

434 **Dataset card**

- 435 1. Our dataset contains 57231 problems in the split of Lean Workbook and 82893 problems
436 in the split of Lean Workbook Plus. We provide the natural language statement, answer,
437 formal statement, and formal proof (if available) for each problem. These data can support
438 autoformalization model training and searching for proofs.
- 439 2. We open-source our code at <https://github.com/InternLM/InternLM-Math> and our
440 data at <https://huggingface.co/datasets/InternLM/Lean-Workbook>.
- 441 3. Croissant metadata URL: [https://huggingface.co/api/datasets/internlm/](https://huggingface.co/api/datasets/internlm/Lean-Workbook/croissant)
442 [Lean-Workbook/croissant](https://huggingface.co/api/datasets/internlm/Lean-Workbook/croissant).
- 443 4. The license of our dataset is Apache 2.0.
- 444 5. We will host our dataset in Huggingface and our code in GitHub. We will maintain this
445 dataset with further improvement.
- 446 6. DOI of dataset: 10.57967/hf/2399