.

Minimum variance

A popular adaptive beamforming method is the minimum variance distortionless response (MVDR), or Capon, beamformer. The MVDR beamformer acts on individual channels, with the optimal weights $\mathbf{W} \in \mathbb{R}^{E \times C \times N_x \times N_y}$ defined as those that, for each transmit event e and location $[r_x, r_y]$, minimize the total signal variance/power while maintaining distortionless response in the direction/focal-point of interest. This amounts to solving:

$$\begin{aligned} \hat{\mathbf{w}}_{mv,e}[r_x, r_y] = \arg \min_{\mathbf{w}} \mathbf{w}^H \mathbf{R}_{x_e[r_x, r_y]} \mathbf{w} \\ \text{s.t. } \mathbf{w}^H \mathbf{1} = 1, \end{aligned} \quad (1.5)$$

where $\hat{\mathbf{w}}_{mv,e}[r_x, r_y] \in \mathbb{R}^C$ is the weight vector for a given transmit event and location, and $\mathbf{R}_{x_e[r_x, r_y]}$ denotes the estimated channel covariance matrix for transmit event e and location $[r_x, r_y]$. Solving (1.5) requires inversion of the covariance matrix, which grows cubically with the number of array channels. This makes MV beamforming computationally much more demanding than DAS, in particular for large arrays. This in practice results in a significantly longer reconstruction time and thereby deprives ultrasound of the interactability that makes it so appealing compared to e.g. MRI and CT. To boost image quality, eigen-space based MV beamforming (Asl & Mahloojifar 2010) performs an eigendecomposition of the covariance matrix and subsequent signal subspace selection before inversion. This further increases computational complexity. While significant progress has been made to decrease the computational time of MV beamforming algorithms (Kim, Park, Kim, Park & Bae 2014)(Bae, Park & Kwon 2016), real-time implementation remains a major challenge. In addition, it relies on accurate estimates of the signal statistics, which (as mentioned above) requires some form of spatiotemporal averaging.

Coherence factor (CF)

Coherence Factor (CF) weighing (Mallart & Fink n.d.) also applies content-adaptive apodization weights \mathbf{W} , however with a specific structure: the weights across different channels are identical/tied. CF weighing thus in practice acts as a post-filter after DAS beamforming. The pixel-wise weighing in this post-filter is based on a ‘‘coherence factor’’: the ratio between the coherent and incoherent energy across the array channels. CF weighting however suffers from artifacts when the SNR is low and estimation of coherent energy is challenging (Nilsen & Holm 2010a). This is particularly problematic for unfocused techniques such as PW and SA imaging.

Wiener

The Wiener beamformer produces a minimum mean-squared-error (MMSE) estimate of the signal amplitude A stemming from a particular direction/location

of interest:

$$\hat{\mathbf{w}}_e[r_x, r_y] = \arg \min_{\mathbf{w}} E(|A - \mathbf{w}^H \mathbf{x}_e[r_x, r_y]|^2). \quad (1.6)$$

The solution is (Nilsen & Holm 2010b):

$$\hat{\mathbf{w}}_e[r_x, r_y] = \frac{|A|^2}{|A|^2 + (\mathbf{w}_{mv,e}[r_x, r_y])^H \mathbf{R}_{n_e[r_x, r_y]} \mathbf{w}_{mv,e}[r_x, r_y]} \mathbf{w}_{mv,e}[r_x, r_y], \quad (1.7)$$

where $\mathbf{R}_{n_e[r_x, r_y]}$ is the noise covariance matrix. Wiener beamforming is thus equivalent to the MVDR beamformer followed by a (CF-like) post-filter that scales the output as a function of the remaining noise power after MVDR beamforming (second term in the denominator). Note that the Wiener beamformer requires estimates of both the signal power and noise covariance matrix. The latter can e.g. be estimated by assuming i.i.d. white noise, i.e., $\mathbf{R}_{n_e[r_x, r_y]} = \sigma_n^2 \mathbf{I}$, and calculating σ_n from the mean squared difference between the MVDR beamformed output and the channel signals.

Iterative Maximum-a-Posteriori (iMAP)

Chernyakova *et al.* propose an iterative maximum-a-posteriori (iMAP) estimator (Chernyakova, Cohen, Shoham & Eldar 2019), which formalizes post-filter-based methods as a MAP problem by incorporating a statistical prior for the signal of interest. Assuming a zero-mean Gaussian random variable of interest with variance σ_a^2 that is uncorrelated to the noise (also Gaussian, with variance σ_n^2), the beamformer can be derived as:

$$\hat{\mathbf{w}}_e[r_x, r_y] = \frac{\sigma_a[r_x, r_y]^2}{M\sigma_a[r_x, r_y]^2 + M\sigma_n[r_x, r_y]^2} \mathbf{1}. \quad (1.8)$$

The signal and noise variances are estimated in an iterative fashion: first a beamformed output is produced according to (1.8), and then the noise variance is estimated based on the mean-squared difference with the individual channels. This process is repeated until a stopping criterion is met. Note that while iMAP performs iterative estimation of the statistics, the beamformer itself shares strong similarity with Wiener postfiltering and CF beamforming.

1.4 Deep learning opportunities

In this section we will elaborate on some of the fundamental challenges of ultrasound beamforming, and the role that deep learning solutions can play in overcoming these challenges (Van Sloun *et al.* 2019).

1.4.1 Opportunity 1: Improving image quality

As we saw in the previous section, classic adaptive beamformers are derived based on specific modeling assumptions and knowledge of the signal statistics.

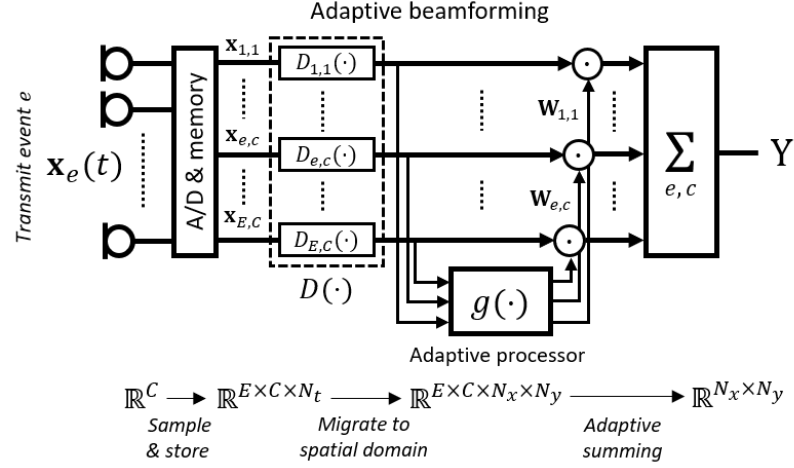


Figure 1.4 Adaptive beamforming

This limits the performance of model-based adaptive beamformers, which are bounded by:

1. Accuracy and precision of the estimated signal statistics using data sampled from a strongly non-stationary ultrasound RF process. In practice, only limited samples are available, making accurate estimation challenging.
2. Adequacy of the simple linear acquisition model (including homogeneous speed of sound) and assumptions on the statistical structure of desired signals and noise (uncorrelated Gaussian). In practice speed of sound is heterogeneous and noise statistics are highly complex and correlated with signal, e.g. via multiple scattering that leads to reverberation and haze in the image.

In addition, specifically for MVDR beamforming, the complexity of the required matrix inversions hinders real-time implementation. Deep learning can play an important role in addressing these issues, enabling:

1. Drastic acceleration of slow model-based approaches such as MVDR, using accelerated neural network implementations as function approximators.
2. High-performant beamforming outputs without explicit estimation of signal statistics by exploiting useful priors learned from previous examples (training data).
3. Data-driven nonlinear beamforming architectures that are not bounded by modelling assumptions on linearity and noise statistics. While analytically formalizing the underlying statistical models of ultrasound imaging is highly challenging, and its optimization likely intractable, deep learning circumvents this by learning powerful models directly from data.

