# TRULY SHIFT-EQUIVARIANT CONVOLUTIONAL NEURAL NETWORKS WITH ADAPTIVE POLYPHASE UPSAMPLING

*Anadi Chaman*

University of Illinois at Urbana-Chamaign
achaman2@illinois.edu

*Ivan Dokmanić*

University of Basel
ivan.dokmanic@unibas.ch

## ABSTRACT

Convolutional neural networks lack shift equivariance due to the presence of downsampling layers. In image classification, adaptive polyphase downsampling (APS-D) was recently proposed to make CNNs perfectly shift invariant. However, in networks used for image reconstruction tasks, it can not by itself restore shift equivariance. We address this problem by proposing adaptive polyphase upsampling (APS-U), a non-linear extension of conventional upsampling, which allows CNNs with symmetric encoder-decoder architecture (for example U-Net) to exhibit perfect shift equivariance. With MRI and CT reconstruction experiments, we show that networks containing APS-D/U layers exhibit state of the art equivariance performance without sacrificing on image reconstruction quality. In addition, unlike prior methods like data augmentation and anti-aliasing, the gains in equivariance obtained from APS-D/U also extend to images outside the training distribution.

***Index Terms***— Shift equivariant CNNs, shift invariance, adaptive polyphase sampling.

## 1. INTRODUCTION

In image-to-image regression problems like MRI and CT reconstruction, shifts in input to a convolutional neural network (CNN) should result in similar shifts in the network's output. This property, called shift equivariance, is a highly desirable inductive bias that allows networks to reconstruct objects accurately irrespective of their position in the frame. However, recent works have shown that despite the presence of convolutions which are shift equivariant, the output of a CNN can be very unstable to shifts in its input [1, 2]. This is because of the use of downsampling layers that lack translation equivariance. For example, Fig. 1(a) shows that shifting a signal can significantly change its downsampled output.

This problem has recently received attention in the context of shift invariance in image classification. Methods like data augmentation [3] and anti-aliasing [4, 5, 1] have been
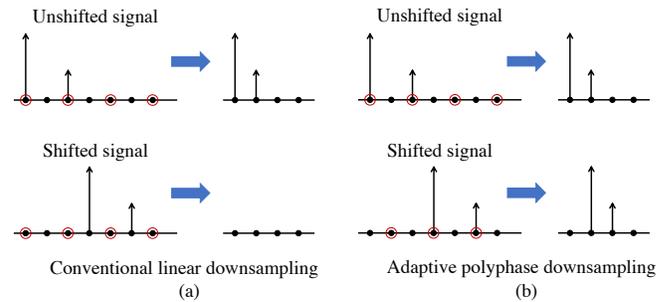
---

Code available at https://github.com/achaman2/truly_shift_invariant_cnns.



**Fig. 1**. (a) Conventional downsampling selects pixels along fixed locations on the sampling grid (shown with red circles). A shift in the signal changes pixels intensities on this fixed grid resulting in a very different output. (b) Adaptive polyphase downsampling selects the sampling grid adaptively such that a shift in input translates the downsampled output.

used to improve robustness to shifts. However, these methods have limitations and do not enable perfect invariance [1, 6]. In our recent work [6], we proposed adaptive polyphase downsampling (APS-D), a non-linear sampling scheme that allows CNN classifiers to exhibit perfect shift invariance. As shown in Fig. 1(b), by selecting its sampling grid in a signal dependent manner, APS-D consistently downsamples the same signal structures irrespective of any shift in input.

While APS-D enables shift invariance in classification, it can not by itself restore shift equivariance in CNNs used for image reconstruction tasks. This is because unlike shift invariance, where the goal is to obtain identical outputs for any shift in input, shift equivariance requires to propagate spatial location information from the network's input and output. Since APS-D reduces signal resolution on downsampling, shifts in the input which are not integer multiples of its stride can be mapped to the same output. Consequently, the shifts in input are lost when these low resolution feature maps pass through upsampling layers to produce the final output. To address this challenge, we propose adaptive polyphase upsampling (APS-U), a non-linear extension of classical upsampling, that allows CNNs with symmetric encoder-decoder structures (eg. the popular U-Net architecture [7]) to exhibit

perfect shift equivariance. APS-U achieves this by upsampling feature maps onto the same grid which was selected by APS-D for sampling in the encoder.

Using MRI and CT reconstruction experiments with U-Net, we show that our approach provides state-of-the art results in shift equivariance without sacrificing image reconstruction performance. While methods like anti-aliasing and data augmentation do improve robustness to shifts on average, there still exist shifts that can adversely impact reconstruction performance. This is a serious problem for applications like medical imaging which require strong robustness guarantees for the reconstructed images. These methods have also been shown to get adversely impacted by non-linear activations like ReLU [6]. On the other hand, our approach is highly robust to shifts and is not impacted by any point-wise non-linearity. In addition, our experiments reveal that networks with APS-D/U continue to remain shift equivariant on images outside the training distribution, which is not the case for anti-aliasing and data augmentation. Similar to APS-D, adaptive polyphase upsampling also does not require additional learnable parameters and can be easily integrated into existing architectures.

