

STEPPROOF: STEP-BY-STEP VERIFICATION OF NATURAL LANGUAGE MATHEMATICAL PROOFS

Anonymous authors

Paper under double-blind review

ABSTRACT

Interactive theorem provers (ITPs) are powerful tools for the formal verification of mathematical proofs down to the axiom level. However, their lack of a natural language interface remains a significant limitation. Recent advancements in large language models (LLMs) have enhanced the understanding of natural language inputs, paving the way for autoformalization—the process of translating natural language proofs into formal proofs that can be verified. Despite these advancements, existing autoformalization approaches are limited to verifying complete proofs and lack the capability for finer, sentence-level verification. To address this gap, we propose StepProof, a novel autoformalization method designed for granular, step-by-step verification. StepProof breaks down complete proofs into multiple verifiable subproofs, enabling sentence-level verification. Experimental results demonstrate that StepProof significantly improves proof success rates and efficiency compared to traditional methods. Additionally, we found that minor manual adjustments to the natural language proofs, tailoring them for step-level verification, further enhanced StepProof’s performance in autoformalization.

1 INTRODUCTION

Mathematics is the basic tool for the development of science, and the reliability of its conclusions affects the stable growth of various disciplines. With the development of the mathematical edifice, mathematical proofs have become more complex and lengthy. The verification of mathematical proof often requires years of careful verification to ensure the accuracy of the work. However, reading a mathematical work requires a large amount of knowledge, and in the face of the many branches of mathematics today, traditional manual verification has become increasingly disastrous. Thus, an idea arose to validate mathematical work written in natural language automatically.

At present, there are two kinds of automatic verification of mathematical proof. One is to write the mathematical certificate into a machine code that can be verified by a specific expert system, which is called the interactive theorem prover (Harrison et al., 2014). After more than 40 years of development, the interactive theorem prover has begun to take shape and has been used in the verification of many mathematical works (Maric, 2015). However, because such machine-verifiable proofs need to be written in a specific programming language, whose learning cost is relatively high, interactive theorem provers are only used by a small number of mathematicians at present (Nawaz et al., 2019).

On the other hand, after the advent of the large language model (Zhao et al., 2023), through prompt engineering (Liu et al., 2023) and few shot learning (Wang et al., 2020b), large language model can be applied to many natural language processing tasks, and achieved considerable performance (Achiam et al., 2023). However, due to the hallucination problem of large language model (Ji et al., 2023), its performance in mathematics and strong logic-related work is not good enough (Huang et al., 2023), so the lack of reliability of large language model to verify mathematical work is easy to cause a lot of errors. In this context, a method that combines large language models and theorem provers gradually comes into people’s view, which is called autoformalization verification (Li et al., 2024).

At present, the existing automatic formal verification work generally adopts a FULL-PROOF strategy. Although such a strategy has achieved some impressive performance in some studies, there are still many problems in the stability of its proof and the fine-grained verification. Aiming at the

054 problems existing in FULL-PROOF, we innovatively propose a new automatic formal proof strat-
055 egy, STEP-PROOF. And the performance of StepProof is better than traditional methods in several
056 aspects.

057 In summary, our contributions mainly include the following points: 1. We pioneered a novel natural
058 language mathematical verification method StepProof, which realizes informal mathematical proof
059 verification at the sentence level. 2. We were the first to realize the test of automatic formalization
060 capabilities on small open-source LLMs. 3. Compared with existing methods, the StepProof method
061 has been significantly improved in all aspects of performance.

062 In this paper, we will first make a brief summary of the existing relevant works and the technical
063 background in Chapter 2. In Chapter 3, we will give a detailed introduction to the two types of
064 strategies, FULL-PROOF and STEP-PROOF, and point out many problems faced by traditional
065 methods. In Chapter 4, we set up a series of experiments to verify the performance of STEP-PROOF
066 in verification tasks and its superiority over FULL-PROOF. At the same time, we also carried out
067 a detailed analysis of some phenomena observed in the experiment to further explain the reason
068 for the advantages. In the end, we analyze and look forward to the current limitations and future
069 development direction.

071 2 RELATED WORK

072
073 **Theorem Prover:** Theorem provers can be roughly divided into two types, interactive theorem
074 provers (ITPs) in which the user can enter and verify the existing proof (Asperti, 2009), and auto-
075 mated theorem provers (ATPs) that try to prove the statements fully automatically (Harrison, 2013).
076 The two types of theorem provers are not mutually exclusive (Nawaz et al., 2019). Most ITPs such
077 as Isabelle (Paulson, 1994), Coq (Huet et al., 1997) and Lean (De Moura et al., 2015) are supported
078 by ATPs that try to automatically prove "obviously" intermediate steps in the proof entered by the
079 user. The entered proofs are rigorously verified back to the axioms of mathematics. Different ITPs
080 use different axiomatic foundations, e.g. set theory, first-order logic, higher-order logic, etc. Each
081 ITP use its own language and syntax, which makes the learning cost of theorem prover high, and
082 precludes ITPs from being widely used (Yushkovskiy, 2018).

083 **Large Language Model:** In recent years, large language models (LLMs) have achieved outstand-
084 ing performance in many downstream tasks of natural language processing. LLMs such as Llama
085 (Dubey et al., 2024), GPT series (Ye et al., 2023) and GLM 4 (GLM et al., 2024) are trained on
086 large databases to understand user input in natural language and produce the related output. While
087 large language models perform well in general tasks such as translation, their performance in deal-
088 ing with logic problems has been limited. A lot of work has also shown that large language models
089 are prone to a variety of problems when dealing with logic problems (Yan et al., 2024; Wan et al.,
090 2024). Although in the subsequent iteration of the model, the developers provided a large amount
091 of high-quality logic-related data to improve the logic capability of the model, the effect that could
092 be achieved was still very limited (Lappin, 2024). Therefore, it has become a trend to add an expert
093 system to the model to improve its accuracy, such as RAG (Fan et al., 2024), which is currently
094 commonly used in question-answering systems. On the other hand, step-by-step reasoning has been
095 proven can improve the reasoning ability of existing LLMs (Wei et al., 2022; Khot et al., 2022),
096 which gives us a hint to apply in autoformalization.