1.4.2 Opportunity 2: Enabling fast and robust compressed sensing

Performing beamforming in the digital domain requires sampling the signals received at the transducer elements and transmitting the samples to a back-end processing unit. To achieve sufficient delay resolution for focusing, hundreds of channel signals are typically sampled at 4-10 times their bandwidth, i.e., the sampling rate may severely exceed the Nyquist rate. This problem becomes even more pressing for 3D ultrasound imaging based on matrix transducer technology, where directly streaming channel data from thousands of sensors would lead to data rates of thousands of gigabits per second. Today's technology thus relies on microbeamforming or time-multiplexing to keep data rates manageable. The former compresses data from multiple (adjacent) transducer elements into a single line, thereby virtually reducing the number of receive channels and limiting the attainable resolution and image quality. The latter only communicates a subset of the channel signals to the backend of the system for every transmit event, yielding reduced frame rates.

To overcome these data-rate challenges without compromising image quality and frame rates, a significant research effort has been focused on compressed sensing for ultrasound imaging. Compressed sensing permits low-data-rate sensing (below the Nyquist rate) with strong signal recovery guarantees under specific conditions (Eldar & Kutyniok 2012, Eldar 2015). In general, one can perform compressed sensing along three axes in ultrasound: 1) fast-time, 2) slow-time, and 3) channels/array elements. We denote the undersampled measured RF data cube as $X_u \in \mathbb{R}^{E_u \times C_u \times N_u}$, with $E_u C_u N_u < ECN$.

Sub-Nyquist fast-time sampling

To perform sampling rate reduction across fast-time, one can consider the received signals within the framework of finite rate of innovation (FRI) (Eldar 2015, Gedalyahu, Tur & Eldar 2011). Tur *et al.* (Tur, Eldar & Friedman 2011) modeled the received signal at each element as a finite sum of replicas of the transmitted pulse backscattered from reflectors. The replicas are fully described by their unknown amplitudes and delays, which can be recovered from the signals' Fourier series coefficients. The latter can be computed from low-rate samples of the signal using compressed sensing (CS) techniques (Eldar 2015, Eldar & Kutyniok 2012). In (Wagner, Eldar, Feuer, Danin & Friedman 2011, Wagner, Eldar & Friedman 2012), the authors extended this approach and introduced compressed ultrasound beamforming. It was shown that the beamformed signal follows an FRI model and thus it can be reconstructed from a linear combination of the Fourier coefficients of the received signals. Moreover, these coefficients can be obtained from low-rate samples of the received signals taken according to the Xampling framework (Mishali, Eldar & Elron 2011, Mishali, Eldar, Dounaevsky & Shoshan 2011, Michaeli & Eldar 2012). Chernyakova *et al.* showed this Fourier domain relationship between the beam and the received signals holds irrespective of the FRI model. This leads to a general concept of

frequency domain beamforming (FDBF) (Chernyakova & Eldar 2014) which is equivalent to beamforming in time. FDBF allows to sample the received signals at their effective Nyquist rate without assuming a structured model, thus, it avoids the oversampling dictated by digital implementation of beamforming in time. When assuming that the beam follows a FRI model, the received signals can be sampled at sub-Nyquist rates, leading to up to 28 fold reduction in sampling rate, i.e. $N_t = 28N_u$ (Chernyakova, Cohen, Mulayoff, Sde-Chen, Fraschini, Bercoff & Eldar 2018, Burshtein, Birk, Chernyakova, Eilam, Kempinski & Eldar 2016, Lahav, Chernyakova & Eldar 2017).

Channel and transmit-event compression

Significant research effort has also been invested in exploration of sparse array designs ($C_u < C$) (Liu & Vaidyanathan 2017, Cohen & Eldar 2018a) and efficient sparse sampling of transmit events ($E_u < E$). It has been shown that with proper sparse array selection and a process called convolutional beamforming, the beampattern can be preserved using far fewer elements than the standard uniform linear array (Cohen & Eldar 2018b). Typical designs include sparse periodic arrays (Austeng & Holm 2002) or fractal arrays (Cohen & Eldar 2020). In (Besson, Carrillo, Perdios, Arditi, Bernard, Wiaux & Thiran 2016) a randomly sub-sampled set of receive transducer elements is used, and the work in (Lorintiu, Liebgott, Alessandrini, Bernard & Friboulet 2015) proposes learned dictionaries for improved CS-based reconstruction from sub-sampled RF lines. In (Huijben, Veeling, Janse, Mischel & van Sloun 2020), the authors use deep learning to optimize channel selection and slow-time-sampling for B-mode and downstream color Doppler processing, respectively. Reduction in both time sampling and spatial sampling can be achieved by Compressed Fourier-Domain Convolutional Beamforming leading to reductions in data of two orders of magnitude (Mamistvalov & Eldar 2020).

Beamforming and image recovery after compression

After compressive acquisition, dedicated signal recovery algorithms are used to perform image reconstruction from the undersampled dataset. Before deep learning became popular, these algorithms relied on priors/regularizers (e.g. sparsity in some domain) to solve the typically ill-posed optimization problem in a model-based (iterative) fashion. They assume knowledge about the measurement PSF and other system parameters. However, as mentioned before, the performance of model-based algorithms is bounded by the accuracy of the modelling assumptions, including the acquisition model and statistical priors. In addition, iterative solvers are time-consuming, hampering real-time implementation. Today, deep learning is increasingly used to overcome these challenges (Perdios, Besson, Arditi & Thiran 2017, Kulkarni, Lohit, Turaga, Kerviche & Ashok 2016): 1) Deep learning can be used to learn complex statistical models (explicitly or implicitly)

directly from training data, 2) Neural-networks can serve as powerful function approximators that accelerate iterative model-based implementations.

1.4.3 Opportunity 3: Beyond MMSE with task-adaptive beamforming

Classically, beamforming is posed as a signal recovery problem under spatially white Gaussian noise. In that context, optimal beamforming is defined as the beamformer that best estimates the signal of interest in a minimum mean squared error (MMSE) sense. However, beamforming is rarely the last step in the processing chain of ultrasound systems. It is typically followed by demodulation (envelope detection), further image enhancement, spatiotemporal processing (e.g. motion estimation), and image analysis. In that regard it may be more meaningful to define optimality of the beamformer with respect to its downstream task. We refer this as task-adaptive beamforming. We can define several such tasks in ultrasound imaging. First, we have tasks that focus on further enhancement of the images after beamforming. This includes downstream processing for e.g. de-speckling, de-convolution, or super-resolution, which all have different needs and requirements from the beamformer output. For instance, the performance of deconvolution or super-resolution algorithms is for a large part determined by the invertibility of the point-spread-function model, shaped by the beamformer. Beyond image enhancement, one can think of motion estimation tasks (requiring temporal consistency of the speckle patterns with a clear spatial signature that can be tracked), or even further downstream applications such as segmentation or computer aided diagnosis (CAD).

One may wonder how to optimize beamforming for such tasks in practice. Fortunately, many of these downstream processing tasks are today performed using convolutional neural networks or derivatives therefrom. If the downstream processor is indeed a neural architecture or another algorithm through which one can easily backpropagate gradients, one can directly optimize a neural-network beamformer with respect to its downstream objective/loss through deep learning methods. In this case, *deep* not only refers to the layers in individual networks, but also the stack of beamforming and downstream neural networks. If backpropagation is non-trivial, one could resort to Monte-Carlo gradient estimators based on e.g. the REINFORCE estimator (Williams 1992) or more generally through reinforcement learning algorithms (Sutton & Barto 2018).

1.4.4 A brief overview of the state-of-the-art

The above opportunities have spurred an ever growing collection of papers from the research community. In Table 1.1, we provide an overview of a selection of these in the context of these opportunities, and the challenges they aim to address. Many of these papers simultaneously address more than one challenge/opportunity.

