

# On the Utility and Limitations of the MSTAR Dataset for Deep Learning–Based SAR Target Recognition

Charles Milligan

Research Institute for Tactical Autonomy, Howard University, Washington D.C USA.

## Abstract

The DARPA Moving and Stationary Target Acquisition and Recognition (MSTAR) dataset has served as a foundational benchmark for synthetic aperture radar (SAR)–based automatic target recognition (ATR) research for more than two decades. In recent years, it has been widely adopted for training and evaluating deep learning models, with reported classification accuracies often exceeding 95%. While these results demonstrate the effectiveness of modern neural architectures, they also raise concerns regarding the representativeness and continued utility of MSTAR as a proxy for operational SAR ATR tasks. In this paper, we examine the strengths and limitations of the MSTAR dataset in the context of deep learning–based model development. We present representative performance results from a convolutional neural network trained on MSTAR imagery, illustrating how dataset characteristics such as centered targets, fixed chip dimensions, limited clutter, and single-object scenes lead to artificially easy classification problems. We further contrast these results with preliminary experiments on the ATRNet-STAR dataset, where modest perturbations including target offset and rotation produce substantial performance degradation. Collectively, these findings highlight the need for more realistic SAR datasets and evaluation protocols. This paper is intended as a work in progress, with ongoing analysis of embedding-space structure and robustness metrics to be reported in future revisions.

## Introduction

Synthetic aperture radar (SAR) remains a critical sensing modality for military and intelligence applications due to its ability to operate in all weather conditions and independent of illumination. Automatic target recognition (ATR) using SAR imagery has therefore been an enduring research focus, spanning model-based, feature-engineered, and, more recently, deep learning approaches.

Among available SAR datasets, the DARPA MSTAR dataset has emerged as the de facto benchmark (Ross et al.

1998; Chen et al. 2016) for algorithm development and comparative evaluation. Originally curated to support classical ATR research, MSTAR provides labeled SAR image chips of military ground vehicles collected under controlled conditions. Its accessibility and consistency have made it attractive for training convolutional neural networks (CNNs) and related deep learning architectures.

However, as the field matures, there is growing recognition that high performance on MSTAR does not necessarily translate to robust operational capability. Many published results implicitly conflate success on MSTAR with general SAR ATR proficiency, despite well-documented dataset constraints. This paper argues that MSTAR should be viewed primarily as a model development and debugging dataset, rather than a reliable indicator of real-world performance. The contributions of this paper are threefold:

1. We summarize key properties of the MSTAR dataset that simplify the learning problem for deep networks.
2. We present representative deep learning results demonstrating near-ceiling performance on MSTAR.
3. We contrast these results with preliminary findings on the ATRNet-STAR dataset (Wagner and Palmer 2020), showing significant sensitivity to relatively minor geometric perturbations.

## Overview of MSTAR Dataset

The MSTAR dataset consists of SAR image chips corresponding to a limited set of ground vehicle targets, typically collected at specific depression angles and across discrete azimuth angles. Each image chip exhibits several defining characteristics:

- Target-centered imagery: The object of interest is consistently located near the center of the chip.
- Fixed spatial dimensions: All image chips share the same pixel resolution and extent.
- Minimal clutter: Backgrounds are relatively homogeneous, with little environmental variation.

- Single-object scenes: Each chip contains exactly one labeled target.

From a machine learning perspective, these properties substantially reduce intra-class variability while amplifying inter-class separability. In effect, the dataset implicitly encodes strong spatial priors that deep networks can exploit without learning robust, physically grounded features. These characteristics introduce significant limitations when MSTAR is used as the sole benchmark for modern deep learning based ATR systems.

### Deep Learning Performance on MSTAR

To illustrate the implications of these dataset properties, we trained a convolutional neural network using standard supervised learning procedures on the MSTAR dataset. The model architecture and training pipeline are consistent with common practices in the literature and are documented in the accompanying Jupyter notebook. Table 1 summarizes representative classification performance on a held out MSTAR test set. The overall precision, recall and f1-scores are all high. Further, these results are consistent with many previously reported studies and confirm that modern CNNs can easily separate MSTAR target classes under standard evaluation protocols.