097 **Autoformalization:** The definition of automatic formalization is very broad, but can be roughly
098 seen as understanding and extracting translation from natural language to obtain the required struc-
099 tured data or formal language, such as entity relation extraction (Nguyen & Grishman, 2015). Early
100 automatic formalization work involved extracting logical propositions from natural language in ad-
101 dition to entity relation extraction (Singh et al., 2020; Lu et al., 2022). However, the main problem in
102 this kind of work is that the extracted logical propositions lack corresponding symbols for derivation
103 and application, so the output of the automatic formalization output cannot be directly applied. With
104 the launch of pre-trained language models such as transformer and BERT, language models have
105 a stronger understanding ability. Wang et al. (2020a) conducted an early automatic formalization
106 attempt for the theorem prover Mizar. However, due to the limited size and training expectations of
107 the models at that time, the effects they could achieve were very limited. With the rapid expansion of
the scale of models and predictions, many new and better performance automatic formalization work
has emerged, such as the Majority Voting method by Lewkowycz et al. (2022), the DSP method by

Jiang et al. (2022)., and the method proposed by Zhou et al. (2024), in DTV (Don't Trust: Verify). They further improved the automatic formalization of the system by combining large models with some syntax-modifying filters. However, on the one hand, all these works are tested in the closed-source large model Minerva, which lacks the testing work of open source model and small model. Meanwhile, all these works adopt the strategy of FULL-PROOF, which has poor controllability for the formal output of the model, and it is difficult to locate the error point of non-formal proof. Qinghua et al. proposed a problem location method based on the FULL-PROOF strategy in SlideRule. However, this method relies heavily on the format and quality of the generated formal content, so it cannot achieve 100% problem locations detected. To solve these problems, we propose StepProof, an automatic formalization strategy that can realize sentence-level verification, to realize the verification of natural language mathematical proofs. Moreover, LEGO-Prover (Wang et al., 2023) proposed a new methodology to decompose the whole proof into several sub-proofs. Although the idea of step-proof is taking shape, it still requires some extra generation of the sub-proof formal statement generation, which increases the error probability of formalization.

3 STEP-PROOF

In this chapter, we will provide a detailed introduction to the workflow design of StepProof. We will also compare the STEP-PROOF strategy used by StepProof with the FULL-PROOF strategy adopted by existing autoformalization systems, highlighting the problems with traditional strategies and the advantages of STEP-PROOF over FULL-PROOF.

3.1 FULL-PROOF

Current research on natural language proofs formal verification predominantly employs the FULL-PROOF strategy, as seen in methods like DSP and DTV. The workflow of FULL-PROOF method can be roughly illustrated as the left of Figure 1. Users submit a provable problem along with its proof process. The problem is first formalized, then the informal problem, the formalized problem, and the entire informal proof are packaged as inputs to a large language model for formalization. After obtaining the formal proof, it is combined with the formalized problem and input into a theorem prover for rule-based formal verification. The verification results are then returned to the user. Despite the clarity and simplicity of the FULL-PROOF workflow, it has significant drawbacks.

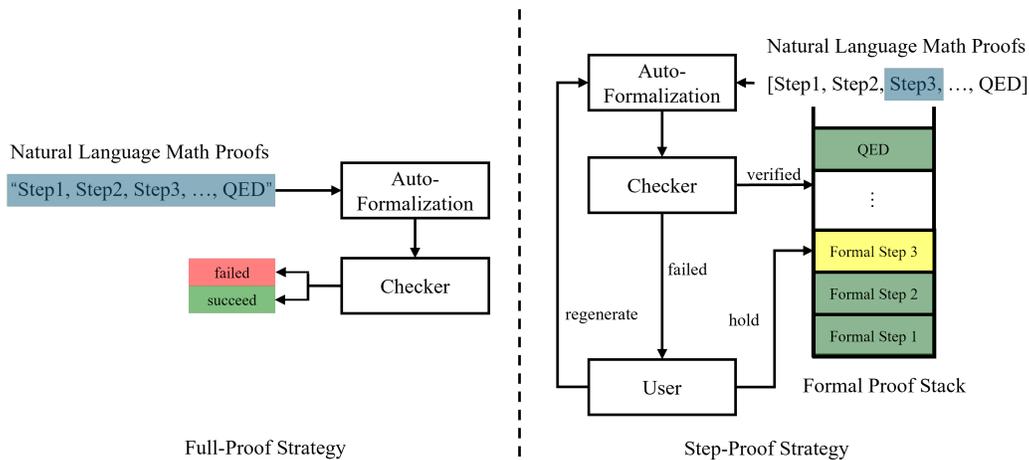


Figure 1: Full-Proof Strategy generate from the whole proof and only provide proof result instead of detailed feedback to help user improve the proof. While Step-Proof separate the whole proof into sub-proof to proof from the bottom to the top and enable the user to get more detailed feedback and fine-grained operation.

First, in the FULL-PROOF automated formalization process, the highly structured and formalized nature of the input and output, coupled with the numerous similarities in solving mathematical

162 equations, often leads to excessive noise in the output. To obtain the desired formalized content,
163 numerous filters must be set up, which results in considerable waste and contamination of generated
164 content. This can also cause generation loops, where the same content is repeatedly generated, a
165 common issue in FULL-PROOF strategies.

166 Additionally, the length of proofs varies, and LLMs in FULL-PROOF struggle to adjust the out-
167 put’s `max_new_tokens` effectively based on input length. This leads to shorter proofs not being
168 truncated in time, thus generating repetitive or noisy content, and longer proofs lacking sufficient
169 `max_new_tokens`.

170 Lastly, the stability of FULL-PROOF generation is poor. For instance, a full proof might be almost
171 entirely correct except for a minor error in a small step, leading to the failure of the entire formal
172 proof. Users attempting regeneration may find previously correct parts presenting erroneous. Thus,
173 FULL-PROOF requires highly accurate one-time generation of the entire content.