References	Opportunity and focus		
	Real-time high image quality (Main goals)	Compressed sensing (Subsampling axis)	Task-adaptive (Considered tasks)
(Luchies & Byram 2018)	Reduce off-axis scattering	-	-
(Yoon, Khan, Huh & Ye 2018)	-	Transmit and channels	-
(Khan, Huh & Ye 2020a)	-	Channels	-
(Khan, Huh & Ye 2020b)	Boost resolution or suppress speckle	-	Deconvolution and speckle reduction
(Luijten, Cohen, de Bruijn, Schmeitz, Mischi, Eldar & van Sloun 2019a)	Boost contrast and resolution	Transmit and Channels	-
(Nair, Washington, Tran, Reiter & Bell 2020)	-	-	Segmentation
(Wiacek, González & Bell 2020)	Accelerate coherence imaging	-	-
(Kessler & Eldar 2020) and (Mamistvalov & Eldar 2021)	-	Channels and fast-time Fourier coefficients	-
(Huijben et al. 2020)	-	Channels and slow-time	Doppler and B-Mode
(Hyun, Brickson, Looby & Dahl 2019)	Suppress speckle	-	Speckle reduction

Table 1.1 Overview of some of the current literature and their main focus in terms of the opportunities defined in section 1.4.

1.4.5 Public datasets and open source code

To support the development of deep-learning-based solutions in the context of these opportunities, a challenge on ultrasound beamforming by deep learning (CUBDL) was organized. For public raw ultrasound channel datasets as well as open source code we refer the reader to the challenge website (Bell, Huang, Hyung, Eldar, van Sloun & Mischi 2020) and the paper describing the datasets, methods and tools (Hyun, Wiacek, Goudarzi, Rothlübbers, Asif, Eickel, Eldar, Huang, Mischi, Rivaz, Sinden, van Sloun, Strohm & Bell 2021).

1.5 Deep learning architectures for ultrasound beamforming

1.5.1 Overview and common architectural choices

Having set the scope and defined opportunities, we now turn to some of the most common implementations of deep learning in ultrasound beamforming architectures. As in computer vision, most neural architectures for ultrasound beamforming are based on 2D convolutional building blocks. With that, they thus rely on translational equivariance/spatial symmetry of the input data. It is worth noting that this symmetry only holds to some extent for ultrasound imaging, as its point spread function in fact changes as a function of location. When operating on raw channel data before time-space migration (TOF correction), these effects become even more pronounced. Most architectures also restrict the receptive field of the neural network, i.e. beamforming outputs for given spatial location (line/pixel) are computed based on a selected subset of X . This is either implicit, through the depth and size of the selected convolutional kernels, or explicit, by for example only providing a selected number of depth slices (Khan et al. 2020a) as an input to the network.

In the following we will discuss a number of architectures. We explicitly specify their input, architectural design choices, and output to clarify what part of the beamforming chain they are acting on and how. We point out that the below is not exhaustive, but with examples selected to illustrate and cover the spectrum of approaches and design choices from an educational perspective.

1.5.2 DNN directly on channel data

In (Nair et al. 2020), deep learning is used directly on the raw channel data. The deep neural network thus has to learn both the classically geometry-based time-to-space migration (TOF correction) as well as the subsequent beamsumming of channels to yield a beamformed image. In particular the former makes this task particularly challenging - the processor is not a typical image-to-image mapping but rather a time-to-space migration. In addition to yielding a beamformed output, the network in parallel also provides a segmentation mask, which is subsequently used to enhance the final image by masking regions in the image that are classified as anechoic.

DNN directly on channel data by Nair et al.

Acquisition type: Single plane wave imaging.

Input: Complex IQ demodulated data cube $X_{in} \in \mathbb{C}^{1 \times C \times N_t}$, reformatted to real inputs $X_{in} \in \mathbb{R}^{C \times N_t \times 2}$ before being fed to the network.

Architecture: U-net variant consisting of a single VGG encoder and two parallel decoders with skip connections that map to B-mode and segmentation outputs respectively. The encoder has 10 convolutional layers (kernel size = 3×3) with batch normalization and downsamples the spatial domain via 2×2 max pooling layers while simultaneously increasing the number of feature channels. The two decoders both comprise 9 convolutional layers and perform spatial upsampling to map the feature space back to the desired spatial domain.

Output: RF B-mode data and pixel-wise class probabilities (segmentation) for the full image $Y_{bf} \in \mathbb{R}^{N_x \times N_y}$, and $Y_{seg} \in \mathbb{R}^{N_x \times N_y \times 1}$, respectively.

1.5.3 DNN for beam-summing

Hybrid architectures: geometry and learning

While the work by (Nair et al. 2020) replaces the entire beamformer by a deep neural network, most of today’s embodiments of ultrasound beamforming with deep learning work on post-delayed channel data. That is, the migration from time-to-space is a deterministic pre-processing step that relies on geometry-based TOF correction. This holds for all of the specific architectures that we will cover in the following. In that sense, they are all hybrid model-based/data-driven beamforming architectures.

Learning improved beam-summing

We will now discuss designs that replace only the beam-summing stage by a deep network, i.e. after TOF correction, as illustrated in Fig. 1.5. In (Yoon et al. 2018, Khan, Huh & Ye 2019, Khan et al. 2020a, Khan et al. 2020b, Vignon, Shin, Meral, Apostolakis, Huang & Robert 2020), the authors apply deep convolutional neural networks to perform improved channel aggregation/beam-summing after TOF correction.

DNN beamsumming by Khan et al.

Acquisition type: Line scanning, so the first dimension (transmit events / lines E) has been mapped to the lateral dimension N_y during the time-space migration/TOF correction.

Input: For each depth/axial location, a TOF-corrected RF data cube $Z_{in} \in \mathbb{R}^{C \times 3 \times N_y}$. The input data cube comprises a stack of 3 axial slices centered around the axial location of interest.

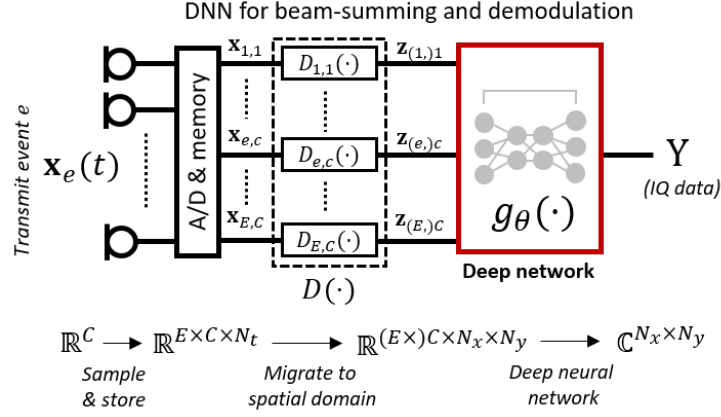


Figure 1.5 DNN replacing the beam-summing processor.

Architecture: 37 convolutional layers (kernel size = 3×3), of which all but the last have batch normalization and ReLU activations.

Output: IQ data for each axial location $\mathbf{y}_{n_x} \in \mathbb{R}^{2 \times 1 \times N_y}$, where the first dimension contains the in-phase (I) and quadrature (Q) components, such that we can also define $\mathbf{y}_{n_x} \in \mathbb{C}^{N_x \times 1}$.

(Kessler & Eldar 2020) use a similar strategy, albeit that pre-processing TOF correction is performed in the Fourier domain. This allows efficient processing when sensing only a small set of the Fourier coefficients of the received channel data. The authors then use a deep convolutional network to perform beamsumming, mapping the TOF-corrected channel data into a single beamformed RF image without aliasing artefacts:

DNN beamsumming by Kessler et al.

Acquisition type: Phased array line scanning, so the first dimension (transmit events / angular lines E) has been mapped to the lateral dimension N_y during the time-space migration/TOF correction. Reconstruction is in the polar domain, i.e. N_x refers to radial position, and N_y to angular position.

Input: TOF-corrected RF data cube $\mathbf{Z}_{in} \in \mathbb{R}^{C \times N_x \times 3}$. The input data cube comprises data corresponding to 3 angles centered around the angular position of interest.

Architecture: U-net variant with 3 contracting blocks and 3 expanding blocks of convolutional layers with parametric ReLU (PReLU) activations.

Output: RF data for each angle of interest $\mathbf{y}_{n_y} \in \mathbb{R}^{N_x \times 1}$.

1.5.4 DNN as an adaptive processor

The architecture we will discuss now is inspired by the MV beamforming architecture. Instead of replacing the beamforming process entirely, the authors in (Luijten, Cohen, de Bruijn, Schmeitz, Mischi, Eldar & van Sloun 2019b, Luijten et al. 2019a) propose to use a deep network as an artificial agent that calculates the optimal apodization weights \mathbf{W} on the fly, given the received pre-delayed channel signals at the array $D(X)$. See Fig. 1.6 for an illustration. By only replacing this bottleneck component in the MVDR beamformer, and constraining the problem further by promoting close-to-distortionless response during training (i.e. $\sum_c w_c \approx 1$), this solution is highly data-efficient, interpretable, and has the ability to learn powerful models from only few images (Luijten et al. 2019a).

