
A fine opportunity to leverage pre-trained models adapted to Solar PV classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Renewable energy such as solar is key to ensuring access to affordable and sus-
2 sustainable energy generation. Surveying its adoption patterns globally is pivotal to
3 measuring and evaluating renewable energy access and creating a more efficient and
4 equitable grid. Leveraging high-resolution imagery to detect solar PVs has proven
5 to be a more exhaustive way of covering all PVs, including residential PVs that can
6 be very challenging to track via conventional surveying methods. While the litera-
7 ture has developed models to classify and segment PV installations, residential PV
8 is still challenging to identify using medium resolution (≈ 30 cm/pixel or above)
9 remote-sensing products. This work explores different fine-tuning (FT) strategies
10 of pre-trained *ViT* models for classification tasks in smaller dataset settings. While
11 FT offers an opportunity for fast and computationally efficient model deployment,
12 practitioners have to be cautious about the effects of fine-tuning on OOD classifi-
13 cation and how advances in text attention mechanisms do not necessarily map to
14 image architectures. Moreover, the LoRA technique (Low-Rank Adaptation) is
15 identified as an efficient method for fine-tuning, enhancing the model’s adaptability
16 to specific tasks while preserving its generalizability. Despite these advancements,
17 achieving robust OOD classification in a foundational model context remains a
18 challenging task.

19 1 Introduction

20 PV adoption deployment and access are crucial for achieving affordable and clean energy access
21 and meeting decarbonization goals across the globe. Understanding the distribution of solar PV
22 allows governments, businesses, and individuals to identify regions with the greatest potential for
23 solar energy generation and efficiently integrate solar PV into the electrical grid. This is particularly
24 relevant as annual growth in the solar industry will average 15% [NREL, 2023]. Hence, there is a
25 need to explore more efficient ways of tracking and monitoring it .

26 Furthermore, PV adoption is not uniform often showcasing underlying race and income inequalities
27 [Sunter et al., 2019, Lukanov and Krieger, 2019, Kwan, 2012]. Measuring adoption helps to alleviate
28 the deepening of “energy privileges” [Stokes et al., 2023] and the misallocation of tax benefits that
29 lead to inequitable adoption of renewable energy. While field surveying is available¹, its time and
30 spatial coverage is often inadequate to capture longitudinal changes in solar adoption, especially in
31 fast-adopting markets like the United States.

32 Existing projects in the literature have built longitudinal PV adoption data with significant spatial
33 and time coverage (i.e. *DeepSolar* [Wang et al., 2022, Yu et al., 2018]). Nonetheless, some of

¹The US Department of Energy’s *OpenPV* is the largest database of PV installation in the US using crowd-sourced data, but it was discontinued in 2019.

34 these projects have some limitations such as high-resolution (≤ 20 cm spatial resolution) data
 35 requirements that are not publicly available. Second, given the sparsity of PV adoption, there are
 36 no country-wide datasets, leading to urban biases [Wang et al., 2022]. Lastly, deploying these
 37 models can be cumbersome and prohibiting for researchers and policy-makers in the developing
 38 world, where computing resources and labels are limited. Although *DeepSolar* [Yu et al., 2018] and
 39 *DeepSolar++* [Wang et al., 2022] models have achieved high performances in tasks such detecting
 40 solar PVs (precision (recall) of 93.1% (87.5%) in residential areas) and predicting installation years
 41 (rate of 93.9 ± 1.0 over a random sample of 23 counties), their coverage is still limited to the US and
 42 dependency on high-resolution satellite/aerial imagery.