174 Moreover, the correlation between formal and informal content generated by FULL-PROOF is weak,
175 making it difficult for users to map formal feedback to corresponding informal content. Although
176 LLMs can generate formal proofs with annotations to map back to informal proofs as Qinghua et
177 al proposed the failure step detection method in SlideRule, this approach is unstable in practice and
178 increases the required token count.

180 181 182 3.2 STEP-PROOF

183 To address the numerous issues faced by the FULL-PROOF strategy, we innovatively propose STEP-
184 PROOF. STEP-PROOF employs a step-by-step generation and verification strategy, offering better
185 performance and stability compared to FULL-PROOF. The workflow of STEP-PROOF is illustrated
186 in the right of Figure 1.

187 Unlike the one-time generation and verification of FULL-PROOF, STEP-PROOF assumes each sen-
188 tence in the proof is a verifiable sub-proposition. Each step is formalized and pushed onto a formal
189 proof stack, where it is verified along with other sub-propositions in the stack. Upon successful
190 verification, the formalized proof and informal proof are packaged as inputs for the next step. For
191 failed steps, the Step Proof allows users to backtrack, retaining previously verified steps and only
192 clearing the erroneous step from the stack. The StepProof can then either re-formalize or optimize
193 the existing informal step as needed.

194 STEP-PROOF offers several advantages over FULL-PROOF. First, it only needs to formalize sin-
195 gle sentences in context, resulting in shorter, less noisy, and more stable output. The step-by-step
196 generation strategy also means that each step’s length is relatively fixed, eliminating the need for
197 adjusting `max_new_tokens` and allowing the use of smaller `max_new_tokens` for formalizing longer
198 theorems.

199 Moreover, STEP-PROOF’s incremental generation and verification tolerate step errors well. Only
200 the specific erroneous step needs to be retracted, rather than regenerating the entire content, en-
201 hancing robustness and efficiency. Finally, Step Proof ensures high correspondence between each
202 informal and formal proof step, providing users with finer operational granularity. For instance,
203 users can suspend a correct but incomplete step and assume it is correct to proceed, a functionality
204 almost impossible under the FULL-PROOF strategy.

205 We also designed a simple and user-friendly interactive interface for the user, as shown in Figure
206 2. Users can complete the natural language proof of the problem through interface interaction, and
207 each step of the proof will be formally verified after submission, and the verified proof step will be
208 marked in green in the background to indicate that the current step has passed the verification and is
209 reliable. If the current step does not pass the automatic validation but the user thinks it is true and
210 wants to continue, you can select HOLD, and the current step will be highlighted in yellow in the
211 background to indicate that the current step is in a suspended state. After completing all the proof
212 steps, the user can input QED to indicate that the system has completed the proof, and the system
213 will combine all the steps to perform the final verification of the proof target. Through the interactive
214 interface, the user can also realize the PDF export of the proof process.