DNN as adaptive processor by Luijten et al.

Acquisition types: Single plane wave imaging and synthetic aperture (intravascular ultrasound). For the latter, a virtual aperture is constructed by combining the received signals of multiple transmits and receives.

Input: TOF-corrected RF data cube $Z \in \mathbb{R}^{1 \times C \times N_x \times N_y}$.

Architecture: Four convolutional layers comprising 128 nodes for the input and output layers, and 32 nodes for the hidden layers. The kernel size of the filters is 1×1 , making the receptive field of the network a single pixel. In practice, this is thus a per-pixel fully-connected layer across the array channels. The activation functions are antirectifiers (Chollet n.d.), which, unlike ReLUs, preserve both the positive and negative signal components at the expense of a dimensionality increase.

Output: Array apodization tensor $\mathbf{W} \in \mathbb{R}^{1 \times C \times N_x \times N_y}$, which is subsequently multiplied (element-wise) with the network inputs Z to yield a beamformed output Y .

Complexity, inference speed, and stability

Since pixels are processed independently by the network, a large amount of training data is available per acquisition. Inference is fast and real-time rates are achievable on a GPU-accelerated system. For an array of 128 elements, adaptive calculation of a set of apodization weights through MV beamforming requires $> N^3 (= 2,097,152)$ floating point operations (FLOPS), while the deep-learning architecture only requires 74656 FLOPS (Luijten et al. 2019a), in practice leading to a more than $400\times$ speed-up in reconstruction time. Compared to MV beamforming, the deep network is qualitatively more robust, with less observed artefactual reconstructions that stem from e.g. instable computations of the inverse autocorrelation estimates in MV.

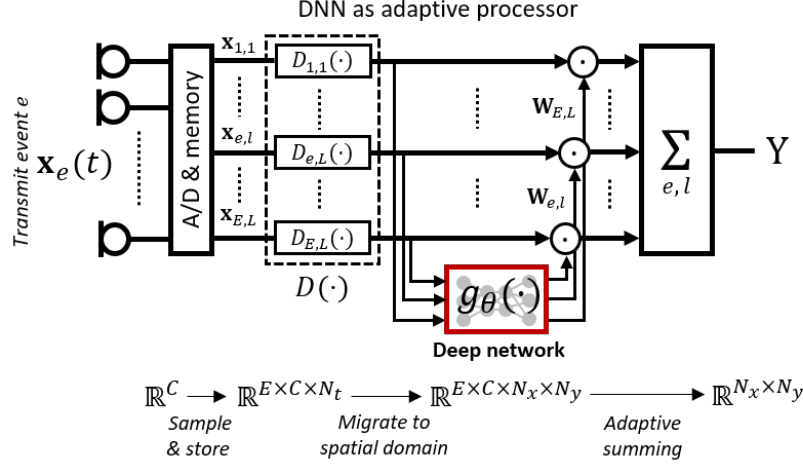


Figure 1.6 DNN replacing the adaptive processor of classical adaptive beamforming methods.

1.5.5 DNN for Fourier-domain beamsummation

Several model-based beamforming methods process ultrasound channel data in the frequency domain (Holfort, Gran & Jensen 2009), (Byram, Dei, Tierney & Dumont 2015). In this spirit, the authors of (Luchies & Byram 2018, Luchies & Byram 2019, Luchies & Byram 2020) use deep networks to perform wideband beamsummation by processing individual DFT bins of the TOF-corrected and axially-windowed RF signals. Each DFT bin is processed using a distinct neural network. After processing in the Fourier domain, the channel signals are summed for each window and a beamformed RF scanline is reconstructed using an inverse short-time Fourier transform (see Fig. 1.7).

DNN for Fourier-domain beamsummation by Luchies et al.

Acquisition type: Line scanning, where each transmit event produces one lateral scanline. The first dimension of $\mathbf{X} \in \mathbb{R}^{E \times C \times N_t}$ is thus directly mapped to the lateral dimension N_y in the TOF-correction step: $\mathbf{Z} = D(\mathbf{X}) \in \mathbb{R}^{C \times N_x \times N_y}$.

Input: Fourier transform of an axially-windowed (window length S) and TOF-corrected RF data cube for a single scanline, i.e. $\tilde{Z}_{n_x, n_y} = \mathcal{F}(Z_{n_x, n_y}) \in \mathbb{C}^{C \times S \times 1}$, with $Z_{n_x, n_y} \in \mathbb{R}^{C \times L_x \times 1}$. Before feeding to the network, \tilde{Z}_{n_x, n_y} is converted to real values by stacking the real and imaginary components, yielding $\hat{Z}_{n_x, n_y} \in \mathbb{R}^{2C \times S \times 1}$.

Architecture: S identical fully-connected neural networks, one for each DFT bin. Each neural network of this stack thus takes the $2C$ channel values corresponding to that bin as its input. The S networks all have

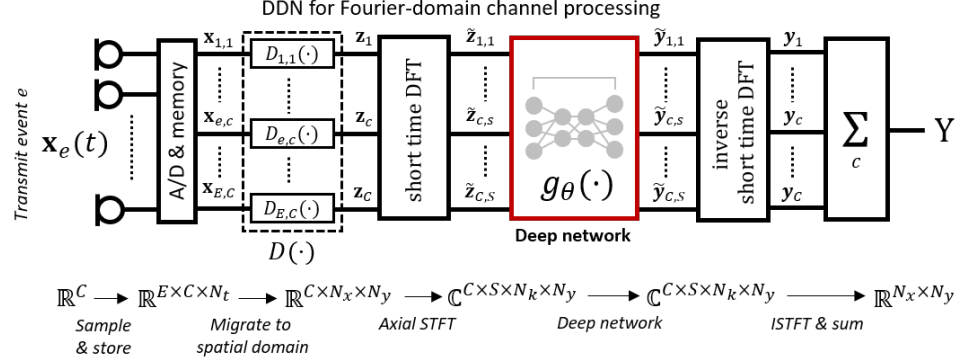


Figure 1.7 DNN for Fourier-domain beams summing

5 fully-connected layers with the hidden layers having 170 neurons and ReLU activations. Each network then returns $2C$ channel values, the real and imaginary components of that frequency bin after processing.

Output: Processed Fourier components of the axially-windowed and TOF-corrected RF data cube for a single scanline: $\tilde{Y}_{n_x, n_y} \in \mathbb{C}^{C \times S \times 1}$. To obtain a beamformed image, the C channels are summed and an inverse short-time Fourier transform is used to compound the responses of all axial windows.

1.5.6 Post-filtering after beams summing

We will now discuss some post-filter approaches, i.e. methods applied after channel beams summing, but before envelope detection and brightness compression. Several post-filtering methods have been proposed for compounding beam-summed RF outputs from multiple transmit events e with some spatial overlap (e.g. multiple plane/diverging waves). Traditionally, multiple transmit events are compounded by coherent summing (i.e. after transmit delay compensation). Today, deep learning is increasingly used to replace the coherent summing step. In (Lu, Millioz, Garcia, Salles, Liu & Friboulet 2020), the authors perform neural network compounding from a small number of transmits, and train towards the image obtained by coherently summing a much larger number of transmit events. In (Chennakeshava, Luijten, Drori, Mischi, Eldar & van Sloun 2020), the authors pose compounding as an inverse problem, which they subsequently solve using a model-based deep network inspired by proximal gradient methods. Post-filtering has also been used to e.g. remove aliasing artifacts due to sub-Nyquist sampling on beamformed 1D RF lines (Mamistvalov & Eldar 2021). The most common method to perform sparse recovery and solve the L1 minimization problem is using compressed sensing algorithms, such as ISTA and NESTA (Eldar 2015). However, they typically suffer from high computational load, and do not always

ensure high quality recovery. Here we discuss the process of unfolding an iterative algorithm as the layers of a deep network for sparse recovery. The authors in (Mamistvalov & Eldar 2021), built an efficient network by unfolding the ISTA algorithm for sparse recovery, based on the previously suggested LISTA (Gregor & LeCun 2010). In their technique, they recover both spatially and temporally sub-Nyquist sampled US data, after delaying it in the frequency domain, using a simple, computationally efficient, and interpretable deep network.