**Related work.** The success of convolutional neural networks inspired research on embedding equivariances to more complex transformations like rotations, scale, reflections and the action of arbitrary groups [8, 9, 10, 11, 12, 13]. Similarly, the related problem of invariance to deformations has been theoretically studied with kernel methods [14, 15] and wavelet filter banks [16, 17]. However, the impact of downsampling on the stability of CNNs to shifts has only recently been analyzed [1, 18]. Manfredi and Wang [19] assessed the lack of shift equivariance in CNN architectures used for object detection. Anti-aliasing and data augmentation were shown to improve shift invariance in classification [1, 4, 5]. However, the improved robustness to shifts from these methods does not extend to image patterns not seen during training [1]. The gains in shift invariance obtained from anti-aliasing were also shown to be limited by the action of non-linear activations in the network [1, 6]. These challenges were addressed in our earlier work [6] where we proposed a new downsampling method called adaptive polyphase sampling that enabled perfect shift invariance in CNN classifiers.

## 2. BACKGROUND

**Shift equivariance.** Let $x \in \mathbb{R}^d$ be an input to an operator $\mathcal{G}$, and $T_k$ represent translation by $k \in \mathbb{R}^d$. Then $\mathcal{G}$ is said to be equivariant to shifts if $\mathcal{G}(T_k(x)) = T_k(\mathcal{G}(x))$.

**Sampling and shift equivariance.** Let $D_2$ and $U_2$ denote downsampling and upsampling[1] operations with stride 2. For

a 1-D signal $x(n)$, $y = D_2(x)$ is given by $y(n) = x(2n)$. $U_2$ upsamples signal $y$ as

$$U_2(y) = z(n) = \begin{cases} y(n/2), & \text{when } n \text{ is even}, \\ 0, & \text{otherwise}. \end{cases} \quad (1)$$

Let $T_k = x(n - k)$ represent a $k$-pixel shift in $x$. For an odd shift $k = (2m + 1)$ with $m \in \mathbb{Z}$, $D_2$ satisfies $D_2(T_{2m+1}(x)) = T_m D_2(T_1(x))$. It is therefore, unsurprisingly, not shift equivariant. Similarly, $\forall k \in \mathbb{Z}$, the action of $U_2$ on $T_k(x)$ is given by

$$U_2(T_k(x)) = T_{2k}(U_2(x)). \quad (2)$$

The lack of equivariance in $D_2$ can not be corrected by anti-aliasing. Even if $x$ and $T_k(x)$ were ideal low pass filtered before downsampling, the subsequent sampled outputs $y_a$ and $y_a^{(k)}$ would have DTFTs

$$Y_a(\omega) = \frac{1}{2} X\left(\frac{\omega}{2}\right) \text{ and } Y_a^{(k)}(\omega) = \frac{1}{2} X\left(\frac{\omega}{2}\right) e^{-\frac{jk\omega}{2}}. \quad (3)$$

One can observe that for any odd shift $k$, there does not exist an $n_0 \in \mathbb{Z}$ such that $Y_a^{(k)}(\omega) = Y_a(\omega) e^{-jn_0\omega}$. This indicates that $y_a^{(k)}$ can not be obtained by translating $y_a$ with any integer shift.

**Adaptive polyphase downsampling**. Originally proposed in [6] to make CNN classifiers shift invariant, adaptive polyphase downsampling (APS-D) samples a 1-D signal $x$ by considering 2 possible sampling grids (illustrated in Fig. 1(b)). The signals supported on these grids are called polyphase components and are given by $y_0(n) = x(2n)$ and $y_1(n) = x(2n + 1)$. APS-D chooses the polyphase component of $x$ with the highest $l_p$ norm as its downsampled output $D_2^A(x)$, i.e. $y_{\text{APS}} = D_2^A(x) = y_i$ where $i = \text{argmax}_j \{\|y_j\|_p\}_{j=0}^1$. We show in [6] that shifting the input of APS-D always results in a shift in its output. More formally, if $y_{\text{APS}} = D_2^A(x)$ and $y_{\text{APS}}^{(k)} = D_2^A(T_k(x))$, then

$$y_{\text{APS}}^{(k)} = \begin{cases} T_{\frac{k}{2}}(y_{\text{APS}}), & \text{when } k \text{ is even}, \\ T_{\frac{k+2i-1}{2}}(y_{\text{APS}}), & \text{when } k \text{ is odd}, \end{cases} \quad (4)$$

where $i \in \{0, 1\}$ represents the polyphase component of $x$ with the highest norm. From (4), since $D_2^A(T_k(x)) \neq T_k(D_2^A(x))$ and the shift between the two signals is dependent on $i$, APS-D is not shift equivariant in the usual sense. Yet it is superior to classical downsampling which can not guarantee a shift in its output when its input is translated. We call this weaker condition as $\sigma$-equivariance.

---

[1] Stride-2 max-pooling layers used in modern CNNs can be decomposed into a dense (stride 1) max-pool followed by $D_2$. Similarly, transposed convolutions and nearest-neighbour upsampling layers can be characterized in terms of $U_2$.