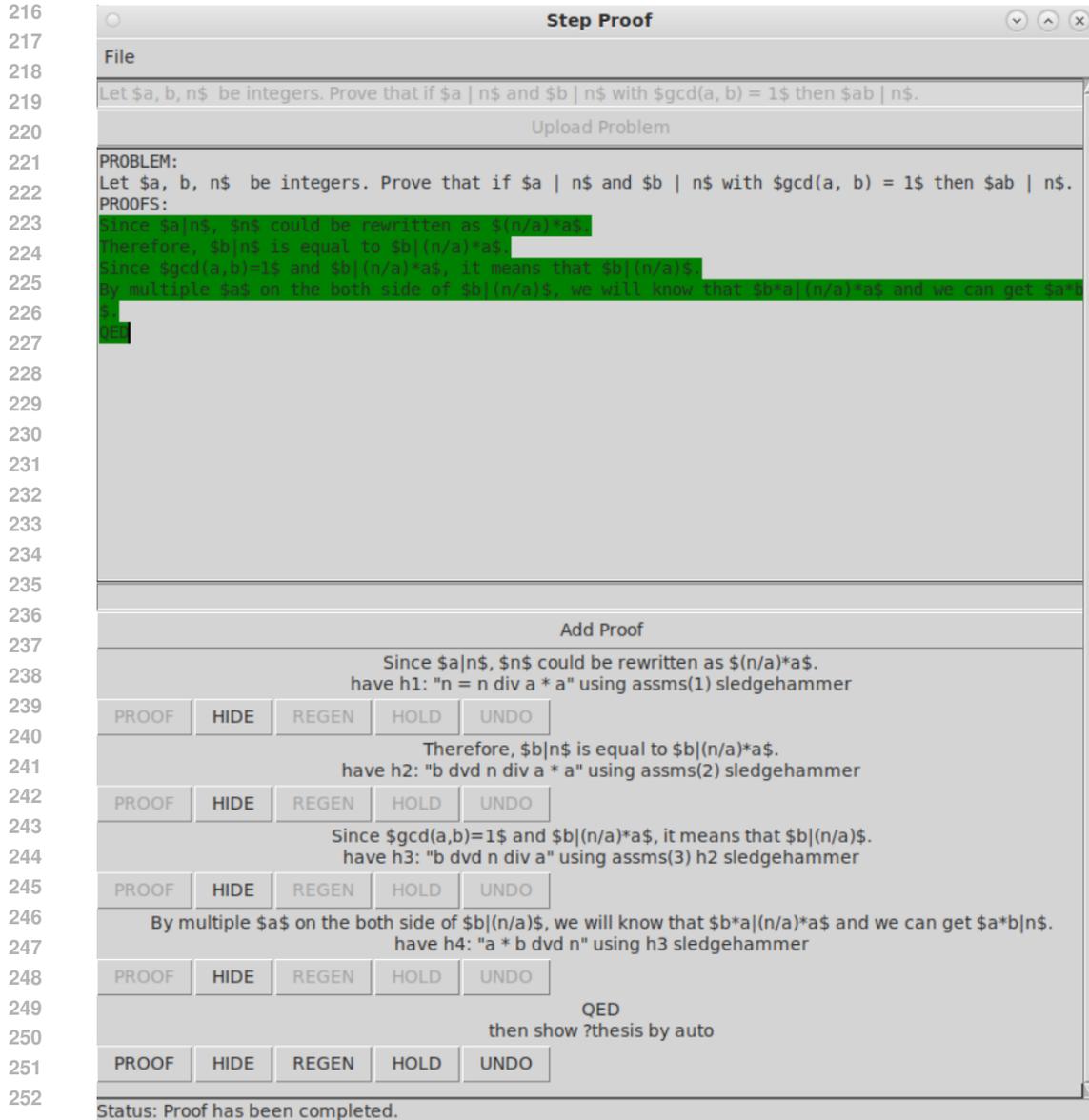


Figure 2: User Interface of StepProof

255 4 EXPERIMENT

256 4.1 EXPERIMENT SETUP

257
258
259
260
261
262
263
264
265
266

To validate the performance advantages of the STEP-PROOF strategy over FULL-PROOF, and to compare StepProof’s overall performance against existing automated formal proof methods, we conducted both strategy performance tests and baseline tests using the same dataset and model settings.

267
268
269

For the test dataset, we used GSM8K, which includes a large number of informal mathematical problems and their correct informal proofs. These informal proofs can be easily segmented into a series of sub-steps. We chose Llama3 8B-Instruct as the large language model for automated formalization. Existing autoformalization tests use closed-source models, and we aim to fill this gap by using an open-source small LLMs.

We set the temperature to 0.3 to balance stability and flexibility. Given the small parameter count of 8B, we used a single example for the few-shot in strategy performance tests. With proof steps in GSM8K (Cobbe et al., 2021) averaging 4-5 steps, we set the max_new_tokens for FULL-PROOF to 1024, and 256 for each step in Step Proof. The test environment consisted of a single NVIDIA A4000 16GB, 8 cores of AMD5800X, and 32GB DDR4 3200 RAM. Isabelle2024 was used as the formal theorem prover with only *Main* prove library as the theorem base¹, with Isabelle-client (Shminke, 2022) as the testing service proxy.

In the strategy performance tests, we evaluated the performance of FULL-PROOF and STEP-PROOF on the GSM8K test set using the following metrics: 1. One-attempt generation proof pass rate r_p . 2. Average formalization time for passed proofs μ_f . 3. Variance in formalization time σ_f^2 . 4. Average proof time for passed proofs μ_p . 5. Variance in proof time σ_p^2 .

In the baseline tests, we compared the multi-attempt test results of GSM8K by Lewkowycz et al. (2022) in Majority Voting and Zhou et al. (2024) in DTV with our results. We evaluated the performance based on the number of attempts and multi-round proof pass rate, allowing up to 10 retries for each failed step in each formalization attempt.

At the same time, considering that StepProof has better granularity than traditional proof tasks, we not only evaluated the overall proof passing rate but also counted the proportion of step-proof passing. For example, if a 6-step proofs can be verified to be true in 3 steps, then we will mark that the step pass rate of the proof is 0.5. In this way, we quantify the formal proof capability of the method more comprehensively rather than the Proof Passing Rate.

In addition, to verify the influence of the writing method of non-formal proof on the passing rate of StepProof formal proof, we extracted 100 questions from the Number theory of MATH (Hendrycks et al., 2021) and made simple manual modifications to make the proof step more consistent with the proof requirement of StepProof.

4.2 EXPERIMENT RESULTS

	Proof Passing Rate	Formalization Time	Proof Time
	r_p	$\mu_f \pm \sigma_f^2$	$\mu_p \pm \sigma_p^2$
FULL-PROOF	5.30%	9.54 ± 12.64s	214.93 ± 20864.97s
STEP-PROOF	6.10%	5.83 ± 4.24s	130.12 ± 5271.65s

Table 1: Performance Test of Full-Proof and Step-Proof

In strategy performance tests, the STEP-PROOF strategy outperformed the FULL-PROOF strategy across the board on the GSM8K test set. As shown in Table 4.2, the STEP-PROOF strategy improved the one-attempt proof pass rate by **15.1%** compared to FULL-PROOF. In terms of average formalization time, STEP-PROOF required **38.9%** less time than FULL-PROOF. For average proof time, STEP-PROOF achieved a **39.5%** performance improvement over FULL-PROOF. Additionally, STEP-PROOF showed more stability in both formalization and proof time compared to FULL-PROOF. These results confirm that our strategy offers better performance, efficiency and stability.

In the baseline test as shown in Table 4.2, Step Proof surpassed DTV in multi-round verification tests on GSM8K, achieving a **10.3%** performance improvement. Moreover, StepProof required fewer attempts compared to DTV², demonstrating its superior proof capability and further validating the advantages of the StepProof methods.

¹Here, we only use Main as the proof library, considering that the use of different libraries will greatly affect the formal verification ability of the theorem prover, in order to provide a relatively standard index. In StepProof, we do not introduce other libraries to further improve the proof ability of the theorem prover. Introducing more libraries in practice will greatly improve the proof ability of the theorem prover, and also improve the proof ability of the whole system to some extent.

²Don’t Trust: Verify (DTV) originally used two close source models—GPT3.5 as the problem generation model and Minerva 8B as the proof generation model, while due to the Minerva being inaccessible and GPT3.5 being costing, we use the same method in DTV, but replace the LLM into Llama3.

