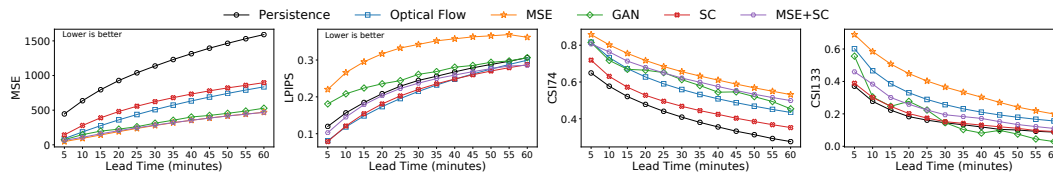


1 Thank you to each reviewer for your helpful feedback on our paper. All comments were taken under consideration and
2 will be used to improve the work. Below we provide our reasoning for several selected points.

3 **Lack of novelty / Use of existing methods** While it might come off as an engineering task, we feel the design and
4 execution of SEVIR required several novel ideas and insights, including recognition of a gap in ML-ready weather
5 datasets, application of the NOAA storm event database for down selection, combining/structuring/formatting multiple
6 data sources together in a way that makes it relevant to a number of tasks in the intersection of meteorology and machine
7 learning. Also, the use of existing methods and metrics in the second half of the paper was an intentional choice by the
8 authors; the goals of the paper are to introduce a new dataset to the community, and to provide simple baselines that
9 are standard, easy to understand, and are reproducible. Inclusion of more complicated or novel architectures/methods
10 would have put us beyond that scope, and moreover would have been unfair to the researchers developing those methods
11 since those results are worthy of their own paper on the topic. Contrary to prior work in this area, we evaluate our
12 models using not only meteorology specific metrics but also with metrics that quantify perceptual quality. This also
13 highlights the fact that a single metric is not always sufficient for evaluating models where image quality perceived by
14 the viewer is also important.

15 **Lack of Benchmarks** Multiple reviewers commented on the weakness of the benchmark used for the nowcasting
16 task, and Reviewer 2 was correct to point out several existing nowcasting methods that are in use that go beyond
17 the "do nothing" persistence model. The choice of only using the persistence benchmark came down to two main
18 reasons (1) Despite it's simplicity, it is a common benchmark used for nowcasting and (2) It is hyperparameter free
19 – unlike other methods for this task (e.g. optical flow). Having said that, based on this reviewer feedback, we have
20 included an optical-flow based benchmark using on SEVIR and intend to include it as a additional baseline if the
21 paper is accepted for publication. Due to page limits, only a portion of the updated figure is shown below. Since we
22 want the focus of the paper to be the utility of the dataset for AI/ML in meteorological applications, we decided an in
23 depth and fair comparison to additional methods (e.g. TrajGRU) would be out of scope (and well over page count).
24 Additionally, evaluation, hyperparameter and architectural tuning of ConvLSTM is part of an ongoing effort and an
25 in-depth comparison of Nowcast models is currently under way.



26 **Discussion surrounding results, metrics and future work** Several reviewers noted weak discussion surrounding
27 metrics and evaluation of the predictions in context of meteorology. This is a fair criticism, and thankfully, is something
28 that could be remedied in the final version (given the extra page allotted). Evaluation metrics for weather prediction is a
29 rich area since the task is multi-objective – in some cases the placement and intensity of weather is most important, and
30 in others the realism and dynamics of the storm are more important. This motivated the range of different losses and
31 objectives, some that stress placement/intensity (MSE/POD/CSI), and others that focus on realism (GAN-based/LPIP).
32 The baselines we provide show that depending on your choice of loss function, certain axes of "goodness" are brought
33 out more than others. We will add more discussion along these lines which address "what is done and why". Also,
34 Reviewer 1 pointed out, our conclusion is simply a summary. We can use the extra space to discuss future and ongoing
35 work, which include a planned AI challenge in the area of radar/satellite nowcasting, applications of SEVIR to few-shot
36 learning that aims to address one of Reviewer 4's questions ("The transferability of the models..."), and approaches that
37 improve robustness in ML models applied to weather.

38 **Clarity / Use of acronyms** For evaluation metrics, we used naming conventions commonly used in the meteorological
39 community (POD vs recall, FAR vs precision, CSI vs IOU). In hindsight, perhaps that wasn't a good choice for NeurIPS.
40 Following Reviewer 2's suggestion, the acronym problem can be alleviated by switching to more standard terminology.

41 *To Reviewer 2:* Agreed about VGG, it's actually surprising how well VGG works in this setting despite being trained on
42 completely different data! Also, we can provide metrics in the actual units of the fields. Our reason for using the 0-255
43 version of radar was to simplify the implementation by not requiring additional post processing, but that's easy to add.

44 *To Reviewer 4:* You're right to point out the word "event" is a bit overloaded, but the NOAA storm events (which are
45 clustered into "episodes") can be linked to particular pixels and frames of the SEVIR events using data in the catalog (so
46 NOT the whole 384² km region). Regarding loss functions: this is very much an open question. Most work uses some
47 L^p loss as it's the simplest, but to your point, there may be better ones. We started to explore that question here with the
48 adversarial and texture losses. Our vision is to spur innovation in this area not only through model architectures, but
49 also new loss functions that are able to simultaneously model weather placement/intensity as well as realism, through a
50 future AI challenge that leverages SEVIR and these baseline implementations.