## 3. OUR APPROACH

We will focus our discussion on shift equivariance in U-Net, a highly popular architecture used for image reconstruction tasks [20, 21, 22]. The general conclusions can be extended to similar architectures containing symmetric encoder-decoder structure. The U-Net's encoder takes a signal $x$ as input and generates a multi-scale decomposition with $L$ scales. Feature map $x_e^{(l)}$ at scale $l$ is used by a convolutional shift equivariant block $F_e^{(l)}$ to generate $s_e^{(l)} = F_e^{(l)}(x_e^{(l)})$ which is then downsampled to produce $x_e^{(l+1)} = D_2(s_e^{(l)})$ for the next scale. Similarly, at scale $l$ in the decoder, a convolutional $F_d^{(l)}$ generates $x_d^{(l)}$ from upsampled feature map $s_d^{(l)} = U_2(x_d^{(l+1)})$ and skip connection $s_e^{(l)}$.

By replacing the downsampling layers of U-Net with APS-D, we can make its encoder $\sigma$-equivariant. Then, a shift $k$ in input shifts feature maps $\{x_e^{(l)}\}_{l=1}^{L}$ by $\{k_l\}_{l=1}^{L}$, where the relation between $k_l$ and $k_{l-1}$ can be obtained via (4).

### 3.1. A case for adaptive upsampling

While APS-D can make a U-Net's encoder $\sigma$-equivariant, it can not restore perfect shift equivariance in the overall U-Net architecture. To understand this better, consider a toy example with a signal $x$ and its 1-pixel shift $\tilde{x} = T_1(x)$ as shown in Fig. 2(a). Downsampling them with APS-D produces $y_{\text{APS}} = D_2^A(x)$ and $\tilde{y}_{\text{APS}} = D_2^A(\tilde{x})$ which correspond to polyphase components with indices $i = 0$ and 1 respectively. Notice from Fig. 2(b) that despite being sampled from different spatially located grids, $y_{\text{APS}}$ and $\tilde{y}_{\text{APS}}$ are identical and the 1-pixel shift between the inputs is lost. Therefore, upsampling them directly results in outputs $z$ and $\tilde{z}$ which are not 1-pixel shifted versions of each other (Fig. 2(c)).

One way to counter this is by upsampling the signals back to the grid from which they were originally sampled. For example, in Fig. 2(c), $\tilde{z}$ obtained after upsampling $\tilde{y}_{\text{APS}}$, can be translated to the grid with index $i = 1$. Similarly, $z$ can be shifted to $i = 0$. One then obtains outputs $z_{\text{APS}}$ and $\tilde{z}_{\text{APS}}$ which from Fig. 2(d) have the same shift between them as the original inputs $x$ and $\tilde{x}$.

The above example shows that by using the index of polyphase component used by APS-D for downsampling, one can upsample signals to the 'correct' sampling grids, and preserve spatial location information. We formalize this notion in the next section.

### 3.2. Adaptive polyphase upsampling

For simplicity, we will focus our discussion to sampling of 1-D signals with stride 2. The conclusions can be easily generalized to higher dimensions and stride. Let $\{y_i\}_{i=0}^{1}$ be the polyphase components of a 1-D signal $x$, given by $y_i(n) = D_2(T_{-i}(x)) = x(2n+i)$, where $i \in \{0,1\}$. We denote $i_x$ to be the index of polyphase component with the highest
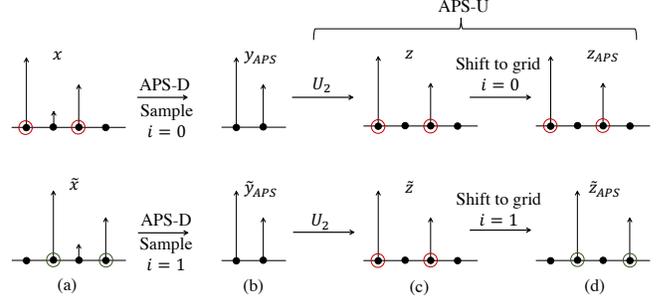


**Fig. 2**. (a)-(b) The relative shift between $x$ and its 1-pixel shift $\tilde{x}$ is lost after they are downsampled. (c) Outputs of conventional upsampling, therefore, are not 1-pixel shifted versions of each other. (d) By shifting the upsampled signal to the grid originally chosen by APS-D, shift equivariance is achieved.

norm. Then the downsampled output obtained from APS-D is $y_{\text{APS}} = D_2^A(x) = D_2(T_{-i_x}(x))$. Let $U_2^A$ denote adaptive polyphase upsampling (APS-U) operator with stride 2. Upsampling $y_{\text{APS}}$ with APS-U yields $U_2^A(y_{\text{APS}}, i_x)$ which is defined as

$$U_2^A(y_{\text{APS}}, i_x) = T_{i_x}\Big(U_2(y_{\text{APS}})\Big). \qquad (5)$$

The example in Fig. 2 illustrates that $U_2^A(D_2^A(x))$ is equivalent to zeroing out all pixels of $x$ except the ones located on the grid $i_x$. Since APS-D selects sampling grid in a manner consistent to translations, one can expect a shift in $x$ to result in the same shift in $U_2^A(D_2^A(x))$. Indeed, we show in Proposition 1 that $U_2^A \circ D_2^A$ is a shift equivariant operator.