DNN as a recovery method by Mamistvalov et al.

Acquisition type: Line scanning, so the first dimension (transmit events / lines E) has been mapped to the lateral dimension N_y during the time-space migration/TOF correction.

Input: Frequency domain delayed and summed data (or frequency domain convolutionally beamformed data), after appropriate inverse Fourier transform and appropriate zero padding, to maintain the desired temporal resolution. The input data is a vector, $Z_{in} \in \mathbb{R}^{N_{st} \times 1}$, where N_{st} is the traditionally used number of samples for beamforming. The recovery is done for each image line separately.

Architecture: Simple architecture of unfolded ISTA algorithm, consisting of 30 layers, each includes two convolutional layers that mimic the matrix multiplications of ISTA and one soft thresholding layer. One last convolutional layer is added to recover the actual beamformed signal from the recovered sparse code.

Output: Beamformed signal for each image line without artifacts caused by sub-Nyquist sampling, $Z_{out} \in \mathbb{R}^{N_{st} \times 1}$.

1.6 Training strategies and data

1.6.1 Training data

The model parameters of the above beamforming networks are optimized using training data that consists of simulations, in-vitro/in-vivo data, or a combination thereof. We will now discuss some of the strategies for selecting training data, and in particular generating useful training targets.

Simulations

Training ultrasound beamformers using simulated data is appealing, since various ultrasound simulation toolboxes, such as Field II (Jensen 2004), k-wave (Treeby & Cox 2010), and the Matlab Ultrasound Toolbox (MUST)², allow for flexible generation of input-target training data. Simulations can be used to generate pairs of RF data, each pair comprising a realistic imaging mode based on the

² <https://www.biomecardio.com/MUST/>, by Damien Garcia

actual hardware and probe (input), and a second mode of the same scene with various more desirable properties (target). One can get creative with the latter, and we here list a number of popular approaches found in the literature:

1. Target imaging mode with unrealistic, yet desired, array and hardware configuration:
 - a. Higher frequencies/shorter wavelengths (without the increased absorption) to improve target resolution (Chennakeshava et al. 2020).
 - b. Larger array aperture to improve target resolution (Vignon et al. 2020).
2. Removal of undesired imaging effects such as off-axis scattering from target RF data (Luchies & Byram 2019).
3. Targets constructed directly from simulation scene/object:
 - a. Point targets on a high-resolution simulation grid (Youn, Ommen, Stuart, Thomsen, Larsen & Jensen 2020).
 - b. Masks of medium properties, e.g. anechoic region segmentations (Nair et al. 2020).

When relying solely on simulations, one has to be careful to avoid catastrophic domain shift when deploying the neural models on real data. Increasing the realism of the simulations, mixing simulations with real data, using domain-adaptation methods, or limiting the neural networks receptive field can help combat domain-shift issues.

Real data

Training targets from real acquisitions are typically based on high-quality (yet computationally complex) model-based solutions such as MV beamforming or extended/full acquisitions in a compressed sensing setup. In the former, training targets are generated offline by running powerful but time-consuming model-based beamformers on the training data set of RF inputs. The goal of deep-learning-based beamforming is then to achieve the same performance as these model-based solutions, at much faster inference rates. In the compressed sensing setup, training targets are generated by (DAS/MV) beamforming the full (not compressed) set of RF measurements X . In this case, the objective of a deep-learning beamformer is to reproduce these beamformed outputs based on compressed/undersampled measurements X_u . As discussed in Sec. 1.4, compression can entail fast-time sub-Nyquist sampling, imaging with sparse arrays, or limiting the number of transmit events in e.g. plane-wave compounding. Real data is available on the aforementioned CUBDL challenge website (Bell et al. 2020).

1.6.2 Loss functions and optimization

In this section, we will discuss typical loss functions used to train deep networks for ultrasound beamforming. Most networks are trained by directly optimizing a

loss on the beamformer output in the RF or IQ domain. Others indirectly optimize the beamformer output in a task-adaptive fashion by optimizing some downstream loss after additional processing. For training, some variant of stochastic gradient descent (SGD), often Adaptive Moment Estimation (ADAM) is used, and some form of learning rate decay is also common. SGD operates on mini batches, which comprise either full input-output image pairs, or some collection of patches/slices/cubes extracted from full images. In the following, we will (without loss of generality) use $Y^{(i)}$ and $Y_t^{(i)}$ to refer to respectively the network outputs and targets for a sample i .

Loss functions for beamformed outputs

Considering image reconstruction as a pixel-wise regression problem under a Gaussian likelihood model, perhaps the most commonly used loss function is the MSE (or ℓ_2 norm) with respect to the target pixel values:

$$\mathcal{L}_{MSE} = \frac{1}{I} \sum_{i=0}^{I-1} \left\| Y^{(i)} - Y_t^{(i)} \right\|_2^2. \quad (1.9)$$

If one would like to penalize strong deviations less stringently, e.g. to be less sensitive to outliers, one can consider a likelihood model that decays less strongly for large deviations, such as the Laplace distribution. Under that model, the negative log likelihood loss function is the mean absolute error (or ℓ_1 norm):

$$\mathcal{L}_{l1} = \frac{1}{I} \sum_{i=0}^{I-1} \left\| Y^{(i)} - Y_t^{(i)} \right\|_1. \quad (1.10)$$

A commonly adopted variant of the MSE loss is the signed-mean-squared-logarithmic-error SMSLE, proposed by Luijten et al. (Luijten et al. 2019b). This metric compresses the large dynamic range of backscattered ultrasound RF signals to promote accurate reconstructions across the entire dynamic range:

$$\mathcal{L}_{SMSLE} = \frac{1}{2I} \sum_{i=0}^{I-1} \left\| \log_{10} \left(Y^{(i)} \right)^+ - \log_{10} \left(Y_t^{(i)} \right)^+ \right\|_2^2 \quad (1.11)$$

$$+ \left\| \log_{10} \left(Y^{(i)} \right)^- - \log_{10} \left(Y_t^{(i)} \right)^- \right\|_2^2, \quad (1.12)$$

where $(\cdot)^+$ and $(\cdot)^-$ yield the magnitude of the positive and negative parts, respectively. Thus far, we have only covered pixel-wise losses that consider every pixel as an independent sample. These losses have no notion of spatial context and do not measure structural deviations. The structural similarity index (SSIM) (Wang, Bovik, Sheikh & Simoncelli 2004) aims to quantify perceived change in structural information, luminance and contrast. In the vein of the SMSLE, Kessler et al. (Kessler & Eldar 2020) propose a SSIM loss for ultrasound beamforming that acts on the log-compressed positive and negative parts of the

beamformed RF signals:

$$\mathcal{L}_{SSIM} = \frac{1}{2I} \sum_{i=0}^{I-1} \left(1 - SSIM \left(\log_{10} \left(Y^{(i)} \right)^+, \log_{10} \left(Y_t^{(i)} \right)^+ \right) \right) \quad (1.13)$$

$$+ \left(1 - SSIM \left(\log_{10} \left(Y^{(i)} \right)^-, \log_{10} \left(Y_t^{(i)} \right)^- \right) \right), \quad (1.14)$$

where, when luminance, contrast and structure are weighed equally, *SSIM* is defined as:

$$SSIM(a, b) = \frac{(2\mu_a\mu_b + \epsilon_1)(2\sigma_{ab} + \epsilon_2)}{(\mu_a^2 + \mu_b^2 + \epsilon_1)(\sigma_a^2 + \sigma_b^2 + \epsilon_2)}, \quad (1.15)$$

with μ_a , μ_b , σ_a , σ_b , and σ_{ab} being the means, standard deviations and cross-correlation of a and b , and ϵ_1 , ϵ_2 being small constants to stabilize the division.