43 Fine-tuning pre-trained models for downstream tasks has demonstrated superior performance when
 44 compared to training from scratch in the context of language models. This approach has emerged as
 45 the prevailing and widely accepted strategy for addressing downstream classification and generation
 46 tasks within the domain of language modeling, as supported by relevant scholarly works [Wei
 47 et al., 2021, Zhang and Bowman, 2018]. Recent advances in fine-tuning, such as *LoRa* [Hu et al.,
 48 2021], have also streamlined the fine-tuning process by reducing the parameter space allowing
 49 faster inferences and task adaptation. This, combined with different optimization techniques such as
 50 contrastive learning [Chen et al., 2020] and self-supervised learning [Chen et al., 2021], have boosted
 51 accuracy on different vision classification benchmarks: a ResNet-50 pre-trained with ImageNet,
 52 improves CIFAR-10 classification from 95% to 98% [Chen et al., 2020]. Fine-tuning also alleviates
 53 some of the financial and environmental limitations related to training from scratch. Patterson et al.
 54 [2021], Schwartz et al. [2019]. These are particularly relevant in the developing world where access
 55 to GPU computing is limited and cloud-based options can be burdensome. Fine-tuning can achieve
 56 competing performances with fewer training labels and shorter computation times.

57 Intuitively, fine-tuning all the layers of a neural network can adapt a pre-trained model faster to a new
 58 task and obtain better results than training only a few layers or reducing its training parameters' feature
 59 space. Nonetheless, previous work in language models (BERT [Devlin et al., 2019] and RoBERTa
 60 [Liu et al., 2019]) have shown that only a fourth of the last layers are required to keep similar levels of
 61 accuracy [Lee et al., 2019], other experiments have shown that linear-probing and freezing might be
 62 better alternatives to naïve transfer and Out-of Distribution (OOD) performance [Kumar et al., 2022]
 63 in contrastive model settings. Nonetheless, there is no evidence of recommendations on supervised
 64 fine-tuning examples and in OOD settings where data sources vary in terms of resolution and color
 65 space. Moreover, while *ViT* architectures have brought new gains, some of the previous rules for
 66 fine-tuning and transfer learning cannot be adopted from the CNN architectures [Chen et al., 2021].

67 In this paper, we test the ability of different fine-tuning
 68 strategies to adapt pre-trained vision transformer (*ViT*)
 69 models to a PV classification downstream task. Addition-
 70 ally, we also evaluate our best models for their ability to
 71 generalize in a OOD setting – geographical and spatial
 72 resolution domains. For this, we use two solar PV datasets
 73 in two areas of interest (AoI), China and the state of Cal-
 74 ifornia (US), in different resolutions. While some recent
 75 literature [He et al., 2022, Goyal et al., 2022] have done
 76 similar experiments to the ones we present in this paper,
 77 this work inscribes into an open problem in Earth Obser-
 78 vation (EO) and uses an applied example that differs in
 79 complexity from vision benchmarks in the literature. We
 80 show that (...)

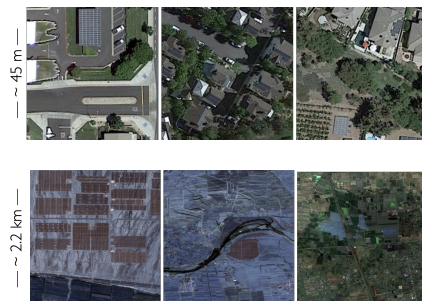


Figure 1: Labeled datasets for PV classification. Both image sets were collected around the same time: 2017 - 2019. The first high-resolution set corresponds to the Google Maps aerial imagery. The second set, Sentinel-2, has 10 times coarser resolution, but a higher revisit time, (better representation of the changes over time). The Sentinel-2 data is aggregated in time and cleaned to remove “bad” pixels and clouds. We have 4179 positive labels and 3966 negative labels across datasets.

81 2 Methodology

82 2.1 Data

83 Experiments will use two in-house remote-sensing
 84 datasets from different resolutions. On one hand, we col-
 85 lected GoogleMaps aerial data imagery for our two AoI,
 86 California (CA), and several Northern provinces of China
 87 (CN). These images have an approximate resolution of 0.2

88 m/pixel with three visible channels (RGB). On the other
89 hand, we collected Sentinel-2 scenes with 10 m/pixel res-
90 olution (50 times coarser than the Google images) for the
91 same locations in CN. While the Google imagery is static
92 in time (latest), the Sentinel-2 data has an average revisit duration of 8 days for our AoI, thus we
93 build a cloud-free median composite image by combining all the 2,560 images collected from 2017
94 to 2019. We adopt '80-20-20' split policy for training, validation and testing sets.