**Proposition 1.** *Let $D_2^A$ and $U_2^A$ denote APS-D and APS-U operators with stride 2. Then $U_2^A \circ D_2^A$ is shift equivariant, i.e.*

$$U_2^A \circ D_2^A(T_k(x)) = T_k(U_2^A \circ D_2^A(x)), \ \forall k \in \mathbb{Z}. \quad (6)$$

*Proof.* Let $x$ and its $k$-pixel shift $\tilde{x} = T_k(x)$ be inputs to $U_2^A \circ D_2^A$. Without loss of generality, we assume $k$ to be an odd integer. Let $i_x$ and $\tilde{i}_x$ be the polyphase component indices of $x$ and $\tilde{x}$ respectively with the highest norm. Then

$$D_2^A(x) = D_2(T_{-i_x}(x)) \text{ and } D_2^A(\tilde{x}) = D_2(T_{-\tilde{i}_x}(\tilde{x})). \quad (7)$$

From [6], shifting a signal with an odd shift $k$, permutes the order of its polyphase components. Therefore, $\tilde{i}_x = 1 - i_x$. We can now write

$$U_2^A \circ D_2^A(x) = T_{i_x} U_2 D_2(T_{-i_x}(x)). \qquad (8)$$

Similarly,

$$U_2^A \circ D_2^A(\tilde{x}) = T_{\tilde{i}_x} U_2 D_2(T_{-\tilde{i}_x}(\tilde{x})) \qquad (9)$$

$$= T_{\tilde{i}_x} U_2 D_2(T_{k-\tilde{i}_x}(x)) \qquad (10)$$

$$= T_{\tilde{i}_x} T_{k-1} U_2 D_2(T_{1-\tilde{i}_x}(x)) \qquad (11)$$

$$= T_k(T_{i_x} U_2 D_2(T_{-i_x}(x))), \qquad (12)$$

where (11) and (12) follow from identities

$$D_2(T_{2m+1}(x)) = T_m(D_2(T_1(x))),$$
$$U_2(T_m(x)) = T_{2m}(U_2(x))),$$

and $\widetilde{i}_x = 1 - i_x$ (for odd $k$). From (8) and (12), $U_2^A \circ D_2^A(T_k(x)) = T_k(U_2^A \circ D_2^A(x))$. The result can similarly be shown for even $k$ in which case $\widetilde{i}_x = i_x$. $\qquad \square$

### 3.3. Restoring shift equivariance with APS-U

We saw in Section 3.2 that the composition of APS-U and APS-D given by $U_2^A \circ D_2^A$ is shift equivariant. We will now show that this property allows a U-Net containing APS-D/U layers to be perfectly shift equivariant.

**Proposition 2.** *A U-Net architecture with downsampling and upsampling layers replaced by APS-D and APS-U layers is shift equivariant.*

*Proof.* We first prove the claim for a U-Net containing a single down and upsampling layer, i.e., with $L = 1$ as shown in Fig. 3. Equivariance for higher $L$ follows by induction.

Let input $x$ to the network produce output $y$ and internal feature maps $s_e^{(0)}$, $x_e^{(1)}$, $x_d^{(1)}$ and $s_d^{(0)}$. Similarly, the output and feature maps for $T_k(x)$ are denoted by $\tilde{y}$ and $\tilde{s}_e^{(0)}$, $\tilde{x}_e^{(1)}$, $\tilde{x}_d^{(1)}$ and $\tilde{s}_d^{(0)}$. Since $F_e^{(0)}$ is convolutional, $\tilde{s}_e^{(0)} = T_k(s_e^{(0)})$. Let $i_x$ and $\widetilde{i}_x$ be the indices of polyphase components of $s_e^{(0)}$ and $\tilde{s}_e^{(0)}$ with the highest norm. Then,

$$x_e^{(1)} = D_2^A(s_e^{(0)}) = D_2(T_{-i_x}(s_e^{(0)})) \qquad (13)$$
$$\tilde{x}_e^{(1)} = D_2^A(\tilde{s}_e^{(0)}) = D_2(T_{-\widetilde{i}_x}(\tilde{s}_e^{(0)})) \qquad (14)$$

From the proof of Proposition 1, $\widetilde{i}_x = 1 - i_x$ when $k$ is odd and $\widetilde{i}_x = i_x$ for even $k$. For odd $k$, we have

$$s_d^{(0)} = U_2^A(F_e^{(1)} x_e^{(1)}) = T_{i_x} U_2(F_e^{(1)} x_e^{(1)}) \qquad (15)$$
$$= T_{i_x} U_2(F_e^{(1)} D_2(T_{-i_x}(s_e^{(0)}))) \qquad (16)$$