Beyond distance measurements between pixel values, some authors make use of specific adversarial optimization schemes that aim to match the distributions of the targets and generated outputs (Chennakeshava et al. 2020). These approaches make use of a second neural network, the adversary or discriminator, that is trained to discriminate between images that are drawn from the distribution of targets, and those that are generated by the beamforming neural network. This is achieved by minimizing the binary cross-entropy classification loss between its predictions and the labels (target or generated), evaluated on batches that contain both target images and generated images. At the same time, the beamforming network is trained to maximize this loss, thereby attempting to fool the discriminator. The rationale here is that the probability distributions of target images and beamformed network outputs match (or strongly overlap) whenever this neural discriminator cannot distinguish images from either distribution anymore. The beamforming network and discriminator thus play a min-max game, expressed by the following optimization problem across the training data distribution $P_{\mathcal{D}}$:

$$\hat{\theta}, \hat{\Psi} = \underset{\psi}{\operatorname{argmin}} \underset{\theta}{\operatorname{argmax}} \left\{ - \mathbb{E}_{(X, Y_t) \sim P_{\mathcal{D}}} [\log(D_{\psi}(Y_t)) + \log(1 - D_{\psi}(f_{\theta}(X)))] \right\}, \quad (1.16)$$

where θ are the parameters of the beamforming network $f_{\theta}(\cdot)$ and ψ are the parameters of the discriminator D_{ψ} . It is important to realize that merely matching distributions does not guarantee accurate image reconstructions. That is why adversarial losses are often applied on input-output and input-target pairs (matching e.g. their joint distributions), or used in combination with additional distance metrics such as those discussed earlier in this section. The relative contributions or weighting of these individual loss terms is typically selected empirically.

Task-adaptive optimization

As discussed in Sec. 1.4, one can also optimize the parameters of beamforming

architectures using a downstream task-based loss, i.e.:

$$\hat{\theta}, \hat{\Phi} = \underset{\theta}{\operatorname{argmin}} \left\{ \mathbb{E}_{(X, s_{task}) \sim P_{\mathcal{D}}} \left[\mathcal{L}_{task} \left\{ g_{\phi} (f_{\theta}(X)), s_{task} \right\} \right] \right\}, \quad (1.17)$$

where s_t denotes some target task, $\mathcal{L}_{task}\{a, b\}$ is a task-specific loss function between outputs a and targets b , and ϕ are the parameters of the task function g_{ϕ} , which can be a neural network. Examples of such tasks include segmentation (Nair et al. 2020), for which \mathcal{L}_{task} is e.g. a Dice loss, or motion estimation (Doppler), for which \mathcal{L}_{task} is e.g. an MSE penalty (Huijben et al. 2020).

1.7 New Research Opportunities

1.7.1 Multi-functional deep beamformer

Although deep beamformers provide impressive performance and ultra-fast reconstruction, one of the downsides of the deep beamformers is that a distinct model is needed for each type of desired output. For instance, to obtain DAS outputs, a model is needed which mimics DAS; similarly, for MVBF a separate model is needed. Although the architecture of the model could be the same, separate weights need to be stored for each output type. Given that hundreds/thousands of B-mode optimizations/settings are used in the current high-end commercial systems, one may wonder whether we need to store thousands of deep models in the scanner to deal with various B-mode settings.

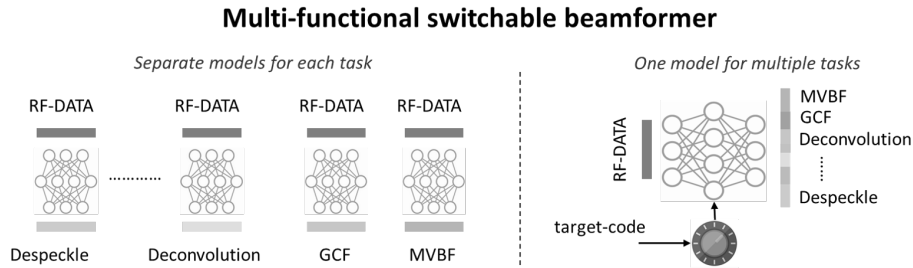


Figure 1.8 An illustration of switchable deep beamformer using AdaIN layer.

To address this issue, (Khan et al. 2020b) recently proposed a *switchable* deep beamformer architecture using adaptive instance normalization (AdaIN) layers as shown in Fig. 1.8. Specifically, AdaIN was originally proposed as an image style transfer method, in which the mean and variance of the feature vectors are replaced by those of the style reference image (Huang & Belongie 2017). Suppose that a multi-channel feature tensor at a specific layer is represented by

$$\mathbf{X} = [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_C] \in \mathbb{R}^{HW \times C}, \quad (1.18)$$

where C is the number of channels in the feature tensor \mathbf{X} , and $\mathbf{x}_i \in \mathbb{R}^{HW \times 1}$ refers to the i -th column vector of \mathbf{X} , which represents the vectorized feature map of size of $H \times W$ at the i -th channel. Then, AdaIN (Huang & Belongie 2017) converts the feature data at each channel using the following transform:

$$\mathbf{z}_i = \mathcal{T}(\mathbf{x}_i, \mathbf{y}_i), \quad i = 1, \dots, C \quad (1.19)$$

where

$$\mathcal{T}(\mathbf{x}, \mathbf{y}) := \frac{\sigma(\mathbf{y})}{\sigma(\mathbf{x})} (\mathbf{x} - m(\mathbf{x})\mathbf{1}) + m(\mathbf{y})\mathbf{1}, \quad (1.20)$$

where $\mathbf{1} \in \mathbb{R}^{HW}$ is the HW -dimensional vector composed of 1, and $m(\mathbf{x})$ and $\sigma(\mathbf{x})$ are the mean and standard deviation of $\mathbf{x} \in \mathbb{R}^{HW}$; $m(\mathbf{y})$ and $\sigma(\mathbf{y})$ refer to the target style domain mean and standard deviation, respectively. Eq. (1.20) implies that the mean and variance of the feature in the input image are normalized so that they can match the mean and variance of the style image feature. Although (1.20) looks heuristic, it was shown that the transform (1.20) is closely related to the optimal transport between two Gaussian probability distributions (Peyré, Cuturi et al. 2019, Villani 2008).

Inspired by this, (Khan et al. 2020b) demonstrate that a *single* deep beamformer with AdaIN layers can learn target images from various styles. Here, a “style” refers to a specific output processing, such as DAS, MVBF, deconvolution image, despeckled images, etc. Once the network is trained, the deep beamformer can then generate various style output by simply changing the AdaIN code. Furthermore, the AdaIN code generation is easily performed with a very light AdaIN code generator, so the additional memory overhead at the training step is minimal. Once the neural network is trained, we only need the AdaIN codes without the generator, which makes the system even simpler.

1.7.2 Unsupervised Learning

As discussed before, most existing deep learning strategies for ultrasound beamforming are based on supervised learning, thus relying predominantly on paired input-target datasets. However, in many real world imaging situations, access to paired images (input channel data and a corresponding desired output image) is not possible. For example, to improve the visual quality of US images acquired using a low-cost imaging system we need to scan exactly the same field of view using a high-end machine, which is not trivial. For denoising or artifact removal, the actual ground-truth is not known in-vivo, so supervised learning approaches are typically left with simulation datasets for training. This challenge has spurred a growing interest in developing an unsupervised learning strategy where channel data from low-end system or artifact corruption can be used as inputs, using surrogate performance metrics (based on high quality images from different machines and imaging conditions, or statistical properties) to train networks.

One possible approach to address this problem is to adopt unpaired style transfer strategies based on e.g. CycleGANs - a technique that has shown successful for many image domain quality improvements, also in ultrasound. For example, the authors in (Jafari, Girgis, Van Woudenberg, Moulson, Luong, Fung, Baltazaar, Jue, Tsang, Nair et al. 2020, Khan, Huh & Ye 2021) employed such a cycleGAN to improve the image quality from portable US image using high-end unmatched image data. In general, approaches that drive training by matching distribution properties (e.g. through discriminator networks as in CycleGAN) rather than strict input-output pairs hold promise for such applications.