	Attempts	Proof Passing Rate	Comments Rate	Model
Majority Voting	64	16.2%	-	Minerva 8B
Don't Trust:Verify*	64	25.3%	31.3%	Llama3 8B
StepProof	10	22.0%	100%	GLM4 9B(4bit)
StepProof	10	27.9%	100%	Llama3 8B

Table 2: Baseline Test in GSM8K

Step Pass Rate	LLAMA3 8B		GLM4 9B(4bit)	
	1 attempt	10 attempts	1 attempt	10 attempts
$0 = r_s$	79.6%	50.5%	83.9%	55.4%
$0 < r_s$	20.4%	49.5%	16.1%	44.6%
$0.5 \leq r_s$	14.6%	38.1%	13.4%	41.9%
$1 = r_s$	6.1%	27.9%	4.8%	22.0%

Table 3: Step Passing Rate Distribution in GSM8K

In the step pass rate test as shown in Table 3, we found that StepProof was able to perform some degree of validation for nearly half of the proofs after 10 rounds of trying. In LLAMA3 8B, 38.1% of the proofs completed more than half of the verification, and 27.9% of the proofs completed all of the verification. Compared with a single attempt, this is a significant improvement. At the same time, we propose a new indicator-step passing rate (r_s) for a more comprehensive evaluation of the proof of automatic formal verification methods.

Step Passing Rate	Original	Modified
$0 = r_s$	35%	32%
$0 < r_s$	65%	68%
$0.5 \leq r_s$	42%	45%
$1 = r_s$	6%	12%

Table 4: Step Passing Rate Distribution in Number Theory

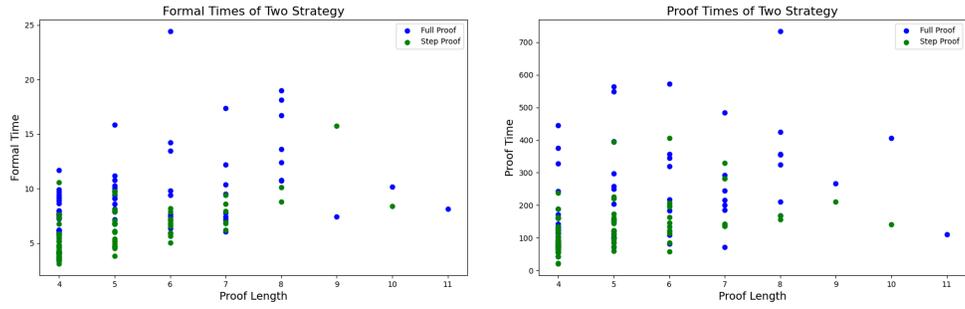
In our test to verify the influence of informal proof writing on the proof pass rate (as shown in the Table 4), we found that the proof pass rate was significantly improved after simple fitting of the informal proof. This shows that compared with FULL-PROOF, STEP-PROOF adopts proof mathematics that is more suitable for step verification, which will be more conducive to improving the pass rate of automatic formal verification.

4.3 EXPERIMENT ANALYSIS

From Figure 3, we find that compared with the FULL-PROOF Strategy, STEP-PROOF takes less time to formalize and prove and is more stable. On the one hand, the generation strategy of Step-Proof makes the content generated in a single attempt less and more stable. Stable step content reduces the number of false proofs and the time to repeatedly prove successful content. Therefore, in terms of both generation efficiency and proof efficiency, STEP-PROOF is superior to FULL-PROOF.

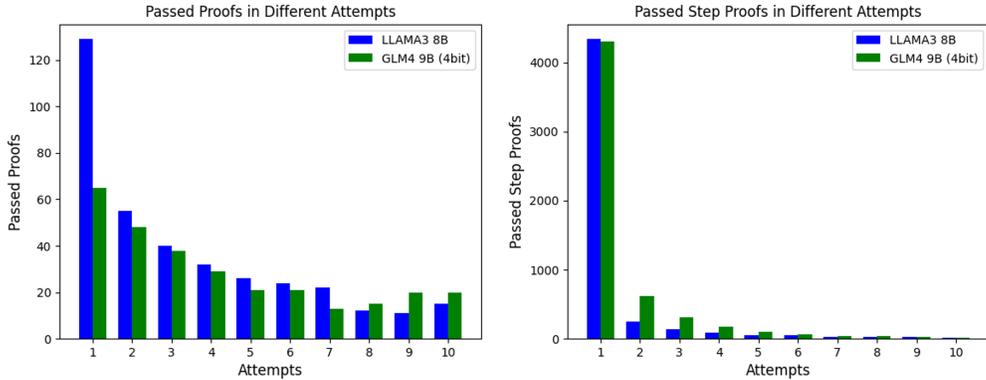
To investigate the relation between step proof pass rate and number of attempts, we plotted Figure 4. It shows that most steps can be proven with relatively few attempts, with only a small fraction requiring multiple tries. We believe the main limitation to step pass rate lies not in the model's formalization ability, but in whether the informal proof steps are suitable for conversion into provable formal steps. We found that many steps in the test set cannot be formalized into provable steps. These unformalizable informal steps significantly limit the further performance improvement of

378
379
380
381
382
383
384
385
386
387
388
389



390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409

Figure 3: Time Cost of Two Strategies in Formalization and Proof



410
411
412
413
414
415
416
417
418
419
420
421
422
423
424

Figure 4: Passed Proofs and Steps in Different Attempts

411 Step Proof. Therefore, to better achieve verification in practical applications of StepProof, each
412 proof sub-step should be ensured to be a formally verifiable sub-proposition as much as possible.
413

414 In addition, we also found that the processing capacity of LLMs is poor when faced with the pro-
415 cessing of continuous equations, and it is easy to fall into the generation loop during the genera-
416 tion process, resulting in the poor quality of the generated formal proof. Therefore, in practical appli-
417 cations, it is necessary to split the steps as much as possible, rather than pursue once completion.
418 At the same time, since the automatic proof ability of automated theorem proving is also limited,
419 resolving the proof process is not only conducive to the formalization process, but also very helpful
420 to the verification process.
421

422 5 LIMITATIONS & FUTURE RESEARCH

423
424
425 At present, the performance of StepProof in small models has been verified in various aspects.
426 However, due to the lack of equipment and funds, the performance of StepProof in larger models has
427 not been verified, although the previous relevant work verified the conclusion that FULL-PROOF
428 strategy has better performance in larger models. However, the STEP-PROOF strategy has not yet
429 been tested on a larger model.