95 The PVs labels for commercial and utility-scale installations come from two sources. For Califor-
96 nia, we used data from the *DeepSolar* database [Yu et al., 2018]. For China, we used data from
97 Kruitwagen et al. [2021], a novel database with more than 60,000 global PV installations. For
98 both data sources, we extracted an image patch of 224×224 pixels around each of the labeled
99 points and fed it to one of the most widely used *ViT* for image classification model pre-trained on
100 `google/vit-base-patch16-224` that uses *ImageNet-21K*.

101 2.2 Experiments

102 To test the best fine-tuning strategy for adaptation to the PV classification downstream task, we run
103 the following experiments:

- 104 1. **Full network retraining (FT)**: The most intuitive way of fine-tuning is retraining all the
105 layers of the model with the downstream dataset. It is well known that this approach leads
106 to better in-distribution accuracies. However, for OOD datasets when the shift is large, these
107 may not perform very well. We use this as a benchmark to compare against other fine-tuning
108 strategies. [~ 85 M parameters]
- 109 2. **Layer freezing (L_k)**: Freezing the first or last layers has become a common practice for
110 fine-tuning in CNN vision models and other language tasks now. In our case, we freeze
111 everything but the first two transformer blocks (L_2) for one set of experiments and similarly
112 for the last two blocks (L_n). [~ 14 M parameters]
- 113 3. **Linear probing (LP)**: Following some of the literature Kumar et al. [2022], Chen et al.
114 [2021] that posits larger gains from linear probing over fine-tuning in the presence of
115 distribution shifts, we perform fine-tuning only the last MLP linear layer of each of the
116 attention heads. [~ 28 M parameters]
- 117 4. **Low-rank Matrix Factorization (LoRA)**: The above fine-tuning methods still pose compu-
118 tational challenges. Approaches in NLP have shown that large-scale pre-trained models used
119 for fine-tuning on different tasks rely on a small intrinsic dimension. In particular, we used
120 *LoRA* Hu et al. [2021], which uses a low-rank decomposition where gradient updates are
121 represented by $h = W_0x + \Delta W_x = W_0 + BAx$, where B and A are matrices in a reduced
122 matrix space: $\mathbb{R}^{d \times r} \times \mathbb{R}^{r \times k}$, which can project the full original parameter space $\mathbb{R}^{d \times k}$. We
123 use rank 4 for our experiments. We expect that we can achieve comparable performance by
124 reducing the trainable parameters by more than 99%. [~ 300 K parameters]

125 Each strategy will use a *ViT* (`google/vit-base-patch16-224`) with an Adam optimizer with the
126 learning rate 1×10^{-5} with a linear-schedule and weight decay (L2 regularization) 1×10^{-3} to
127 minimize over-fitting. All models are trained for 20 epochs.

128 3 Results

129 To compare the experiment results, we use a fine-tuned ResNet50 (pre-trained on *ImageNet-21K*)
130 with our dataset as a baseline. This CNN achieved a F1 score of 75%. However, the score improved
131 by 20% with the *ViT* model. The initial set of experiments indicated some overfitting so we added
132 regularization to minimize that.

133 We summarize the experiment results in Table 1. We observe that *LoRA* generally performs better
134 with the best trade-off in terms of computational cost. The F1 scores for high-resolution datasets are
135 higher than Sentinel as expected. The performance on the California HR dataset is lower than that
136 of China's HR dataset even though the former has more samples. The reason is California's dataset
137 includes commercial and utility-scale PVs that are harder to detect as compared to China's dataset
138 which has only utility-scale (commercial-scale PVs are smaller in size than utility-scale PVs).

139 Another interesting observation is the compar-
 140 ison of L_1 and L_n for China HR and S2. L_1
 141 for China S2 performs better as opposed to high
 142 resolution datasets. This is because the data dis-
 143 tribution domain of Sentinel 2 is significantly
 144 different than the high-resolution aerial imagery.