References

- Asl, B. M. & Mahloojifar, A. (2010), 'Eigenspace-based minimum variance beamforming applied to medical ultrasound imaging', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **57**(11), 2381–2390.
- Austeng, A. & Holm, S. (2002), 'Sparse 2-d arrays for 3-d phased array imaging-design methods', *IEEE Trans. Ultrason., Ferroelectr., Freq. Control* **49**(8), 1073–1086.
- Bae, M., Park, S. B. & Kwon, S. J. (2016), 'Fast minimum variance beamforming based on legendre polynomials', *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* **63**(9), 1422–1431.
- Baran, J. M. & Webster, J. G. (2009), Design of low-cost portable ultrasound systems, *in* 'Annual International Conference of the IEEE Engineering in Medicine and Biology Society', IEEE, pp. 792–795.
- Bell, M., Huang, J., Hyung, D., Eldar, Y., van Sloun, R. & Mischi, M. (2020), 'Challenge on Ultrasound Beamforming by Deep Learning', cubd1.jhu.edu. [Online; accessed 15-September-2021].
- Bercoff, J., Tanter, M. & Fink, M. (2004), 'Supersonic shear imaging: a new technique for soft tissue elasticity mapping', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **51**(4), 396–409.
- Besson, A., Carrillo, R. E., Perdios, D., Arditi, M., Bernard, O., Wiaux, Y. & Thiran, J.-P. (2016), A compressed beamforming framework for ultrafast ultrasound imaging, *in* '2016 IEEE Int. Ultrasonics Symposium (IUS)', IEEE, pp. 1–4.
- Burshtein, A., Birk, M., Chernyakova, T., Eilam, A., Kempinski, A. & Eldar, Y. C. (2016), 'Sub-nyquist sampling and fourier domain beamforming in volumetric ultrasound imaging', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **63**(5), 703–716.
- Byram, B., Dei, K., Tierney, J. & Dumont, D. (2015), 'A model and regularization scheme for ultrasonic beamforming clutter reduction', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **62**(11), 1913–1927.
- Chennakeshava, N., Luijten, B., Drori, O., Mischi, M., Eldar, Y. C. & van Sloun, R. J. (2020), High resolution plane wave compounding through deep proximal learning, *in* '2020 IEEE International Ultrasonics Symposium (IUS)', IEEE, pp. 1–4.
- Chernyakova, T., Cohen, D., Shoham, M. & Eldar, Y. C. (2019), 'imap beamforming for high quality high frame rate imaging', *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* pp. 1–1.
- Chernyakova, T., Cohen, R., Mulayoff, R., Sde-Chen, Y., Frascini, C., Bercoff, J. & Eldar, Y. C. (2018), 'Fourier-domain beamforming and structure-based reconstruc-

- tion for plane-wave imaging', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **65**(10), 1810–1821.
- Chernyakova, T. & Eldar, Y. C. (2014), 'Fourier-domain beamforming: the path to compressed ultrasound imaging', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **61**(8), 1252–1267.
- Chollet, F. (n.d.), 'Antirectifier', <https://github.com/keras-team/keras/blob/master/examples/antirectifier.py>. Accessed: 14-03-2021.
- Christensen-Jeffries, K., Couture, O., Dayton, P. A., Eldar, Y. C., Hynynen, K., Kiessling, F., O'Reilly, M., Pinton, G. F., Schmitz, G., Tang, M.-X. et al. (2020), 'Super-resolution ultrasound imaging', *Ultrasound in medicine & biology* **46**(4), 865–891.
- Cohen, R. & Eldar, Y. C. (2018a), 'Sparse convolutional beamforming for ultrasound imaging', *IEEE Trans. Ultrason., Ferroelectr., Freq. Control* **65**(12), 2390–2406.
- Cohen, R. & Eldar, Y. C. (2018b), 'Sparse doppler sensing based on nested arrays', *IEEE Trans. Ultrason., Ferroelectr., Freq. Control* **65**(12), 2349–2364.
- Cohen, R. & Eldar, Y. C. (2020), 'Sparse array design via fractal geometries', *IEEE transactions on signal processing* **68**, 4797–4812.
- Demené, C., Deffieux, T., Pernot, M., Osmanski, B.-F., Biran, V., Gennisson, J.-L., Sieu, L.-A., Bergel, A., Franqui, S., Correas, J.-M. et al. (2015), 'Spatiotemporal clutter filtering of ultrafast ultrasound data highly increases doppler and ultrasound sensitivity', *IEEE transactions on medical imaging* **34**(11), 2271–2285.
- Eldar, Y. C. (2015), *Sampling theory: Beyond bandlimited systems*, Cambridge University Press.
- Eldar, Y. C. & Kutyniok, G. (2012), *Compressed sensing: theory and applications*, Cambridge University Press.
- Errico, C., Pierre, J., Pezet, S., Desailly, Y., Lenkei, Z., Couture, O. & Tanter, M. (2015), 'Ultrafast ultrasound localization microscopy for deep super-resolution vascular imaging', *Nature* **527**(7579), 499.
- Gedalyahu, K., Tur, R. & Eldar, Y. C. (2011), 'Multichannel sampling of pulse streams at the rate of innovation', *IEEE Transactions on Signal Processing* **59**(4), 1491–1504.
- Gregor, K. & LeCun, Y. (2010), Learning fast approximations of sparse coding, in 'Proceedings of the 27th International Conference on International Conference on Machine Learning', Omnipress, pp. 399–406.
- Holfort, I. K., Gran, F. & Jensen, J. A. (2009), 'Broadband minimum variance beamforming for ultrasound imaging', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **56**(2), 314–325.
- Huang, X. & Belongie, S. (2017), Arbitrary style transfer in real-time with adaptive instance normalization, in 'Proceedings of the IEEE International Conference on Computer Vision', pp. 1501–1510.
- Huijben, I. A., Veeling, B. S., Janse, K., Mischi, M. & van Sloun, R. J. (2020), 'Learning sub-sampling and signal recovery with applications in ultrasound imaging', *IEEE Transactions on Medical Imaging* **39**(12), 3955–3966.
- Hyun, D., Brickson, L. L., Looby, K. T. & Dahl, J. J. (2019), 'Beamforming and speckle reduction using neural networks', *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*.
- Hyun, D., Wiacek, A., Goudarzi, S., Rothlübbers, S., Asif, A., Eickel, K., Eldar, Y., Huang, J., Mischi, M., Rivaz, H., Sinden, D., van Sloun, R., Strohm, H. & Bell, M.

- (2021), ‘Deep learning for ultrasound image formation: Cubdl evaluation framework and open datasets’, *IEEE Transactions on Ultrasonics, Ferroelectrics and Frequency Control* (accepted) .
- Jafari, M. H., Girgis, H., Van Woudenberg, N., Moulson, N., Luong, C., Fung, A., Balthazaar, S., Jue, J., Tsang, M., Nair, P. et al. (2020), ‘Cardiac point-of-care to cart-based ultrasound translation using constrained cyclegan’, *International journal of computer assisted radiology and surgery* **15**(5), 877–886.
- Jensen, J. A. (2004), Simulation of advanced ultrasound systems using field ii, in ‘2004 2nd IEEE International Symposium on Biomedical Imaging: Nano to Macro (IEEE Cat No. 04EX821)’, IEEE, pp. 636–639.
- Kessler, N. & Eldar, Y. C. (2020), ‘Deep-learning based adaptive ultrasound imaging from sub-nyquist channel data’, *arXiv preprint arXiv:2008.02628* .
- Khan, S., Huh, J. & Ye, J. C. (2019), Deep learning-based universal beamformer for ultrasound imaging, in D. Shen, T. Liu, T. M. Peters, L. H. Staib, C. Essert, S. Zhou, P.-T. Yap & A. Khan, eds, ‘Medical Image Computing and Computer Assisted Intervention – MICCAI 2019’, Springer International Publishing, Cham, pp. 619–627.
- Khan, S., Huh, J. & Ye, J. C. (2020a), ‘Adaptive and compressive beamforming using deep learning for medical ultrasound’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **67**(8), 1558–1572.
- Khan, S., Huh, J. & Ye, J. C. (2020b), ‘Switchable deep beamformer’, *arXiv preprint arXiv:2008.13646* .
- Khan, S., Huh, J. & Ye, J. C. (2021), ‘Variational formulation of unsupervised deep learning for ultrasound image artifact removal’, *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* pp. 1–1.
- Kim, K., Park, S., Kim, J., Park, S. & Bae, M. (2014), ‘A fast minimum variance beamforming method using principal component analysis’, *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* **61**(6), 930–945.
- Kulkarni, K., Lohit, S., Turaga, P., Kerviche, R. & Ashok, A. (2016), Reconnet: Non-iterative reconstruction of images from compressively sensed measurements, in ‘Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition’, pp. 449–458.
- Lahav, A., Chernyakova, T. & Eldar, Y. C. (2017), ‘Focus: Fourier-based coded ultrasound’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **64**(12), 1828–1839.
- Liu, C.-L. & Vaidyanathan, P. (2017), Maximally economic sparse arrays and cantor arrays, in ‘2017 IEEE 7th Int. Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)’, IEEE, pp. 1–5.
- Lorintiu, O., Liebgott, H., Alessandrini, M., Bernard, O. & Friboulet, D. (2015), ‘Compressed sensing reconstruction of 3D ultrasound data using dictionary learning and line-wise subsampling’, *IEEE Trans. Med. Imag.* **34**(12), 2467–2477.
- Lu, J., Millioz, F., Garcia, D., Salles, S., Liu, W. & Friboulet, D. (2020), ‘Reconstruction for diverging-wave imaging using deep convolutional neural networks’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **67**(12), 2481–2492.
- Luchies, A. C. & Byram, B. C. (2018), ‘Deep neural networks for ultrasound beamforming’, *IEEE transactions on medical imaging* **37**(9), 2010–2021.
- Luchies, A. C. & Byram, B. C. (2019), ‘Training improvements for ultrasound beamforming with deep neural networks’, *Physics in Medicine & Biology* **64**(4), 045018.