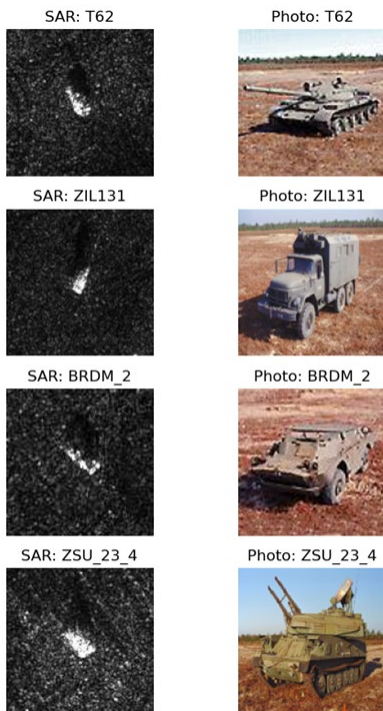


Figure 1: Example of MSTAR SAR image chips with reference photographs.

Class	Precision	Recall	F1-Score
2S1	0.98	0.96	0.97
BRDM-2	0.96	0.98	0.97
D7	0.72	1.00	0.84
SLICY	1.00	1.00	1.00
T62	0.98	0.89	0.93
ZIL131	0.98	0.87	0.92
ZSU-23-4	1.00	0.82	0.90
<b>ACCURACY</b>			<b>0.93</b>
<b>AVERAGE</b>	<b>0.95</b>	<b>0.93</b>	<b>0.93</b>

Table 1: Representative CNN Performance on MSTAR dataset

Beyond scalar metrics, we examined internal feature embeddings extracted from the penultimate layer of the network. Preliminary analysis suggests that embeddings corresponding to different target classes form tight, well-separated clusters when restricted to a narrow range of azimuth angles. The observed structure further supports the conclusion that MSTAR presents a low-complexity classification problem for deep models.

### Sensitivity Analysis Using ATRNet-STAR

To examine model robustness, we evaluated the same trained network on data derived from the ATRNet-STAR dataset, which introduces controlled but operationally relevant variations. Key differences relative to MSTAR include:

- Targets offset from the center of the image chip
- Target rotations relative to the sensor
- Increased geometric variability without changing target identity.

Class	Precision	Recall	F1-Score
2S1	0.28	0.27	0.28
BRDM-2	0.48	0.75	0.59
D7	0.29	0.23	0.25
T62	0.32	0.36	0.34
T72	0.57	0.6	0.58
ZIL131	0.29	0.25	0.27
ZSU-23-4	0.33	0.46	0.47
BMP-2	0.49	0.46	0.47
BTR-60	0.2	0.17	0.18
BTR-70	0.41	0.13	0.19
<b>ACCURACY</b>			<b>0.41</b>
<b>AVERAGE</b>	<b>0.37</b>	<b>0.37</b>	<b>0.36</b>

Table 2: CNN Performance on ATRNet-STAR.

Table 2 summarizes preliminary performance results on ATRNet-STAR. Despite representing the same target classes, these modest perturbations result in substantial degradation. This finding highlights a key risk that models that appear highly capable under MSTAR evaluation may in fact be brittle when exposed to even mild departures from the dataset’s implicit assumptions.

### Embedding-Space Structure for a Narrow Azimuth Slice (170°–190°)

To further characterize the internal representations learned by the network, we performed an embedding-space analysis restricted to a narrow azimuth slice (170°–190°). Embeddings were extracted from the dense penultimate layer of the trained CNN and analyzed using principal component analysis (PCA), inter-/intra-class scatter metrics, and the Bhattacharyya distance to quantify class separability.

Figure 2 (PCA visualization for 170°–190°) shows the first two principal components of the embedding vectors. Even under this restricted viewing geometry, the learned representations exhibit strong clustering behavior: samples from each target class form compact, tightly grouped clusters, with clear separation between most classes in the low-dimensional projection. This qualitative structure supports the hypothesis that, within a constrained azimuth regime, MSTAR presents a highly separable classification problem for modern deep networks.

