1 We thank all the reviewers for their constructive comments. The following are our point-by-point replies.

2 **Rev. 1.**

3 **1)** *this is only for 1HL networks but the Theorems do not specify this*: The presentation in the submission might cause
4 confusion. While Theorems 5, 9 and 10 consider the three-layer networks for simplicity, their extension to $L$-layer
5 networks is easy by replacing the derivatives of the loss with the back-propagated delta. We briefly mentioned this at
6 the beginning of Sec. 3.2. But, we will make a clearer argument in an update.

7 **2)** *No definition of saddle and local minimum points*: We will include their standard definitions in Sec. 2. See also 13)
8 for semi-flat minima.

9 **3)** *Strict or non-strict saddles and what order?*: It is easy to derive the index of the saddle point based on the Hessian in
10 Lemma 4, but the details were skipped by space limitation. If $G$ is positive (or negative) definite and the off-diagonal $F$
11 is of full rank, the index is easily given by Eq.(37) in Supplements. In the other cases, the index may not be explicit,
12 depending on the eigenspaces of $\tilde{F}$ and $\tilde{G}$. We will add a brief comment and show the details in Supplements.

13 **4)** *Implications to learning with GD*: We totally agree that this is an important topic, and are currently working on it.

14 **5)** *Saddle by inactive units*: Theorems 2 and 4 in [4] consider a special parameter, which corresponds to $\boldsymbol{v}_{H_0} = 0$ in
15 the current paper. For general $\boldsymbol{v}_{H_0}$, they are not critical points, and this fact is mentioned in lines 111-112 with full
16 description in Supplements. If $\boldsymbol{v}_{H_0} = 0$, we can show similar results to [4]. In an update, we will clarify this at the
17 paragraph in line 111.

18 **6)** *Embedding of a saddle point*: In the unit replication, Lemma 4 tells that the Hessian of the narrower network is
19 recovered in the wider one. So, embedding of a saddle point with positive and negative eigenvalues gives saddle points.
20 For the embedding by inactive units and propagation, a critical point is not in general embedded into a critical point.

21 **7)** *Refer Section 2 in Introduction*: We will reflect this.

22 **Rev. 2.**

23 **8)** *Embeddings are artificial*: We do not think the three methods are artificial for the reason described in the paragraph,
24 line 88-93. Overparameterization is an important issue in this field as we describe in Sec. 1, and it refers to the situation
25 where a network has more sizes than the one necessary to realize a function. We rigorously formulate and discuss this
26 situation in the paper. It is also important to note that, as discussed in lines 88-93, the classical results [10,14] already
27 proved that the three methods of embedding are the ONLY ways of realizing a network function by a wider network in
28 the case of three-layer models. These embeddings are thus the essential ones in discussing overparameterization.

29 **9)** *Embedding methods are not clear in Theorems*: Theorem 5 does specify the embedding by $\theta_\lambda^{(H)}$, which is defined as
30 a symbol of the Unit Replication in Eq.(4). Additionally, Sec. 3.2 starts with a statement that we consider this type.
31 Theorem 9 clearly states that the embedded point is "defined by Eq(10)", which is the definition of Inactive Units (line
32 173). In an update, we will explicitly place the name of embedding at the theorems to avoid confusion.

33 **10)** *Difference of considerations for the smooth and ReLU activation*: For smooth activation, the inactive units and
34 propagation do not embed a critical point to a critical one in general (lines 111-112), so we discuss only the unit
35 replication for local minima. For ReLU, note that the definition of inactive units is different from that of smooth
36 activation, and discussing an embedding with inactive units is meaningful. Inactive propagation is the same as the
37 smooth case, and a critical point does not give a critical point in general. We will clarify this in an update.

38 **11)** *Different choices of P and Q in PAC-Bayes*: Using the same $P$ and $Q$ does not necessarily give a fair comparison,
39 but they should be chosen so that the bounds are tight. The distributions of the parameters that give similar loss values
40 are different, and the choice of the posterior $Q$ must reflect this difference for meaningful bounds. The prior $P$ is
41 arbitrary, as long as it does not depend on the training data. The choices in the paper reflect these conditions.

42 **12)** *Inconsistency between experiment and theory*: Sec. 5 highlights the difference between the smooth and ReLU
43 models, and the form Eq.(11) holds also in 1-dimensional output (see [4]), which means the difference exits also. While
44 the theory in Sec.5 are based on Theorem 5 or Lemma 4, the experiments use 1-dimensional network for simplicity of
45 realizing zero error. We will make clearer explanations on this point in the paragraph on the experiments.

46 **Rev. 3.**

47 **13)** *Definition of semi-flatness*: In Sec. 1, we say "semi-flat minima, at which a lower dimensional affine subset in the
48 parameter space gives a constant value of error", and use it as a definition. We will make it clearer. Flat-minima often
49 refer to points at which all the directions are flat. "Semi-flat" allows the case that only part of the directions are flat.

50 **14)** *PAC-Bayes analysis already exits*: Our motivation in Sec. 5 is to compare the smooth activation and ReLU in
51 overparaemterized cases. We will clarify this point better and cite the reference.

52 **15)** *How does the number of surplus units affect the landscape?*: Thank you for pointing out this important question.
53 We must admit that the detailed discussions were just left in Supplements by space limitation. For smooth networks,
54 some of the surplus parameters make sharp directions so the bound is linearly increasing to the surplus number of
55 units, while in ReLU there are no sharp directions so that no such major increasing terms exist in the bounds. This is
56 illustrated in the experimental results (Fig.2(b)). We will emphasize this in an update.