-
- Luchies, A. C. & Byram, B. C. (2020), ‘Assessing the robustness of frequency-domain ultrasound beamforming using deep neural networks’, *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* **67**(11), 2321–2335.
- Luijten, B., Cohen, R., de Bruijn, F. J., Schmeitz, H. A., Mischi, M., Eldar, Y. C. & van Sloun, R. J. (2019a), Deep learning for fast adaptive beamforming, in ‘ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)’, IEEE, pp. 1333–1337.
- Luijten, B., Cohen, R., de Bruijn, F. J., Schmeitz, H. A., Mischi, M., Eldar, Y. C. & van Sloun, R. J. (2019b), Deep learning for fast adaptive beamforming, in ‘ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)’, IEEE, pp. 1333–1337.
- Mallart, R. & Fink, M. (n.d.), ‘Adaptive focusing in scattering media through sound-speed inhomogeneities: The van cittert zernike approach and focusing criterion’, *The Journal of the Acoustical Society of America* **96**(6), 3721–3732.
- Mamistvalov, A. & Eldar, Y. C. (2020), ‘Compressed fourier-domain convolutional beamforming for wireless ultrasound imaging’, *arXiv preprint arXiv:2010.13171* .
- Mamistvalov, A. & Eldar, Y. C. (2021), ‘Deep unfolded recovery of sub-nyquist sampled ultrasound image’, *arXiv preprint arXiv:2103.01263* .
- Michaeli, T. & Eldar, Y. C. (2012), ‘Xampling at the rate of innovation’, *IEEE Transactions on Signal Processing* **60**(3), 1121–1133.
- Mishali, M., Eldar, Y. C., Dounaevsky, O. & Shoshan, E. (2011), ‘Xampling: Analog to digital at sub-Nyquist rates’, *IET circuits, devices & systems* **5**(1), 8–20.
- Mishali, M., Eldar, Y. C. & Elron, A. J. (2011), ‘Xampling: Signal acquisition and processing in union of subspaces’, *IEEE Transactions on Signal Processing* **59**(10), 4719–4734.
- Nair, A. A., Washington, K. N., Tran, T. D., Reiter, A. & Bell, M. A. L. (2020), ‘Deep learning to obtain simultaneous image and segmentation outputs from a single input of raw ultrasound channel data’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **67**(12), 2493–2509.
- Nilsen, C. C. & Holm, S. (2010a), ‘Wiener beamforming and the coherence factor in ultrasound imaging’, *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control* **57**(6), 1329–1346.
- Nilsen, C.-I. C. & Holm, S. (2010b), ‘Wiener beamforming and the coherence factor in ultrasound imaging’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **57**(6), 1329–1346.
- Perdios, D., Besson, A., Arditi, M. & Thiran, J.-P. (2017), A deep learning approach to ultrasound image recovery, in ‘IEEE International Ultrasonics Symposium (IUS)’, Ieee, pp. 1–4.
- Peyré, G., Cuturi, M. et al. (2019), ‘Computational optimal transport: With applications to data science’, *Foundations and Trends® in Machine Learning* **11**(5-6), 355–607.
- Provost, J., Papadacci, C., Arango, J. E., Imbault, M., Fink, M., Gennisson, J.-L., Tanter, M. & Pernot, M. (2014), ‘3D ultrafast ultrasound imaging in vivo’, *Physics in Medicine & Biology* **59**(19), L1.
- Sutton, R. S. & Barto, A. G. (2018), *Reinforcement learning: An introduction*, MIT press.
- Szabo, T. L. (2004), *Diagnostic ultrasound imaging: inside out*, Academic Press.

- Tanter, M. & Fink, M. (2014), ‘Ultrafast imaging in biomedical ultrasound’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **61**(1), 102–119.
- Treeby, B. E. & Cox, B. T. (2010), ‘k-wave: Matlab toolbox for the simulation and reconstruction of photoacoustic wave fields’, *Journal of biomedical optics* **15**(2), 021314.
- Tur, R., Eldar, Y. C. & Friedman, Z. (2011), ‘Innovation rate sampling of pulse streams with application to ultrasound imaging’, *IEEE Transactions on Signal Processing* **59**(4), 1827–1842.
- Van Sloun, R. J., Cohen, R. & Eldar, Y. C. (2019), ‘Deep learning in ultrasound imaging’, *Proceedings of the IEEE* **108**(1), 11–29.
- Vignon, F., Shin, J. S., Meral, F. C., Apostolakis, I., Huang, S.-W. & Robert, J.-L. (2020), Resolution improvement with a fully convolutional neural network applied to aligned per-channel data, in ‘2020 IEEE International Ultrasonics Symposium (IUS)’, IEEE, pp. 1–4.
- Villani, C. (2008), *Optimal transport: old and new*, Vol. 338, Springer Science & Business Media.
- Wagner, N., Eldar, Y. C., Feuer, A., Danin, G. & Friedman, Z. (2011), Xampling in ultrasound imaging, in ‘Medical Imaging 2011: Ultrasonic Imaging, Tomography, and Therapy’, Vol. 7968, International Society for Optics and Photonics, p. 796818.
- Wagner, N., Eldar, Y. C. & Friedman, Z. (2012), ‘Compressed beamforming in ultrasound imaging’, *IEEE Transactions on Signal Processing* **60**(9), 4643–4657.
- Wang, Z., Bovik, A. C., Sheikh, H. R. & Simoncelli, E. P. (2004), ‘Image quality assessment: from error visibility to structural similarity’, *IEEE transactions on image processing* **13**(4), 600–612.
- Wiacek, A., González, E. & Bell, M. A. L. (2020), ‘Coherenet: A deep learning architecture for ultrasound spatial correlation estimation and coherence-based beamforming’, *IEEE transactions on ultrasonics, ferroelectrics, and frequency control* **67**(12), 2574–2583.
- Williams, R. J. (1992), ‘Simple statistical gradient-following algorithms for connectionist reinforcement learning’, *Machine learning* **8**(3-4), 229–256.
- Yoon, Y. H., Khan, S., Huh, J. & Ye, J. C. (2018), ‘Efficient b-mode ultrasound image reconstruction from sub-sampled rf data using deep learning’, *IEEE transactions on medical imaging* **38**(2), 325–336.
- Youn, J., Ommen, M. L., Stuart, M. B., Thomsen, E. V., Larsen, N. B. & Jensen, J. A. (2020), ‘Detection and localization of ultrasound scatterers using convolutional neural networks’, *IEEE Transactions on Medical Imaging* **39**(12), 3855–3867.