430 In addition, StepProof is strict for users to enter proof steps, while FULL-PROOF is more flexible
431 and can incorporate some non-proof explanatory language or prompt words into the proof. However,
StepProof will make mistakes in the step proof because the prompt word is not provable.

432 Finally, StepProof pays more attention to the sequential proof with steps, but when faced with some
 433 structured proof methods, its performance is still limited, while FULL-PROOF can better capture
 434 the overall proof structure.

435 In the future, our work will mainly start from the following two points: 1. At present, we are writing
 436 a corpus for StepProof for automatic formal verification tasks oriented to step verification, and we
 437 hope that a more targeted corpus will help improve the step formalization ability of the model. 2.
 438 We will implement specific structured proof based on the structured design of the system to further
 439 improve the integrity and flexibility of StepProof proof.
 440

441 6 CONCLUSION

442 In this paper, we innovatively propose a new automatic formal proof method called StepProof. Step-
 443 Proof implements sentence-level formal verification of natural language mathematical proofs, al-
 444 lowing users to conduct more flexible formal verification. At the same time, compared with the
 445 traditional FullProof strategy, StepProof has been significantly improved in formalization, proof
 446 efficiency and proof accuracy.
 447

448 In addition, StepProof can preserve the contents of the proof that has already been verified, providing
 449 more information than the traditional Full-Proof strategy that can only indicate the passage and
 450 failure of the proof. We used a small model on the GSM8K data set to test the proof pass rate, and
 451 its performance reached the level of state-of-the-art.
 452

453 We also conducted a test on the Number Theory dataset of MATH to test the effect of formal content
 454 writing on the proof pass rate. The experimental results show that by optimizing the non-formal
 455 proof for step verification, the passing rate of the non-formal proof can be significantly improved.
 456 We will further optimize the architecture of StepProof in future work to make it more flexible to
 457 handle formal verification of various non-formal mathematical proofs.
 458

459 REFERENCES

- 460 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-
 461 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
 462 report. *arXiv preprint arXiv:2303.08774*, 2023.
 463
- 464 Andrea Asperti. A survey on interactive theorem proving. *URL: [http://www.cs.unibo.it/~as-
 466 perti/SLIDES/itp.pdf](http://www.cs.unibo.it/~as-

 465 perti/SLIDES/itp.pdf)*, 2009.
- 467 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 468 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
 469 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
 470 2021.
- 471 Leonardo De Moura, Soonho Kong, Jeremy Avigad, Floris Van Doorn, and Jakob von Raumer. The
 472 lean theorem prover (system description). In *Automated Deduction-CADE-25: 25th International
 473 Conference on Automated Deduction, Berlin, Germany, August 1-7, 2015, Proceedings 25*, pp.
 474 378–388. Springer, 2015.
- 475 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 476 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
 477 *arXiv preprint arXiv:2407.21783*, 2024.
 478
- 479 Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang, Hengyun Li, Dawei Yin, Tat-Seng Chua, and
 480 Qing Li. A survey on rag meeting llms: Towards retrieval-augmented large language models. In
 481 *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,
 482 pp. 6491–6501, 2024.
- 483 Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
 484 Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng,
 485 Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu,
 Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao,

- 486 Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu,
487 Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan
488 Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang,
489 Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language
490 models from glm-130b to glm-4 all tools, 2024.
- 491 John Harrison. A survey of automated theorem proving. *Sat*, 17:20–18, 2013.
- 492 John Harrison, Josef Urban, and Freek Wiedijk. History of interactive theorem proving.
493 In Jörg H. Siekmann (ed.), *Computational Logic*, volume 9 of *Handbook of*
494 *the History of Logic*, pp. 135–214. North-Holland, 2014. doi: [https://doi.org/10.1016/](https://doi.org/10.1016/B978-0-444-51624-4.50004-6)
495 [B978-0-444-51624-4.50004-6](https://doi.org/10.1016/B978-0-444-51624-4.50004-6). URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/B9780444516244500046)
496 [article/pii/B9780444516244500046](https://www.sciencedirect.com/science/article/pii/B9780444516244500046).
- 497 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
498 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
499 2021.
- 500 Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song,
501 and Denny Zhou. Large language models cannot self-correct reasoning yet. *arXiv preprint*
502 *arXiv:2310.01798*, 2023.
- 503 Gérard Huet, Gilles Kahn, and Christine Paulin-Mohring. The coq proof assistant a tutorial. *Rapport*
504 *Technique*, 178, 1997.
- 505 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang,
506 Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM*
507 *Computing Surveys*, 55(12):1–38, 2023.
- 508 Albert Q Jiang, Sean Welleck, Jin Peng Zhou, Wenda Li, Jiacheng Liu, Mateja Jamnik, Timothée
509 Lacroix, Yuhuai Wu, and Guillaume Lample. Draft, sketch, and prove: Guiding formal theorem
510 provers with informal proofs. *arXiv preprint arXiv:2210.12283*, 2022.
- 511 Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish
512 Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. *arXiv*
513 *preprint arXiv:2210.02406*, 2022.
- 514 Shalom Lappin. Assessing the strengths and weaknesses of large language models. *Journal of Logic,*
515 *Language and Information*, 33(1):9–20, 2024.
- 516 Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ra-
517 masesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam
518 Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with lan-
519 guage models, 2022. URL <https://arxiv.org/abs/2206.14858>.
- 520 Zhaoyu Li, Jialiang Sun, Logan Murphy, Qidong Su, Zenan Li, Xian Zhang, Kaiyu Yang, and Xujie
521 Si. A survey on deep learning for theorem proving. *arXiv preprint arXiv:2404.09939*, 2024.
- 522 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-
523 train, prompt, and predict: A systematic survey of prompting methods in natural language pro-
524 cessing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- 525 Xuantao Lu, Jingping Liu, Zhouhong Gu, Hanwen Tong, Chenhao Xie, Junyang Huang, Yanghua
526 Xiao, and Wenguang Wang. Parsing natural language into propositional and first-order logic with
527 dual reinforcement learning. In *Proceedings of the 29th International Conference on Computa-*
528 *tional Linguistics*, pp. 5419–5431, 2022.
- 529 Filip Maric. A survey of interactive theorem proving. *Zbornik radova*, 18(26):173–223, 2015.
- 530 M Saqib Nawaz, Moin Malik, Yi Li, Meng Sun, and M Lali. A survey on theorem provers in formal
531 methods. *arXiv preprint arXiv:1912.03028*, 2019.