Similarly,

$$\tilde{s}_d^{(0)} = U_2^A(F_e^{(1)} \tilde{x}_e^{(1)}) = T_{\widetilde{i}_x} U_2(F_e^{(1)} D_2(T_{k-\widetilde{i}_x}(s_e^{(0)}))). \qquad (17)$$

Using the shift equivariance of $F_e^{(1)}$ and arguments similar to the proof of Proposition 1, one can show that

$$\tilde{s}_d^{(0)} = T_{k+i_x} U_2(F_e^{(1)} D_2(T_{-i_x}(s_e^{(0)}))) = T_k(s_d^{(0)}) \qquad (18)$$

Convolutional $F_d^{(0)}$ combines $s_d^{(0)}$ and skip connection $s_e^{(0)}$ to generate $y$. Since $\tilde{s}_e^{(0)} = T_k(s_e^{(0)})$ and $\tilde{s}_d^{(0)} = T_k(s_d^{(0)})$, this results in $\tilde{y} = T_k(y)$. We can similarly prove the result for even shift $k$ by using $\widetilde{i}_x = i_x$. The case of $L > 1$ follows by induction. $\qquad \square$
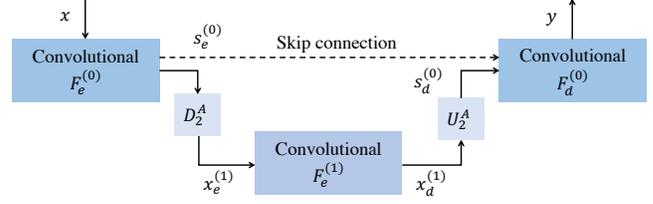


**Fig. 3**. U-Net containing a single APS-D and APS-U layer.

The proofs of Propositions 1 and 2 rely on the fact that $\widetilde{i}_x = 1 - i_x$ when $k$ is odd. This condition is true as long as APS-D performs downsampling consistently. However, as pointed in [6], in the event when two polyphase components have exactly identical $l_p$ norms, this might fail to occur. However, in theory, if we assume polyphase components to be drawn from a continuous distribution, this is an event of probability zero. It is also very less likely in practice. As suggested in [6], one could avoid this rare event by selecting polyphase components using methods more robust than simple norm maximization.

## 4. EXPERIMENTS

We evaluate the performance of U-Net on MRI and CT reconstruction tasks with conventional (baseline) and APS-D/U sampling layers. For the two tasks, zero filled and filtered back projected images respectively are provided as input to the U-Net which then generates the final reconstructions. We also combine both the baseline and APS layers with anti-aliased filters of size $2 \times 2$, $3 \times 3$ and $5 \times 5$ similar to [4]. Baseline models with low pass filters of size $j \times j$ are labelled as LPF-j, and those containing APS are denoted by APS-j. The baseline U-Net's encoder contains filters [64, 128, 256, 512, 1024] with 4 strided-maxpool layers, and the decoder uses transposed convolutions for upsampling. Similar architecture is used for the LPF and APS based U-Net variants except that their stride layers are replaced by anti-aliased and APS based sampling respectively. All the networks were trained with MSE loss function and without any random shifts, unless mentioned otherwise. Networks trained with random shifts are denoted by DA. Further details on training and implementation are available in the code provided in `https://github.com/achaman2/truly_shift_invariant_cnns`.

In addition to downsampling, CNNs can lose shift equivariance/invariance due to boundary effects as well [6]. Therefore, to separate the impact of boundary artifacts and sampling we use circular padded convolutions and evaluate the networks for equivariance with circular shifts.

We compared the models on two categories of metrics.

- **Equivariance metrics**: For inputs $x$ and $T_k(x)$ to a U-Net $G$, we use the SSIM and NMSE between

4

| | Equivariance metrics | | | | Reconstruction metrics (unshifted) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | NMSE | | SSIM | | NMSE | | PSNR | | SSIM | |
| **Model** | PD | PDFS | PD | PDFS | PD | PDFS | PD | PDFS | PD | PDFS |
| Baseline | 0.0014 | 9.56e-4 | 0.9965 | 0.9975 | 0.016 | 0.053 | 33.83 | 29.92 | 0.8093 | 0.6301 |
| Baseline + DA | 1.37e-4 | 1.26e-4 | 0.9987 | 0.9990 | 0.016 | 0.053 | 33.58 | 29.86 | 0.8034 | 0.6272 |
| LPF-2 | 7.43e-4 | 5.34e-4 | 0.9978 | 0.9985 | 0.016 | 0.052 | 33.93 | 29.96 | 0.8122 | 0.6323 |
| LPF-3 | 4.97e-4 | 4.13e-4 | 0.9984 | 0.9988 | **0.016** | **0.052** | **33.95** | **29.96** | **0.8125** | **0.6325** |
| LPF-5 | 5.05e-5 | 3.47e-5 | 0.9997 | 0.9998 | 0.018 | 0.055 | 33.19 | 29.76 | 0.7951 | 0.6225 |
| APS | **1.21e-7** | **7.37e-15** | **1.0** | **1.0** | 0.017 | 0.054 | 33.4 | 29.79 | 0.8013 | 0.6244 |
| APS-2 | **1.25e-7** | **5.86e-8** | **1.0** | **1.0** | 0.017 | 0.054 | 33.51 | 29.83 | 0.8023 | 0.6255 |
| APS-3 | **3.10e-7** | **1.36e-7** | **1.0** | **1.0** | **0.016** | **0.052** | **33.95** | **29.96** | **0.8124** | **0.6325** |
| APS-5 | **6.49e-7** | **2.74e-7** | **1.0** | **1.0** | 0.016 | 0.054 | 33.87 | 29.88 | 0.8088 | 0.6282 |

