

1 We thank the reviewers for their very helpful remarks. Overall, we accepted all small corrections, cited the additional
2 literature, and added missing details pointed out by the reviewers. In the following, we address common concerns
3 and then individual reviewers' questions, but to our regret, space limits dictate that we cannot respond to every point.
4 Multiple reviewers pointed out the lack of attached code, for which we apologize. Our code is not proprietary, and
5 we plan to fully release it with the camera-ready version. Unfortunately attachments are not allowed for rebuttals.
6 Reviewers also commented that our test data is a small set. From a chemistry perspective, the tested lignin molecule
7 with multiple linkages is already a complex system that involves many viable conformers, wherein good sampling
8 performance can provide chemically rich information toward studying depolymerization reaction pathways. Even this
9 single large lignin molecule required *weeks of compute time* to gather enhanced MD results. In fact, it is this very
10 compute cost that our method attempts to address. Moreover, note that TorsionNet *generalizes over many molecules*
11 *during the training curriculum*. Nonetheless, we are adding more alkane eval environment data for the final submission.

12 **Response to Reviewer 1.** 1. *Concerns on Gibbs score.* The normalization constants (Z_0 , E_0) are *not* assumed to
13 be "ground truth". Normalization is used for (1) numerical stability of rewards and (2) interpretability of results.
14 Furthermore, Gibbs free energy is always **relative** and force field methods mostly concern themselves with energy
15 *differences* rather than absolutes. As Reviewer 2 comments, this metric is intuitive and grounded in the physics.

16 **Response to Reviewer 2.** First, we graciously accept Reviewer 2's positive remarks and thank them for their support.
17 1. *Connection of Gibbs score to other statistics/ bias to thermal accessibility.* We thank Reviewer 2 for pointing out that
18 (log) Gibbs score is deeply grounded in physics and related to the Boltzmann distribution as relative population and
19 relative free energy. We present it as "new" as it has not been used in conformer generation literature, and as a "score"
20 to be clear to an ML audience. In general, drug screening moves in the direction of finding more accessible conformers;
21 good bindings must have relatively low Gibbs free energy. Furthermore, while the Gibbs score is designed as a generic
22 metric, we can bias the reward for specific tasks and continue using the developed methods.

23 2. *De-duplication of ETKDG/symmetry concerns.* The output of all compared methods (incl. ETKDG) goes through a
24 minimization step, and then a distance exclusion function to remove all near duplicates. Similarly, while Gibbs may
25 be undercounted due to symmetries, the comparison is fair since we apply the same scoring uniformly. We agree the
26 symmetry issue is important, and must be investigated when dealing with highly symmetrical molecules (lignin is not).

27 3. *Concerns about Confab/random search over rotations.* To clarify a misconception: Confab essentially **is** a random
28 search over rotations. It iterates in a random order over all possible conformers (O'Boyle et al., 2011). Therefore,
29 sampling the first N of this exhaustive search is similar to TorsionNet with uniform distribution policy, except without
30 replacement. Running it to completion on the T-alkanes experiment is not necessary, as in this case we are not looking
31 for the entire partition function, but merely the lowest energy conformation (which can be identified by eye).

32 4. *Number of calls to MMFF.* It takes around 500,000 evaluations on **non-test** lignins and 1000 at inference time to
33 achieve the score presented in the paper. To compare, SGMD takes 25 million Charmm evals at 2 fs steps on test lignin.

34 **Response to Reviewer 3.** We thank Reviewer 3 for their positive remarks.

35 **Response to Reviewer 4.** 1. *Unfair CPUtime/walltime comparison.* TorsionNet is not using GPU speedup at inference
36 time. *We specifically ran it on CPU to achieve fairer test results.* We will add more experiment detail to appendix.

37 2. *Gibbs score vs. all-by-all distance.* It is not clear to us what metric should be used to compare metrics, and we
38 consider our Gibbs score an attempt to import the well established thermodynamics principle of free energy into the
39 context of conformer generation as observed by Reviewer 2. Nonetheless, we agree that analyzing the diversity of our
40 generated conformers via all to all distance in the supplement would be valuable, and will do so.

41 3. *Ground truth for lignin.* There are only a handful of experimental crystal structures that have been published
42 (Vermaas et al., 2019) and they are mostly limited to small dimeric structures only. Data from MD serve as the ground
43 truth dataset, comprised of physics-based conformers, similarly to (Simm and Hernández-Lobato, 2019).

44 4. *Are MD simulations of lignin useful for energy production and any paper that talks about it?* Scientists employing
45 mechanochemical experimental processes (where MD's ability to model many conformers is important) to extract
46 renewable energy from plant biomass argue that: "Understanding wood and lignin processing on a molecular level
47 appears essential for improving their degradation efficiencies." (Kleine et al., 2013)

48 5. *Why CL experiments? How is CL used?* The CL experiments are simple models to demonstrate the validity of our
49 theoretical work. Fig 3. y-axis shows the distance in energy between the best sampled conformer of the current model
50 and the global best conformer. We observe how errors drop for unseen test molecules as more examples are added to the
51 train set. We train TorsionNet sequentially from small molecules to large as described in Secs. 3.2,3.3 and App. C.1.

52 6. *Protein folding IS working on MD.* We agree, and are unaware of anything in our paper that suggests otherwise.
53 Protein folding commonly uses enhanced sampling MD methods (we do as well in our benchmarks).