- 540 Thien Huu Nguyen and Ralph Grishman. Relation extraction: Perspective from convolutional neural
541 networks. In *Proceedings of the 1st workshop on vector space modeling for natural language*
542 *processing*, pp. 39–48, 2015.
- 543 Lawrence C Paulson. *Isabelle: A generic theorem prover*. Springer, 1994.
- 544 Mohammadreza Pourreza and Davood Rafiei. Din-sql: Decomposed in-context learning of text-to-
545 sql with self-correction. *Advances in Neural Information Processing Systems*, 36, 2024.
- 546 Boris Shminke. Python client for isabelle server. *arXiv preprint arXiv:2212.11173*, 2022.
- 547 Hrituraj Singh, Milan Aggrawal, and Balaji Krishnamurthy. Exploring neural models for parsing
548 natural language into first-order logic. *arXiv preprint arXiv:2002.06544*, 2020.
- 549 Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen-tse Huang, Pinjia He, Wenxiang
550 Jiao, and Michael R Lyu. A & b== b & a: Triggering logical reasoning failures in large language
551 models. *arXiv preprint arXiv:2401.00757*, 2024.
- 552 Haiming Wang, Huajian Xin, Chuanyang Zheng, Lin Li, Zhengying Liu, Qingxing Cao, Yinya
553 Huang, Jing Xiong, Han Shi, Enze Xie, et al. Lego-prover: Neural theorem proving with growing
554 libraries. *arXiv preprint arXiv:2310.00656*, 2023.
- 555 Qingxiang Wang, Chad Brown, Cezary Kaliszyk, and Josef Urban. Exploration of neural machine
556 translation in autoformalization of mathematics in mizar. In *Proceedings of the 9th ACM SIG-*
557 *PLAN International Conference on Certified Programs and Proofs*, pp. 85–98, 2020a.
- 558 Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. Generalizing from a few examples:
559 A survey on few-shot learning. *ACM computing surveys (csur)*, 53(3):1–34, 2020b.
- 560 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
561 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
562 *neural information processing systems*, 35:24824–24837, 2022.
- 563 Junbing Yan, Chengyu Wang, Jun Huang, and Wei Zhang. Do large language models understand
564 logic or just mimick context? *arXiv preprint arXiv:2402.12091*, 2024.
- 565 Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhan Cui, Zeyang Zhou,
566 Chao Gong, Yang Shen, et al. A comprehensive capability analysis of gpt-3 and gpt-3.5 series
567 models. *arXiv preprint arXiv:2303.10420*, 2023.
- 568 Artem Yushkovskiy. Comparison of two theorem provers: Isabelle/hol and coq. *arXiv preprint*
569 *arXiv:1808.09701*, 2018.
- 570 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
571 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*
572 *preprint arXiv:2303.18223*, 2023.
- 573 Jin Peng Zhou, Charles Staats, Wenda Li, Christian Szegedy, Kilian Q Weinberger, and Yuhuai Wu.
574 Don’t trust: Verify-grounding llm quantitative reasoning with autoformalization. *arXiv preprint*
575 *arXiv:2403.18120*, 2024.

584 A APPENDIX

585 A.1 ADDITIONAL CASE OF MANUAL MODIFICATION FOR STEP PROOF FITTING

586 For the following informal problem and proof:

587 **Problem:** What is the average of the two smallest positive integer solutions to the congruence
588 $14u \equiv 46 \pmod{100}$. Show it is 64.

589 **Solution:** Note that 14, 46, and 100 all have a common factor of 2, so we can divide it out: the
590 solutions to

$$14u \equiv 46 \pmod{100}$$

594 are identical to the solutions to

$$595 \quad 7u \equiv 23 \pmod{50}.$$

596 Make sure you see why this is the case.

597 Now we can multiply both sides of the congruence by 7 to obtain

$$598 \quad 49u \equiv 161 \pmod{50},$$

600 which also has the same solutions as the previous congruence, since we could reverse the step above
601 by multiplying both sides by 7^{-1} . (We know that 7^{-1} exists modulo 50 because 7 and 50 are
602 relatively prime.)

604 Replacing each side of $49u \equiv 161$ by a $\pmod{50}$ equivalent, we have

$$605 \quad -u \equiv 11 \pmod{50},$$

606 and thus

$$607 \quad u \equiv -11 \pmod{50}.$$

608 This is the set of solutions to our original congruence. The two smallest positive solutions are
609 $-11 + 50 = 39$ and $-11 + 2 \cdot 50 = 89$. Their average is $\boxed{64}$.

611 In normal cases, we can cut the content according to the period to get the following sequence of
612 non-formal proofs.

613
614 1 Note that 14 , 46 , and 100 all have a common factor of 2 ,
615 \rightarrow so we can divide it out: the solutions to $14u \equiv 46$
616 $\pmod{100}$ are identical to the solutions to $7u \equiv 23$
617 $\pmod{50}$.

618 2
619 3 Make sure you see why this is the case.

620 4
621 5 Now we can multiply both sides of the congruence by 7 to obtain
622 $\rightarrow 49u \equiv 161 \pmod{50}$, which also has the same
623 \rightarrow solutions as the previous congruence, since we could reverse
624 \rightarrow the step above by multiplying both sides by 7^{-1} .

625 6 We know that 7^{-1} exists modulo 50 because 7 and 50 are
626 \rightarrow relatively prime.

627 7
628 8 Replacing each side of $49u \equiv 161$ by a $\pmod{50}$
629 \rightarrow equivalent, we have $-u \equiv 11 \pmod{50}$, and thus u
630 $\rightarrow \equiv -11 \pmod{50}$.

631 9 This is the set of solutions to our original congruence.