**Table 1**. Equivariance and reconstruction metrics obtained with different variants of U-Net on fastMRI validation set.

| Model | Baseline | APS | LPF-2 | APS-2 | LPF-3 | APS-3 | LPF-5 | APS-5 | Baseline + DA |
|---|---|---|---|---|---|---|---|---|---|
| ($\Delta$PSNR) | 4.03 | **8.6e-4** | 2.58 | **5.87e-4** | 2.36 | **3.81e-3** | 0.062 | **3.69e-3** | 0.037 |

**Table 2**. Worst case decline in PSNR of MRI reconstructions caused by randomly shifting the images in fastMRI validation set.

| Model | NMSE | SSIM |
|---|---|---|
| Baseline | 8.67e-3 | 0.9722 |
| Baseline+DA | 1.87e-3 | 0.9882 |
| LPF-2 | 4.65e-3 | 0.9816 |
| LPF-3 | 3.20e-3 | 0.9861 |
| LPF-5 | 4.40e-4 | 0.9979 |
| APS | **3.09e-14** | **1.0** |
| APS-2 | **3.20e-14** | **1.0** |
| APS-3 | **3.39e-7** | **1.0** |
| APS-5 | **2.39e-6** | **1.0** |

**Table 3**. Equivariance metrics for networks trained on fastMRI training set but evaluated on ImageNet validation set.

$T_k(G(x))$ and $G(T_k(x))$ averaged over the dataset to evaluate shift equivariance of the U-Net.

We also examine the possible decline in PSNR ($\Delta PSNR$) of image reconstructions caused by shifting the U-Net's input.

- **Reconstruction metrics**: To ensure that equivariance gains do not cause any sacrifice in reconstruction performance, we measure NMSE, PSNR and SSIM of the reconstructions for unshifted images.

### 4.1. MRI reconstruction

We train and evaluate different variants of U-Net on FastMRI single coil knee reconstruction task [22]. The networks were trained on a dataset containing 34742 images and evaluated on the validation set with 7135 images. The images were of size $320 \times 320$. Equivariance metrics were evaluated for each image with random shifts between $-16$ to 16 and an average was computed over the 2 partitions provided in the dataset—'PDFS' and 'PD'. Similar to [22], the reconstruction and equivariance metrics were computed over the provided image volumes rather than individual images.

Table 1 shows that models which use APS layers exhibit orders of magnitude lower equivariance errors than all other downsampling variants while still observing reconstruction performance comparable to baseline. In fact, they even surpass networks trained with random shifts (DA) on equivariance metrics by a large margin.

In addition to the equivariance metrics averaged over the entire dataset, we also assess worst case metrics, i.e. we shift each image in the validation set with 10 different random shifts and measure the worst absolute change in PSNR of the corresponding reconstructions obtained from U-Nets with different sampling modules. The worst possible decline for any image in the dataset is reported in Table 2. The results indicate that there exist shifts which can significantly impact the PSNR of reconstructions with baseline and LPF models. Networks with APS on the other hand are significantly more robust to shifts. An example illustration is provided in Fig. 4. We can observe in Fig. 4(c) that an (8, 5) pixel shift in the baseline network's input changes the pixel intensity distribution of its output and results in a decline in its PSNR with respect to the shifted ground truth. Consequently, the new reconstruction of the network is not a shifted version of its previous output. In contrast, as shown in Fig. 4(d), the two reconstructions obtained using APS U-Net are (8, 5) pixel shifted versions of each other.
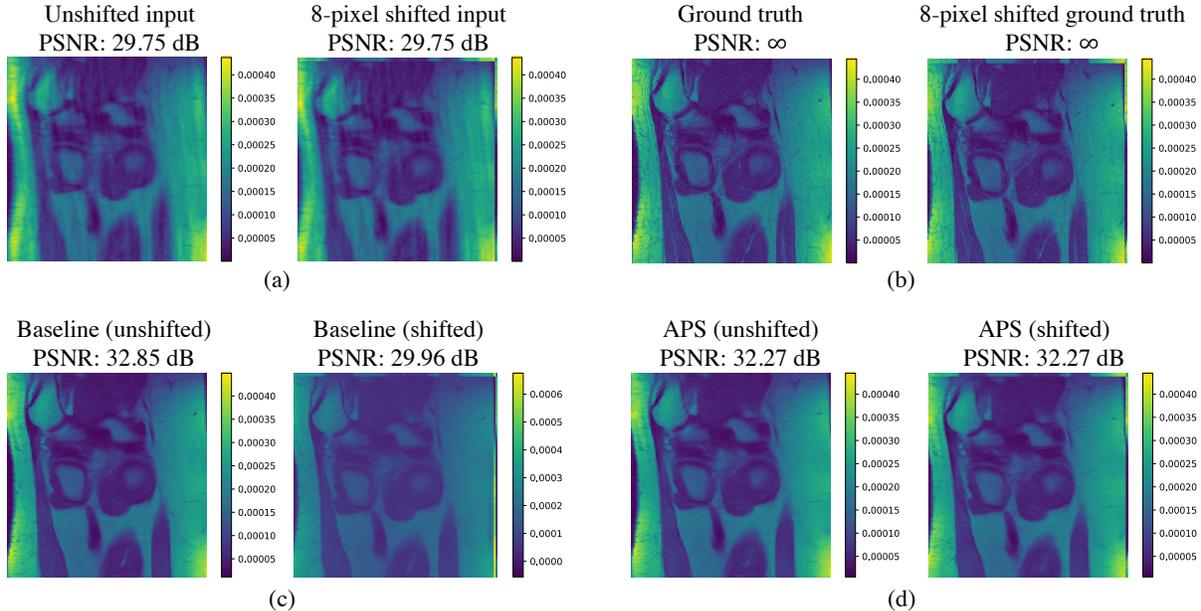
**Fig. 4**. Decline in PSNR of MRI reconstruction caused by shift in U-Net's input. (a) Input to the network and its (8, 5) pixel shifted version. (b) Ground truth. (c) Shifting the input to baseline U-Net results in a change in the pixel intensity distribution of its output and decline in PSNR. (d) Output of U-Net containing APS-D/U layers is highly robust to shifts.

