Table 1: Requested additional comparisons. SAT: Satellite, LAQN: Ground stations. Random seed of 0.

| Model | Data Sources | sRMSE ($\mu \pm \sigma$) | RMSE ($\mu \pm \sigma$) | NLPL ($\mu \pm \sigma$) |
|---|---|---|---|---|
| Single GP | LAQN only | $1.04 \pm 0.04$ | $23.02 \pm 11.26$ | $12.0 \pm 12.22$ |
| MR-GP | SAT only | $0.72 \pm 0.41$ | $14.87 \pm 9.34$ | $16.7 \pm 23.14$ |
| VBAgg-Normal | LAQN & SAT | $0.82 \pm 0.48$ | $16.24 \pm 9.15$ | $9.78 \pm 11.97$ |
| MR-GPRN w/o CL | LAQN & SAT | $0.69 \pm 0.43$ | $14.03 \pm 8.93$ | $9.24 \pm 14.35$ |
| MR-GPRN w/ CL | LAQN & SAT | $0.69 \pm 0.42$ | $14.45 \pm 9.09$ | $8.83 \pm 12.92$ |
| MR-DGP | LAQN & SAT | $\mathbf{0.39 \pm 0.13}$ | $\mathbf{8.65 \pm 4.93}$ | $\mathbf{4.54 \pm 4.12}$ |

1 We thank the reviewers for their time and detailed, constructive feedback. We are glad to see our application-motivated
2 methodological contributions and narrative were well-received. We denote with e.g. **R1.2.3** our response to Reviewer
3 1, Section 2, Paragraph 3. (**Joint**): As requested we are offering additional baselines (Table 1) that strengthen our
4 results and are discussed below. As suggested by **R3** we will move inducing point material to the appendix. The
5 additional page will allow us to improve clarity: we will expand on the MR-DGP model, we will add the additional
6 baselines, suggested by **R1** and **R4**, with further discussions of results and the uncertainty quantification benefits of the
7 CL corrections. We will also improve the description of the experiments and lighten the use of inline equations.

8 **R1.2.4**: We will improve the motivation for the composite likelihood (CL). The estimated CL ensures that the asymptotic
9 posterior $p(\mathbf{Y}|\mathbf{f}, \mathbf{X}, \theta)$ converges to the misspecified asymptotic MLE distribution [24, 30]. The CL cannot be set as a
10 free parameter because otherwise we would not obtain this theoretical guarantee. The MR-DGP learns dependencies
11 between the layers $p(\mathbf{f}_1|\mathbf{f}_2, \ldots)$ and hence between resolutions. **R1.2.5**: We have rerun all our experiments and
12 additional requested baselines on the real world example, see [**Joint**] above. **R1.2.7** The size of the std is due to
13 variability across the 42 sites in London, we will also offer site-standardized results (e.g. sRMSE) in the appendix. We
14 are happy to follow alternative standardizations if reviewers express a preference for the final version.

15 **R3.2.3**: We agree that the distinction between multi-fidelity and multi-resolution would be beneficial and we will offer
16 that. See [**Joint**] above to see how we will improve the explanation of MR-DGP (including the choice of kernel).
17 MR-DGP arises very naturally in real world examples and is able to successfully handle data from biased sensor
18 networks; as shown in our experiments the performance is significant. **R3.2.10-11**: We will reduce the number of inline
19 equations by moving the standard results into the appendix. **R3.2.13**: This is an interesting paper that simply takes the
20 formulation of [27] and applies it to the multi-task setting through the LCM formulation. The model proposed by the
21 authors is a special case of MR-GPRN where the latent GPs $\mathbf{W}$ are constant. We have also presented a very natural
22 and principled method to dealing with bias whereas they consider a very ad-hoc solution through data normalization.
23 **R3.5.1**: It is indeed natural that specific contributions will be more or less interesting to different readers. In this paper
24 we have tackled some of the underlying issues of previous approaches and offer the state-of-the-art.

25 **R4.1.1**: We respectfully disagree with "generalization of multi-task ... straightforward", in fact we have challenged
26 two very common assumptions by correcting (MR-GPRN) or accounting (MR-DGP) for dependent observations and
27 have provided a principled way to deal with biases and multi-resolution. Through these contributions we have shown
28 impressive results on a very complex problem. **R4.1.2 + R4.2.Originality**: Both are latent variable models and with
29 MR-DGP we use the latent structure to model the obs.dependency instead of correcting via CL. Further extensions
30 or special cases of this framework we leave for future work. **R4.2.3**: Indeed the composite weight comes out of the
31 expectation (See Eqn 14 in appendix). We do not model the cross-resolution dependencies via CL in the MR-GPRN
32 model, the weight corrects for posterior contraction due to loss of that dependency as done in similar settings [24]. We
33 will further demonstrate the UQ benefits of CL in the extra page through coverage and pred.densities. (**R4.2.4**) VBAgg
34 is the only published work that is a suitable baseline. We have shown that both handle the same types of multi-res
35 data in Sec. F of the appendix. As also suggested by **R3** we will merge this into the main text. We do not use [6] as
36 baseline because, despite the name, it is unable to handle multiple observation processes, see ($\ell : 170 - 173$). **R4.2.5**:
37 See [**R1.2.4**] above. **R4.2.6**: The dimension of the Hessian is $|\hat{\theta}| \times |\hat{\theta}|$ where $\theta$ are the hyper parameters. The size is
38 very small and is dominated. We will clarify this in the main text. (**R4.2.Clarity.3**): Thank you for pointing out the
39 typo in bolding MR-DGP, we will fix this. MAPE is an asymmetric loss that penalizes overestimation. As shown in Fig.
40 2 the prediction from MR-DGP is slightly over estimating whereas VBAgg-Normal is severely underestimating, hence
41 MAPE over penalizes MR-DGP. **R4.2.Clarity.4**: Different sensor networks are calibrated differently, hence comparing
42 raw values is not viable. The information theoretic corrections are from the composite weights in Sec. 3, we will clarify.
43 Hyper-local is higher resolution than typical LSOA area estimates. (**R4.2.Clarity.5-7**): Thank you for pointing out
44 these typos, we will fix them. The subscript $m_k$ is meant to be $m_a$ which represents the output for each layer. The $p$ is
45 used to denote the multiple tasks and because we have presented MR-DGP in the general case the ordering of tasks
46 between layers is a user-choice. We will improve clarity throughout the paper based on all reviewers' suggestions.