632 10
633 11 The two smallest positive solutions are $-11 + 50 = 39$ and
634 $\rightarrow -11 + 2 \cdot 50 = 89$.

635 12
636 13 Their average is $\boxed{64}$.

637 However, there are many problems with such a simple decomposition. Such as the following sen-
638 tences are unable to be formalized.

640 1 Make sure you see why this is the case.

641 2
642 3 This is the set of solutions to our original congruence.

643
644 In addition, from the point of view of the order of proof, I can see that "We know that 7^{-1} exists
645 modulo 50 because 7 and 50 are relatively prime." is a prerequisite for "Now we can multiply both
646 sides of the congruence by 7 to obtain $49u \equiv 161 \pmod{50}$, which also has the same solutions as
647 the previous congruence, since we could reverse the step above by multiplying both sides by 7^{-1} .",
so the order of the two should be reversed.

648 Therefore, in order to meet the requirements of StepProof, I deleted the statements that could not be
 649 formalized and corrected the sequence to obtain the following proof sequence.
 650

651 1 Note that 14 , 46 , and 100 all have a common factor of 2 ,
 652 \rightarrow so we can divide it out: the solutions to $14u \equiv 46$
 653 $\pmod{100}$ are identical to the solutions to $7u \equiv 23$
 654 $\pmod{50}$.

655 2
 656 3 We know that 7^{-1} exists modulo 50 because 7 and 50 are
 657 \rightarrow relatively prime.

658 4
 659 5 Now we can multiply both sides of the congruence by 7 to obtain
 660 $\rightarrow 49u \equiv 161 \pmod{50}$, which also has the same
 661 \rightarrow solutions as the previous congruence, since we could reverse
 662 \rightarrow the step above by multiplying both sides by 7^{-1} .

663 6
 664 7 Replacing each side of $49u \equiv 161$ by a $\pmod{50}$
 665 \rightarrow equivalent, we have $-u \equiv 11 \pmod{50}$, and thus u
 666 $\rightarrow \equiv -11 \pmod{50}$. This is the set of solutions to our
 667 \rightarrow original congruence.

668 8
 669 9 The two smallest positive solutions are $-11+50 = 39$ and
 670 $\rightarrow -11+2 \cdot 50 = 89$.

671 10
 672 11 Their average is $\boxed{64}$.

672 A.2 FORMAL CONTENT COMPARISON BETWEEN FULL-PROOF AND STEP-PROOF

673
 674 **Problem:** Given that $1 < x < 2$, show that $\sqrt{(x-1)^2} + |x-2|$ can be reduced to 1.

675
 676 **Solution:** Since $x > 1$, we know that $\sqrt{(x-1)^2} = x-1$. And since $x < 2$, we know that
 677 $|x-2| = 2-x$. Therefore, it can be rewritten as $x-1+2-x$, and we know that it can be reduced
 678 to 1.

679 After formalize the problem we get following Isabelle statement.

```
680
681 1 theorem sqrt_abs_simplify:
682 2   assumes h0: "1 < x"
683 3   and h1: "x < 2"
684 4   shows "sqrt((x - 1)^2) + abs(x - 2) = 1"
```

685 Using Full-Proof, we can directly obtain the following formal proof.

```
686
687 1 proof -
688 2   have h2: "sqrt((x - 1)^2) = x - 1"
689 3     using h0 by (simp add: abs_of_nonneg)
690 4   have h3: "abs(x - 2) = 2 - x"
691 5     using h1 by (simp add: abs_of_nonpos)
692 6   have "sqrt((x - 1)^2) + abs(x - 2) = (x - 1) + (2 - x)"
693 7     using h2 h3 by simp
694 8   also have "... = 1"
695 9     by simp
696 10  finally show ?thesis .
697 11 qed
```

698 In the StepProof approach, we would firstly input the first step

699
 700 1 Since $x > 1$, we know that $\sqrt{(x-1)^2} = x - 1$.

701

Then after formalization, we can get corresponding formal proof.

```
702 1 have h2: "sqrt((x - 1)^2) = x - 1" using h0 by simp
```

```
703
```

After this formal step has been verified, we will input the next steps in the same way.

```
704
```

```
706 1 And since  $x < 2$ , we know that  $|x-2| = 2-x$ .
```

```
707
```

```
708 1 have h3: "abs(x - 2) = 2 - x" using h1 by simp
```

```
709
```

```
710 1 Therefore, it can be rewritten as  $x - 1 + 2 - x$ , and we know
711    $\leftrightarrow$  that it can be reduced to 1.
```

```
712
```

```
713 1 show ?thesis using h2 h3 by simp
```

```
714
```

After the global goal has been proven, the system will automatically align the formal proof with the informal proof. Then we will get the over all formal proof of StepProof as following

```
715
```

```
717 1 proof-
```

```
718 2
```

```
  (*Since  $x > 1$ , we know that  $\sqrt{(x-1)^2} = x - 1$ .*)
```

```
719 3
```

```
  have h2: "sqrt((x - 1)^2) = x - 1" using h0 by simp
```

```
720 4
```

```
  (*And since  $x < 2$ , we know that  $|x-2| = 2-x$ .*)
```

```
721 5
```

```
  have h3: "abs(x - 2) = 2 - x" using h1 by simp
```

```
722 6
```

```
  (*Therefore, it can be rewritten as  $x - 1 + 2 - x$ , and
723    $\leftrightarrow$  we know that it can be reduced to 1.*)
```

```
724 7
```

```
  show ?thesis using h2 h3 by simp
```

```
725 8
```

```
726 9 qed
```

```
727
```

```
728
```

```
729
```

```
730
```

```
731
```

```
732
```

```
733
```

```
734
```

```
735
```

```
736
```

```
737
```

```
738
```

```
739
```

```
740
```

```
741
```

```
742
```

```
743
```

```
744
```

```
745
```

```
746
```

```
747
```

```
748
```

```
749
```

```
750
```

```
751
```

```
752
```

```
753
```

```
754
```

```
755
```