### 4.1.1. Out-of-distribution equivariance

In image classification, gains in shift invariance obtained from data augmentation and anti-aliasing are known to not extend well on out-of-distribution images [1, 6]. Here, we observe a similar phenomenon for shift equivariance with U-Net. We take networks trained on the fastMRI dataset and evaluated equivariance metrics on the first 1000 images from ImageNet validation set [23]. Table 3 shows that SSIM equivariance metric on ImageNet for baseline network trained with data augmentation is $0.9882$, when it was $0.9990$ on the fastMRI dataset. We observe a similar decline for anti-aliased models as well. On the other hand, APS continues to provide SSIM of $1.0$ and orders of magnitude lower NMSE on the ImageNet dataset as well. The gains in shift equivariance provided by APS-D/U are therefore far more generalizable to out-of-dataset distributions in comparison to data augmentation and anti-aliasing.

### 4.2. CT reconstruction

We trained U-Net with different downsampling variants on the LoDoPaB-CT dataset [24] to perform CT reconstruction. The networks were trained on a dataset containing 35820 images and evaluated on the test set with 3553 images. We crop each image in the training and test set to size $352 \times 352$ to ensure feature maps with even dimensions inside the U-Net. This was done to avoid boundary artifacts that arise when downsampling an odd length signal and its circular shifted version [6]. Table 4 shows that similar to MRI reconstruction,

networks with APS-D/U layers outperform the other models on shift equivariance, while performing comparably on reconstruction performance for unshifted images.

## 5. CONCLUSIONS

Convolutional neural networks lose shift equivariance due to the presence of downsampling layers. While classical methods like data augmentation and anti-aliasing can improve shift equivariance on average, we show that they are not effective against all shifts. In addition, equivariance gains obtained with these methods are limited by the action of non-linear activations and do not necessarily extend well to image patterns not seen during training. In this work, we propose adaptive polyphase upsampling (APS-U) and combine it with our recently proposed adaptive polyphase downsampling (APS-D) scheme to enable perfect shift equivariance in symmetric encoder-decoder CNNs. Using experiments on MRI and CT reconstruction with U-Net architecture, we show that our approach significantly outperforms prior methods in improving translation equivariance. We also observe that the equivariance gains extend to out-of-distribution images and do not cause any sacrifice in reconstruction performance.

## Acknowledgement

| Equivariance metrics | | | Reconstruction metrics | | |
|---|---|---|---|---|---|
| **Model** | **NMSE** | **SSIM** | **NMSE** | **PSNR** | **SSIM** |
| Baseline | 1.2e-4 | 0.9961 | 0.0056 | 34.76 | 0.8309 |
| Baseline+DA | 1.97e-5 | 0.9989 | 0.0056 | 34.7 | 0.8252 |
| LPF-2 | 2.95e-5 | 0.9991 | **0.0052** | **35.35** | **0.836** |
| LPF-3 | **2.83e-7** | **1.0** | 0.0059 | 34.31 | 0.8191 |
| LPF-5 | 1.67e-06 | 0.9999 | 0.0058 | 34.43 | 0.8239 |
| APS | **3.55e-8** | **1.0** | 0.0055 | 34.84 | 0.8214 |
| APS-2 | **1.53e-8** | **1.0** | **0.0052** | **35.31** | **0.8334** |
| APS-3 | **6.495e-10** | **1.0** | 0.0054 | 34.95 | 0.8206 |
| APS-5 | **8.76e-9** | **1.0** | 0.0056 | 34.69 | 0.8308 |

**Table 4**. Equivariance and reconstruction metrics evaluated over LoDoPaB-CT dataset.

# 6. REFERENCES

[1] Aharon Azulay and Yair Weiss, "Why do deep convolutional networks generalize so poorly to small image transformations?," *Journal of Machine Learning Research*, vol. 20, no. 184, pp. 1–25, 2019.