145 3.1 OOD evaluation

146 We also evaluated our best models on OOD
 147 datasets. The results are shown in Table 1.
 148 The model trained on China Sentinel 2 dataset
 149 (the best being LoRA) is used to test the other
 150 two datasets (China HR and Cal HR) and so
 151 on. China S2 trained model performs better on
 152 China HR than California HR by $\sim 30\%$ F1
 153 score. This is because the PVs are more simi-
 154 lar in China HR and S2 than PVs in China and
 155 California.

156 4 Discussion

157 We have analyzed various fine-tuning strategies
 158 for *ViT* architectures and found that full param-
 159 eter retraining is often not required to achieve
 160 baseline performance. We also find that dif-
 161 fering datasets and image characteristics (*i.e.*
 162 resolution, color space, etc.) call for different
 163 transfer learning methods, with *LoRA* being the
 164 best-performing strategy across datasets in our experiments. These results line up with some of the
 165 findings in the literature [Kumar et al., 2022, Chen et al., 2021]. As suggested by other papers [Goyal
 166 et al., 2022, Raghunathan et al., 2020], it is also relevant to think in terms of OOD examples and how
 167 different optimization processes between pre-training and fine-tuning can lead to different results
 168 (i.e. contrastive loss in pre-training, and cross-entropy during fine-tuning leads to sub-optimal results
 169 Goyal et al. [2022]). We would like to work towards improving the OOD evaluations as this can
 170 help us translate our PV classification parameters across different geographical regions, a especially
 171 crucial in the sustainability domain.

172 While other works have explored the use of few-shot meta-learners to achieve OOD performance in
 173 remote-sensing scenarios [Wang et al., 2020], this works differs from those approaches by caring
 174 solely about domain adaptation and wanting to explore a streamlined way of fine-tuning strategy
 175 for a common problem in PV detection. In this work, we have posed the opportunity of FT as a
 176 computationally simpler and less data-greedy alternative to training from scratch. We have shown
 177 that we can achieve comparable performance to our baselines, with a comparably smaller set of labels
 178 (6,200 labels). Despite our experiments combining different resolutions, we find that detecting small
 179 PV installations is still challenging, and fine-tuned pre-trained models are still not able to outperform
 180 models trained exclusively on high-resolution imagery.

	FT	L_2	L_n	<i>LoRA</i>	<i>LP</i>
CA [HR]	0.90	0.84	0.89	0.89	0.92
CN [HR]	0.96	0.95	0.96	0.97	0.95
CN [S2]	0.83	0.84	0.79	0.85	0.75

Table 1: F1 Scores of fine-tuning experiments. Each column corresponds to a different fine-tuning strategy: FT (full fine-tuning), L_2 (First two attention blocks), L_n (Last two attention blocks), *LoRA* (Low-rank adaptation), and *LP* Linear Probing.

Trained	OOD	Acc	F1
CN-S2 [LoRa]	CN HR	0.64	0.75
CN-S2 [LoRa]	CA HR	0.46	0.46
CN-HR [LoRa]	CN S2	0.48	0.29
CN-HR [LoRa]	CA HR	0.51	0.33
CA-HR [Linear]	CN S2	0.44	0.17
CA-HR [Linear]	CN HR	0.51	0.37

Table 2: OOD inference performance for different fine-tuned models. Each row corresponds to a different combination of fine-tuned model and OOD dataset.