Quantitatively, the Fisher Discriminant Ratio (FDR = 1.2856) indicates that between-class scatter exceeds within-class scatter in the learned feature space. The Davies–Bouldin Index (DBI = 1.3729) and Dunn Index (0.3134) further confirm that clusters are generally compact and well-separated, though not perfectly isolated. These values are consistent with the visually apparent separation in the PCA projection.

The Bhattacharyya analysis provides a more granular view of pairwise class confusability. Most class pairs exhibit large Bhattacharyya distances, indicating strong statistical separability. For example, pairs such as D7–SLICY and SLICY–ZIL131 show very high separation values of 575.9 and 5416.3, respectively. However, the smallest Bhattacharyya values occur for the T62–ZIL131 (226.9) and 2S1–T62 (231.7) pairs, indicating relatively higher overlap in feature space. This is consistent with the PCA visualization, where T62 and ZIL131 display a modest degree of cluster proximity and limited overlap relative to the other classes.

Importantly, even this most confusable pair remains largely separable in the full embedding space, suggesting that the overlap observed in two-dimensional PCA projections reflects dimensionality reduction effects rather than fundamental inseparability. Nevertheless, the T62–ZIL131

interaction highlights that residual ambiguity persists even in the otherwise simplified MSTAR setting.

Collectively, these azimuth-sliced embedding results reinforce a central claim of this paper: when target pose is narrowly constrained, MSTAR induces highly structured and well-separated feature manifolds that make the classification task comparatively easy for deep networks. The tight intra-class grouping and strong inter-class separation observed here provides further evidence that high accuracy on MSTAR, particularly within restricted azimuth regimes, should be interpreted as performance on a low-complexity, bias-structured dataset rather than as proof of operational robustness.

### Discussion

Recent advances in machine learning have highlighted the risks of dataset bias and shortcut learning (Torralba and Efros 2011; Geirhos et al. 2020), where models achieve strong performance by exploiting spurious correlations or structural regularities rather than learning task-relevant, invariant features. In computer vision, this phenomenon has been observed in datasets where object position, background statistics, or acquisition geometry are inadvertently correlated with class labels.

The MSTAR dataset exhibits several characteristics that make it particularly susceptible to shortcut learning. Consistent target centering, fixed chip dimensions, limited background clutter, and single-object scenes create strong spatial and contextual priors. Deep networks trained on MSTAR can therefore achieve high accuracy by relying on positional cues, silhouette regularities, or background homogeneity rather than learning robust scattering-based or geometry-invariant representations of target structure.

From this perspective, MSTAR performance should be interpreted cautiously. High accuracy does not necessarily indicate that a model has learned physically meaningful SAR features, but rather that it has successfully internalized the dataset’s implicit biases. This framing helps explain the sharp performance degradation observed when the same models are evaluated on datasets such as ATRNet-STAR, where these priors are intentionally violated through target offset and rotation. By situating MSTAR within the broader context of dataset bias and shortcut learning (Torralba and Efros 2011; Geirhos et al. 2020), we argue that it serves best as a controlled sandbox for algorithm development, not as a stand-alone benchmark for assessing operational readiness or generalization.

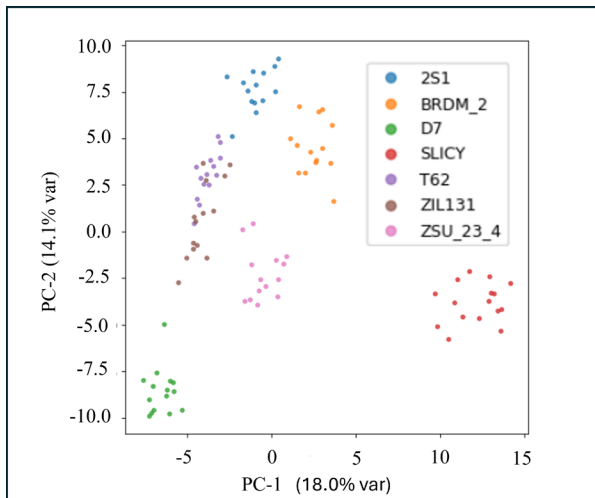


Figure 2: First two principal components of the embedding vectors.