[2] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry, "A rotation and a translation suffice: Fooling CNNs with simple transformations," 2019.

[3] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2015.

[4] Richard Zhang, "Making convolutional networks shift-invariant again," in *Proceedings of the 36th International Conference on Machine Learning*. 09–15 Jun 2019, vol. 97 of *Proceedings of Machine Learning Research*, pp. 7324–7334, PMLR.

[5] Xueyan Zou, Fanyi Xiao, Zhiding Yu, and Yong Jae Lee, "Delving deeper into anti-aliasing in convnets," in *BMVC*, 2020.

[6] Anadi Chaman and Ivan Dokmanic, "Truly shift-invariant convolutional neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 3773–3783.

[7] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Cham, 2015, pp. 234–241, Springer International Publishing.

[8] Taco Cohen and Max Welling, "Group equivariant convolutional networks," in *Proceedings of The 33rd International Conference on Machine Learning*, New York, New York, USA, 20–22 Jun 2016, vol. 48 of *Proceedings of Machine Learning Research*, pp. 2990–2999, PMLR.

[9] Taco Cohen, Maurice Weiler, Berkay Kicanaoglu, and Max Welling, "Gauge equivariant convolutional networks and the icosahedral CNN," in *Proceedings of the 36th International Conference on Machine Learning*. 09–15 Jun 2019, vol. 97 of *Proceedings of Machine Learning Research*, pp. 1321–1330, PMLR.

[10] Taco S. Cohen, Mario Geiger, Jonas Köhler, and Max Welling, "Spherical CNNs," in *International Conference on Learning Representations*, 2018.

[11] Maurice Weiler, Fred A. Hamprecht, and Martin Storath, "Learning steerable filters for rotation equivariant cnns," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[12] Ivan Sosnovik, Michał Szmaja, and Arnold Smeulders, "Scale-equivariant steerable networks," in *International Conference on Learning Representations*, 2020.

[13] Siamak Ravanbakhsh, Jeff Schneider, and Barnabás Póczos, "Equivariance through parameter-sharing," in *Proceedings of the 34th International Conference on Machine Learning*. 06–11 Aug 2017, vol. 70 of *Proceedings of Machine Learning Research*, pp. 2892–2901, PMLR.

[14] Alberto Bietti and Julien Mairal, "Group invariance, stability to deformations, and complexity of deep convolutional representations," *J. Mach. Learn. Res.*, vol. 20, no. 1, pp. 876–924, Jan. 2019.

[15] Alberto Bietti and Julien Mairal, "Invariance and stability of deep convolutional representations," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2017, NIPS'17, p. 6211–6221, Curran Associates Inc.

[16] Stéphane Mallat, "Group invariant scattering," *Communications on Pure and Applied Mathematics*, vol. 65, no. 10, pp. 1331–1398, 2012.

[17] Joan Bruna and Stephane Mallat, "Invariant scattering convolution networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1872–1886, 2013.

[18] Logan Engstrom, Brandon Tran, Dimitris Tsipras, Ludwig Schmidt, and Aleksander Madry, "Exploring the landscape of spatial robustness," in *Proceedings of the 36th International Conference on Machine Learning*. 09–15 Jun 2019, pp. 1802–1811, PMLR.

[19] Marco Manfredi and Yu Wang, "Shift equivariance in object detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 32–45.

[20] Kyong Hwan Jin, Michael T. McCann, Emmanuel Froustey, and Michael Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE Transactions on Image Processing*, vol. 26, no. 9, pp. 4509–4522, 2017.

[21] George Barbastathis, Aydogan Ozcan, and Guohai Situ, "On the use of deep learning for computational imaging," *Optica*, vol. 6, no. 8, pp. 921–943, Aug 2019.

[22] Jure Zbontar, Florian Knoll, Anuroop Sriram, Tullie Murrell, Zhengnan Huang, Matthew J. Muckley, Aaron Defazio, Ruben Stern, Patricia Johnson, Mary Bruno, Marc Parente, Krzysztof J. Geras, Joe Katsnelson, Hersh Chandarana, Zizhao Zhang, Michal Drozdzal, Adriana Romero, Michael Rabbat, Pascal Vincent, Nafissa Yakubova, James Pinkerton, Duo Wang, Erich Owens, C. Lawrence Zitnick, Michael P. Recht, Daniel K. Sodickson, and Yvonne W. Lui, "fastMRI: An open dataset and benchmarks for accelerated MRI," 2018.

[23] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[24] Johannes Leuschner, Maximilian Schmidt, Daniel Otero Baguer, and Peter Maaß, "The lodopab-ct dataset: A benchmark dataset for low-dose ct reconstruction methods," *arXiv preprint arXiv:1910.01113*, 2019.