181 References

- 182 T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A Simple Framework for Contrastive Learning
183 of Visual Representations. In *Proceedings of the 37th International Conference on Machine*
184 *Learning*, pages 1597–1607. PMLR, Nov. 2020. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v119/chen20j.html)
185 [v119/chen20j.html](https://proceedings.mlr.press/v119/chen20j.html). ISSN: 2640-3498.
- 186 X. Chen, S. Xie, and K. He. An Empirical Study of Training Self-Supervised Vision Transformers,
187 Aug. 2021. URL <http://arxiv.org/abs/2104.02057>. arXiv:2104.02057 [cs].
- 188 J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional
189 Transformers for Language Understanding, May 2019. URL [http://arxiv.org/abs/1810.](http://arxiv.org/abs/1810.04805)
190 [04805](http://arxiv.org/abs/1810.04805). arXiv:1810.04805 [cs].
- 191 S. Goyal, A. Kumar, S. Garg, Z. Kolter, and A. Raghunathan. Finetune like you pretrain: Improved
192 finetuning of zero-shot vision models, Dec. 2022. URL <http://arxiv.org/abs/2212.00638>.
193 arXiv:2212.00638 [cs].
- 194 X. He, C. Li, P. Zhang, J. Yang, and X. E. Wang. Parameter-efficient Model Adaptation for Vision
195 Transformers, Dec. 2022. URL <http://arxiv.org/abs/2203.16329>. arXiv:2203.16329 [cs].
- 196 E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-Rank
197 Adaptation of Large Language Models, Oct. 2021. URL <http://arxiv.org/abs/2106.09685>.
198 arXiv:2106.09685 [cs].
- 199 L. Kruitwagen, K. T. Story, J. Friedrich, L. Byers, S. Skillman, and C. Hepburn. A global inventory
200 of photovoltaic solar energy generating units. *Nature*, 598(7882):604–610, Oct. 2021. ISSN
201 1476-4687. doi: 10.1038/s41586-021-03957-7. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41586-021-03957-7)
202 [s41586-021-03957-7](https://www.nature.com/articles/s41586-021-03957-7). Number: 7882 Publisher: Nature Publishing Group.
- 203 A. Kumar, A. Raghunathan, R. Jones, T. Ma, and P. Liang. Fine-Tuning can Distort Pretrained Features
204 and Underperform Out-of-Distribution, Feb. 2022. URL <http://arxiv.org/abs/2202.10054>.
205 arXiv:2202.10054 [cs].
- 206 C. L. Kwan. Influence of local environmental, social, economic and political variables on the spatial
207 distribution of residential solar PV arrays across the United States. *Energy Policy*, 47:332–344, Aug.
208 2012. ISSN 0301-4215. doi: 10.1016/j.enpol.2012.04.074. URL [https://www.sciencedirect.](https://www.sciencedirect.com/science/article/pii/S0301421512003795)
209 [com/science/article/pii/S0301421512003795](https://www.sciencedirect.com/science/article/pii/S0301421512003795).
- 210 J. Lee, R. Tang, and J. Lin. What Would Elsa Do? Freezing Layers During Transformer Fine-Tuning,
211 Nov. 2019. URL <http://arxiv.org/abs/1911.03090>. arXiv:1911.03090 [cs].
- 212 Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and
213 V. Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach, July 2019. URL
214 <http://arxiv.org/abs/1907.11692>. arXiv:1907.11692 [cs].
- 215 B. R. Lukanov and E. M. Krieger. Distributed solar and environmental justice: Exploring the
216 demographic and socio-economic trends of residential PV adoption in California. *Energy Policy*,
217 134:110935, Nov. 2019. ISSN 0301-4215. doi: 10.1016/j.enpol.2019.110935. URL <https://www.sciencedirect.com/science/article/pii/S0301421519305221>.
- 219 NREL. Sumer 2023 Solar Industry Update. Technical report, National Renewable Energy Laboratory,
220 Aug. 2023. URL <https://www.nrel.gov/docs/fy23osti/87189.pdf>.
- 221 D. Patterson, J. Gonzalez, Q. Le, C. Liang, L.-M. Munguia, D. Rothchild, D. So, M. Texier, and
222 J. Dean. Carbon Emissions and Large Neural Network Training, Apr. 2021. URL [http://arxiv.](http://arxiv.org/abs/2104.10350)
223 [org/abs/2104.10350](http://arxiv.org/abs/2104.10350). arXiv:2104.10350 [cs].
- 224 A. Raghunathan, S. M. Xie, F. Yang, J. Duchi, and P. Liang. Understanding and Mitigating the
225 Tradeoff Between Robustness and Accuracy, July 2020. URL [http://arxiv.org/abs/2002.](http://arxiv.org/abs/2002.10716)
226 [10716](http://arxiv.org/abs/2002.10716). arXiv:2002.10716 [cs, stat].
- 227 R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni. Green AI, Aug. 2019. URL [http://arxiv.](http://arxiv.org/abs/1907.10597)
228 [org/abs/1907.10597](http://arxiv.org/abs/1907.10597). arXiv:1907.10597 [cs, stat].