The results presented here reinforce a critical point: MSTAR is an excellent dataset for learning how to build SAR classifiers, but a poor dataset for validating operational robustness. Deep networks trained on MSTAR achieve high accuracy not solely because they learn invariant target signatures, but because they exploit dataset regularities that do not generalize.

This does not diminish the historical or pedagogical value of MSTAR. Rather, it suggests that continued reliance on MSTAR as a primary benchmark risk overstating progress in SAR ATR. More realistic datasets featuring clutter, multi-object scenes, variable framing, and broader pose distributions are essential for meaningful evaluation.

## Conclusion & Near-Term Extensions

This paper examined the utility and limitations of the MSTAR dataset for deep learning-based SAR ATR, with particular emphasis on learned representation structure. While modern CNNs achieve near-ceiling classification accuracy on standard MSTAR splits, deeper analysis of embedding-space geometry reveals why this performance is attainable.

Low-dimensional visualizations of dense-layer embeddings, restricted to a narrow azimuth slice ( $170^\circ$ – $190^\circ$ ), show tightly clustered intra-class groupings and strong inter-class separation. Quantitative scatter metrics—including Fisher Discriminant Ratio, Davies–Bouldin Index, Dunn Index, and pairwise Bhattacharyya distances—confirm that between-class variance generally exceeds within-class variance in the learned feature space. With the exception of

modest overlap between T62 and ZIL131, the classes are cleanly separable under constrained pose conditions. These findings reinforce the central argument of this paper: MSTAR induces highly structured, low-complexity decision boundaries that make the classification problem comparatively easy for modern deep networks.

Importantly, this embedding-space structure should not be interpreted as evidence of operational robustness. Rather, it reflects the dataset’s strong priors - consistent target centering, limited clutter, fixed chip dimensions, and restricted geometric variability. When these priors are relaxed, as in preliminary experiments on ATRNet-STAR, performance degrades substantially. Together, the accuracy results, embedding geometry, and cross-dataset sensitivity analysis demonstrate that high performance on MSTAR largely reflects exploitation of dataset regularities rather than fully invariant, physically grounded SAR feature learning. This work remains ongoing. Near-term extensions include:

- Cross-dataset embedding comparison: Evaluating how embedding compactness and separability shift when models are exposed to ATRNet-STAR or synthetically perturbed imagery.
- Pose-sensitivity profiling: Systematically measuring representation drift under controlled translation, rotation, and framing perturbations.

Collectively, these extensions aim to move SAR ATR evaluation beyond scalar accuracy metrics toward representation-aware, robustness-centered assessment. By explicitly analyzing learned feature geometry and its sensitivity to controlled perturbations, we seek to establish more meaningful criteria for judging progress in SAR-based automatic target recognition.

## References

- Ross, T. D.; Worrell, S. W.; Velten, V. J.; Mossing, J. C.; and Bryant, M. L. 1998. Standard SAR ATR Evaluation Experiments Using the MSTAR Public Release Data Set. *Lincoln Laboratory Journal* 11(2): 161–180.
- Keydel, E. R.; Lee, S. W.; and Moore, J. T. 1996. MSTAR Extended Operating Conditions: A Tutorial. In *Proceedings of SPIE 2757, Algorithms for Synthetic Aperture Radar Imagery III*.
- Chen, S.; Wang, H.; Xu, F.; and Jin, Y. Q. 2016. Target Classification Using the Deep Convolutional Networks for SAR Images. *IEEE Transactions on Geoscience and Remote Sensing* 54(8): 4806–4817.
- Wagner, R.; and Palmer, J. 2020. ATRNet-STAR: A SAR Dataset for Robust Automatic Target Recognition. Air Force Research Laboratory Technical Report.
- Torralba, A.; and Efros, A. A. 2011. Unbiased Look at Dataset Bias. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Geirhos, R.; Jacobsen, J. H.; Michaelis, C.; Zemel, R.; Brendel, W.; Bethge, M.; and Wichmann, F. A. 2020. Shortcut Learning in Deep Neural Networks. *Nature Machine Intelligence* 2: 665–673.