

1 **R1: The experimental setup** of our paper is designed to both corroborate our theoretical results which are non-trivial,  
2 while also demonstrating the possible application of the hypernetwork induced prior in some practical use cases. We  
3 concede that our experiments are by no means thorough enough to consider the hypernetwork kernels as a go-to  
4 algorithm for image completion tasks, however we feel they do serve their purpose demonstrating the usefulness of  
5 hypernetworks induced priors, explaining the validity and inductive bias of the architecture. **Size of  $f$ :** While in some  
6 of the hypernetworks in the literature  $f$  is not very large, this is the case in many of the recent networks, e.g., [B,C,D,E].  
7 **R2: We accept the suggested terminology** and would use “hypernet” to refer to network  $f$  only. **There are two types**  
8 **of hypernetwork** architectures (including the hypernetwork  $f$  and the primary network  $g$ ). Type A has a much larger  
9 network  $f$  than  $g$ , and the input to  $f$  is larger than the input to  $g$ . Type B has a smaller  $f$  and larger  $g$ . In the references  
10 that the reviewer mentioned, where one optimizes the input of  $f$ , it is more natural to use type B (smaller inputs).  
11 However, type A is at least as prevalent in the literature as type B. Examples of type A include cases where  $g$  is a  
12 single convolutional layer in a deeper network [17,15,5,A]. This is also very prominent in recent work in which the  
13 perceptual task is done by a resnet  $f$  and the solution is parameterized by a small network  $g$  [B,C,D,E]. In such cases,  $f$   
14 observes the entire context, while  $g$  is a local model, see also [F]. Our analysis holds also for type B networks, given  
15 that the network  $f$  is wide enough to be approximated by a GP. However, since type B networks are used to find optimal  
16 hyperparameters, they require training by definition, and NTKs, which are studied at initialization, are less relevant.  
17 Indeed, there is currently no theoretical machinery for understanding the dynamics of optimizing the input with a fixed,  
18 trained network in the NTK regime. Our work here presents a first step at a new understanding of hypernetworks through  
19 the new machinery of NTK, which covers the type A scenario. We hope to extend this analysis to type B scenarios in  
20 the future as well by using the recent Tensor Programs framework. However, this is out of the scope of this contribution.  
21 **The MNIST experiment** is a typical type A hypernet setting. The perceptual input is processed by  $f$ , and  $g$  is a model  
22 of the “scene”. It directly follows [B,E] (E is a paper R4 pointed to as an example for realistic settings). The reviewer’s  
23 suggestion to condition  $f$  on the digit label is equivalent to learning 10 different denoising networks, which is not  
24 utilizing the full power of hypernets. **The computational advantage (L 208)** is meaningful when compared to other  
25 kernel methods. From the composition of the hyperkernel (Eq. 12),  $\Theta^f(x, x')$  can be evaluated separately for all pairs  
26  $x, x'$ , instead of evaluating kernel values for all pair of tuples  $(x, z), (x', z')$  when considering other kernel methods.  
27 When  $f$  is a convnet, this can represent a significant reduction in computational cost. **R3: Our theoretical results are**  
28 **non-trivial.** In particular, Thm. 1 provides the asymptotic behavior of high order NTK terms which hold for ReLU  
29 hypernetworks, as well as regular ReLU MLPs. Our contribution here is both a technical novelty (in the proof) and the  
30 significance of the final result. On the technical level, as noted in remark 1 and in L 280-283, we have proven (and  
31 improved upon) a conjecture on the asymptotic rates of various correlation functions arising in neural network dynamics  
32 (see [5]). As for the result itself, we are the first to arrive at these tight bounds which relate to both hypernetworks and  
33 MLPs. Thms. 2 and 3 describe the conditions in which GP behavior emerges in hypernetworks (again nontrivial), and  
34 describe the composition of the GP and NTK kernels. We feel these theoretical results are of interest to the community.  
35 **The case of a finite  $f$  and an infinite  $g$**  is left for future work. Note that an infinite  $g$  would require  $f$  to output an  
36 infinite number of parameters. We do discuss the case of both  $f$  and  $g$  being infinite in Sec. 4. **Reporting variance** we  
37 regret not reporting error bars, which will be added. The results were averaged over 10 different training and test splits.  
38 **HN outperforms HK in Tab. 2** For small data regimes, the HK outperforms,  
39 while for larger datasets, the HN is better. This is consistent with prior observa-  
40 tions that kernel methods with NTK tend to outperform in low data regimes. In  
41 **Fig. 2**, the input of  $f$  is depicted in Row 2, the input of  $g$  is a pixel coordinate.  
42 The output of  $g$  is the pixel intensity in the corresponding coordinate. **Typos**  
43 We apologize for the typos and would provide more background on NTK and  
44 GP. **R4:** The paper is theoretically driven. **The experiment setup** is very similar  
45 to that of the mentioned paper (arXiv 1902.10404) only done on MNIST.  
46 To demonstrate this, Fig. I has interpolation results similar to that paper by inter-  
47 polating between  $[\Theta^h(u, u_1), \dots, \Theta^h(u, u_N)]$  and  $[\Theta^h(v, u_1), \dots, \Theta^h(v, u_N)]$   
48 for two images  $u$  and  $v$  in Eq. 16. **We regret the lack of details** in the exper-  
49 imental section. The hypernet  $f$  in our setup is a convolutional neural network  
50 operating on sparse images containing the context points, similarly to [G].

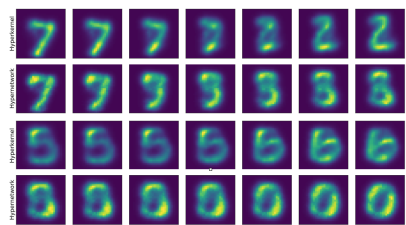


Figure I: Interpolation between 7  $\rightarrow$  5 and 5  $\rightarrow$  6 using hyperkernel (rows 1,3) and hypernetwork (rows 2,4). Both methods used merely 200 samples for training. The hypernetwork trained with sgd clearly underperforms in this low data regime

51 REFERENCES: [A] Wu et al. Pay Less Attention with lightweight and Dy-  
52 namic Convolutions. ICLR 2019. [B] Littwin et al. Deep Meta Functionals for  
53 Shape Representation. ICCV, 2019. [C] Rotman et al. Electric Analog Circuit  
54 Design with Hypernetworks and a Differential Simulator. ICASSP, 2020. [D]  
55 Bergman et al. Implicit neural representations with periodic activation func-  
56 tions. arXiv:2006.09661, 2020. [E] Klocek et al. Hypernetwork functional image representation. ICANN, 2019. [F]  
57 Nachmani et al. Hyper-Graph-Network Decoders for Block Codes. NeurIPS, 2019. [G] Sitzmann et al. Implicit Neural  
58 Representations with Periodic Activation Functions. arXiv:2006.09661, 2020.