1 We sincerely thank the reviewers for their insightful and constructive comments. First, we address the common
2 questions of **R1,R2,R3:** *Definition of GCN:* We define GCN in eq (4), which not only includes Kipf & Welling [16]
3 $D^{-1/2}AD^{-1/2}$, but also $D^{-1}A$ (e.g. [12] GraphSAGE mean aggregator) among others. Other papers (e.g. JK paper by
4 Xu et al, 2018, p3) also consider these as variants of GCN. We will make this explicit in the final version. *Contribution:*
5 the main contribution of this work is to develop deep theoretical understanding of GCNs, which would inform efficient
6 GCN architectures design. We do not intend to achieve the state-of-the-art graph classification model. We will tone
7 down the emphasis on the model architecture design, as similar work may already exist. *Real-world Experiment:* To
8 verify theoretical results, we needed the ground truth model, which is unknown for real data. We are happy to add
9 real-world experiments in the final version. Next, we address the specific concerns raised by each reviewer below.

10 **R1** *graph attention network (GAT)* multi head attention in GAT concatenated $K$ independent attention mechanism
11 with the *same* propagation rule (PR) while ours utilizes modules with *different* PRs. We will consider incorporating
12 attention mechanism as a future work. *jumping knowledge network* Both JK-network and ours use residuals, but ours
13 uses residual connections for GCN while JK-network is designed for Message Passing Graph Neural Networks. *graph*
14 *classification setup* We used mean-pooling to aggregate node-level representations, after which a single number is
15 passed to a classification layer. We will clarify this and include discussions with GAT and JK in the final version.

16 **R2** *authors try to use GNN for graph generation* There seem to be some misunderstandings about the paper. We do *not*
17 learn graph generation. Instead, we use GCN to learn graph moments and to classify graphs. *... $D^{-1}A$ has limitations*
18 *in learning degree of a graph.* We are *not* claiming the limitation of a particular GCN variant (Fig 2 is only an example).
19 Our key message is that unlike fully connected neural networks, GCNs are not universal approximators. Therefore,
20 choosing the right PR (for example, $A$ vs $D^{-1}A$) is crucial in the GCN's ability to learn graph moments. We provide
21 theoretical analysis and offer a solution that can alleviate this issue. Note that $D^{-1}A$ is *not* our definition of GCN (see
22 eq (4)), nor is it used in our theoretical analysis. *... why $D^{-1}A$ is associated with node permutation invariant.* This is
23 a misunderstanding. We are *not* claiming that $D^{-1}A$ is related to node permutation invariance. In fact, any GCN in
24 the form of eq. (4) $F(A,h) = \sigma(f(A) \cdot h \cdot W + b)$ is permutation invariant, regardless of the function $f(A)$ (see sec
25 2.1 and 2.3). The purpose of Proposition 1 and Theorem 2 is to argue that GCN can be restrictive due to permutation
26 invariance, thus having the right PR, activation and number of layers is crucial in its ability to learn graph moments.

27 *did not see these papers used these three different graph operations* It is not explicit but easy to derive. [12] GraphSAGE
28 with MEAN aggregator, averages over $h_i + \sum_{j \in \mathcal{N}_i} h_j$, which is equivalent to the operator $\tilde{D}^{-1}\tilde{A}$ where $\tilde{A} = A + I_N$ is
29 an adjacency with self-loops. Kipf-Welling GCN uses $\tilde{D}^{-1/2}\tilde{A}\tilde{D}^{-1/2}$ ([16] eqs (7),(8)). [18] uses the graph module $A$
30 ([18] eq (2)) *simple variant of original GNN* Our theoretical analysis demonstrates the importance of PR (e.g. using $A$,
31 GCN cannot learn $\sum_j (D^{-1}A)^p_{ij}$ or vice versa), as well as to have sufficient number of layers with residual connections.
32 Hence, having a modular design with sufficiently many layers and residual connections would be useful for learning
33 graph moments. We have no intention to claim the novelty of our proposed GCN, but only to validate our theoretical
34 findings. In fact, we pointed out the similarities our model with GIN in lines 173-176. *graph generation baselines* As
35 we are *not* generating graphs in this work, our method is unrelated to graph generation baselines.

36 **R3** *better not to claim...expressiveness of GCNs as the first contribution.* While it may be known in the literature,
37 we are not aware of any rigorous theoretical analysis for the exact same setting. We will add more references and
38 tone down the claim in the final version. *bridge the theoretical analysis with the proposed design.* Our theoretical
39 analysis shows the limitations of having a single PR in learning graph moments. It also points out the importance of
40 having sufficient number of layers and residual connections. Our modular design combines these results and arrive at
41 an architecture with multiple PR modules, layers and residual connections. We will improve the writing and include
42 more descriptions in the final version. *compare with some existing GCN designs* Our GCN already includes several
43 existing GCN designs, which we refer as different PRs. We are happy to include other GCN designs for comparison as
44 well. *Could not find the code...* We apologize for the confusion. We will release the code in the updated version.

45 **R4** *lack of related work concerning graph moments.* Graph moments, or "Graph Power" is a concept from graph
46 theory ( see the Lin and Skiena (1995) ) and has been used extensively in network science. Graph moments encodes
47 topological information of a graph and is closely related to graph coloring and Hamiltonicity. We will include references
48 from graph theory and network science. Note that we are not aware of any other work that learns graph moments using
49 GCNs. *Notation: lack of explicit introduction of notational convenience* Thank you, will fix. *Line 112* yes, for brevity
50 $M_p$ is $M_p(A)$. *an index such as in $f(A)_i$.* Yes, $f(A)_i$ is the $i$th component. *The moments are matrix valued?* No, our
51 definition eq. (1) is vector valued, (summed over one node index). *residual connections are very strongly stated...* We
52 agree with the reviewer regarding the positioning of our contribution. We will tone down the emphasis on modular
53 design. The restrictions found in our theory show a need for including multiple propagation rules, hence the module.