- 229 L. C. Stokes, E. Franzblau, J. R. Lovering, and C. Miljanich. Prevalence and predictors of wind
230 energy opposition in North America. *Proceedings of the National Academy of Sciences*, 120(40):
231 e2302313120, Oct. 2023. doi: 10.1073/pnas.2302313120. URL [https://www.pnas.org/doi/](https://www.pnas.org/doi/abs/10.1073/pnas.2302313120)
232 [abs/10.1073/pnas.2302313120](https://www.pnas.org/doi/abs/10.1073/pnas.2302313120). Company: National Academy of Sciences Distributor: National
233 Academy of Sciences Institution: National Academy of Sciences Label: National Academy
234 of Sciences Publisher: Proceedings of the National Academy of Sciences.
- 235 D. A. Sunter, S. Castellanos, and D. M. Kammen. Disparities in rooftop photovoltaics deployment
236 in the United States by race and ethnicity. *Nature Sustainability*, 2(1):71–76, Jan. 2019. ISSN
237 2398-9629. doi: 10.1038/s41893-018-0204-z. URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41893-018-0204-z)
238 [s41893-018-0204-z/](https://www.nature.com/articles/s41893-018-0204-z). Number: 1 Publisher: Nature Publishing Group.
- 239 S. Wang, M. Rußwurm, M. Körner, and D. B. Lobell. Meta-Learning For Few-Shot Time Series
240 Classification. In *IGARSS 2020 - 2020 IEEE International Geoscience and Remote Sensing*
241 *Symposium*, pages 7041–7044, Sept. 2020. doi: 10.1109/IGARSS39084.2020.9441016. URL
242 <https://ieeexplore.ieee.org/document/9441016>. ISSN: 2153-7003.
- 243 Z. Wang, M.-L. Arlt, C. Zanocco, A. Majumdar, and R. Rajagopal. DeepSolar++: Understanding
244 residential solar adoption trajectories with computer vision and technology diffusion models. *Joule*,
245 6(11):2611–2625, Nov. 2022. ISSN 2542-4785, 2542-4351. doi: 10.1016/j.joule.2022.09.011. URL
246 [https://www.cell.com/joule/abstract/S2542-4351\(22\)00477-9](https://www.cell.com/joule/abstract/S2542-4351(22)00477-9). Publisher: Elsevier.
- 247 C. Wei, S. M. Xie, and T. Ma. Why Do Pretrained Language Models Help in Downstream Tasks?
248 An Analysis of Head and Prompt Tuning. Nov. 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=MDMV2SxCboX)
249 [id=MDMV2SxCboX](https://openreview.net/forum?id=MDMV2SxCboX).
- 250 J. Yu, Z. Wang, A. Majumdar, and R. Rajagopal. DeepSolar: A Machine Learning Framework
251 to Efficiently Construct a Solar Deployment Database in the United States. *Joule*, 2(12):2605–
252 2617, Dec. 2018. ISSN 2542-4351. doi: 10.1016/j.joule.2018.11.021. URL [https://www.](https://www.sciencedirect.com/science/article/pii/S2542435118305701)
253 [sciencedirect.com/science/article/pii/S2542435118305701](https://www.sciencedirect.com/science/article/pii/S2542435118305701).
- 254 K. Zhang and S. Bowman. Language Modeling Teaches You More than Translation Does: Lessons
255 Learned Through Auxiliary Syntactic Task Analysis. In *Proceedings of the 2018 EMNLP Workshop*
256 *BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 359–361, Brussels,
257 Belgium, Nov. 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5448.
258 URL <https://aclanthology.org/W18-5448>.