
Autoformalization with Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

Autoformalization is the process of automatically translating from natural language mathematics to formal specifications and proofs. A successful autoformalization system could advance the fields of formal verification, program synthesis, and artificial intelligence. While the long-term goal of autoformalization seemed elusive for a long time, we show large language models provide new prospects towards this goal. We make the surprising observation that LLMs can correctly translate a significant portion (25.3%) of mathematical competition problems perfectly to formal specifications in Isabelle/HOL. We demonstrate the usefulness of this process by improving a previously introduced neural theorem prover via training on these autoformalized theorems. Our methodology results in a new state-of-the-art result on the MiniF2F theorem proving benchmark, improving the proof rate from 29.6% to 35.2%.

1 Introduction

Autoformalization refers to the task of automatically translating from natural language mathematics to a formal language [46, 42]. The implication of a successful autoformalization tool is huge in both practical and philosophical terms. It would reduce the currently excessive cost of formalization efforts [27], and in the long-term it could connect the various research fields that automate aspects of mathematical reasoning, such as automated theorem proving and computer algebra, to the vast body of mathematical knowledge exclusively written up in natural language. Moreover, autoformalization would be a true testament to machine understanding, grasping both the fuzziness of natural language and the preciseness of formal language.

Recent advances in large language models [8, 10] showed promising capabilities of understanding formal languages [9, 32]. However, the existing successes are limited to formal languages where there exists a large body of corpus on the web (e.g., Python language). Formal mathematics data is very scarce. For example, one of the largest formal mathematics libraries, the [Archive of Formal Proofs](#), is only 180MB in size, that is less than 0.18% of the training data for the large language model Codex [9]. Moreover, unlike in the case of commonly used programming languages, where natural language docstrings are broadly available, there is almost zero aligned data between natural language and formal mathematics. Therefore, it is unclear the recent successes can directly contribute to the development of autoformalization.

In this work, we explore the prospects of autoformalization with large language models. To our surprise, we find that large language models already have a decent capability of formalizing natural language mathematics in an interactive theorem prover. See Figure 1 for a perfect autoformalization example. The model not only translates into syntactically correct Isabelle code, but also grasps the non-trivial reasoning in natural language. We randomly pick 150 formalizations and manually evaluate their correctness. Among them, LLMs are capable of producing 38 perfect formalizations! As an application, we further demonstrate that autoformalization can provide useful training data for neural theorem provers. We use autoformalized statements as targets for proof search with a

Case Study 1 Question:

"Prove that there is no function f from the set of non-negative integers into itself such that $f(f(n)) = n + 1987$ for every n ."

Codex Output:

```
theorem
  fixes f :: "nat \ $\rightarrow$  nat"
  assumes "\ $\forall$  n. f (f n) = n + 1987"
  shows False
```

Figure 1: Case study 1: An example of a perfect translation from natural language to Isabelle code.

neural theorem prover for Isabelle/HOL. After fine-tuning our neural theorem prover on the proofs it found, its success rate on the MiniF2F benchmark [52] increases significantly, achieving a new state-of-the-art result of 35.2% theorems proven.

2 Related Work

Early applications of machine learning in theorem proving include the works by Schulz [40] and Urban [43], and later, directly guiding interactive proof assistants using machine learning techniques [15]. The revolution of deep learning then kicked off a new wave of interest in the topic starting with DeepMath [1, 33].

Several approaches have been suggested to address data scarcity: Imitation-free reinforcement learning was used to avoid the need for training on human proofs [31, 6, 15, 49]. Also, hindsight experience replay [2] was used to generate additional training data [5]. Hahn et al. [19], Schmitt et al. [39], Kreber & Hahn [29] and Wu et al. [50] have shown that training on synthetic formulas can be successful for temporal logics and inequalities. Rabe et al. [37] masked out different subexpressions from formal mathematical statements and generated 100 training examples for each source statement. Skip-tree data can also be used to improve the performance of neural theorem provers [22].

Wang et al. [46] explored the use of supervised and unsupervised translation techniques for autoformalization. Supervised translation yielded interesting results, but relied on synthetic (natural-looking) data that was generated by the Mizar theorem prover, while we rely on models trained via self-supervised language modeling, not trained for this particular purpose.

3 Background

Formal Mathematics A few important and complex results of mathematics and computer science have been formalized manually using *interactive theorem provers*, such as the four color theorem [16], the Kepler conjecture [20], the odd-order theorem [17] and the verification of a microkernel [27]. This gives us almost complete certainty about the correctness of proofs, which can be of great value to resolve doubt about the correctness of complicated mathematical proofs or proving certain properties of software used in safety-critical applications, such as aircraft components [28].

These projects relied on interactive theorem provers, such as Isabelle [48], Coq [12], HOL Light [23], and Lean [13], which are essentially programming languages that enable users to enter their statements and proofs in a formal language, and which can then be checked automatically for correctness. Interactive theorem provers offer a limited amount of automation, but projects that formalize complex problems typically span many years of tedious work by specialists. Only in narrow domains like chip design and the verification of drivers in operating systems has the automation of logic made sufficient progress to find commercial applications.

Progress in autoformalization and the automation of proofs might eventually make mathematics a universally available tool and enable a paradigm shift in science and the development of (safety-critical) software. Our interest in formalizing mathematics, however, has an additional aspect. We believe that autoformalization will serve a dual purpose and will not only accelerate the development

76 of tools for mathematical reasoning, but also provide a means to ground machine learning systems,
77 enabling a positive feedback loop between machine learning and formal systems (cf. [42]).

78 **Large Language Models** Our work relies heavily on large language models (LLMs), in particular
79 on PaLM [10] and Codex [9]. The training goal of these models is to predict the next word given
80 some prefix. This allows us to train these models on arbitrary text, which is available in vast quantities.
81 After training the models on hundreds of billions of words (cf. [25]), they are often able to generate
82 high-quality text. We can also give these models an arbitrary prefix (the *prompt*) that they are
83 then supposed to continue, which gives us some control over what they generate. This has been
84 demonstrated with news articles, conversations, summaries, jokes, and poems. LLMs have also been
85 evaluated on natural language word problems on datasets such as GSM8K [11] and MATH [24], and
86 have been shown to make progress on these benchmarks with increasing scale [10].

87 **In-context Learning** Large language models have shown a remarkable ability to learn patterns
88 and tasks within the current input (context) that they are given [8]: this is called *in-context learning*
89 or *few-shot learning*. For example, if we prompt a language model with a few pairs of English and
90 matching French sentences, and end with a new English sentence, then the language model is very
91 likely to pick up on the translation task and attempt a translation of the last English sentence. This
92 observation has been used, for example, to achieve strong translation performance without access to
93 large corpora of matching sentence pairs [21].

94 This allows us to specify the task of autoformalization simply by giving a couple of example formal-
95 izations. In Section 4 we will detail how exactly we use in-context learning for autoformalization.

96 4 Autoformalization for Mathematical Competition Problems

97 Inspired by the success of LLMs for synthesizing computer code by co-training on both natural
98 language and code on web-scale data, we explore the capabilities of LLMs to turn natural language
99 mathematics into formalized theorems for the interactive theorem prover Isabelle. This can be seen
100 as a machine translation task (cf. [47]) in which the input language is English and output language is
101 formal code used by the interactive proof assistant Isabelle [48].

102 We first study autoformalization in a constrained setting – formalizing mathematical competition
103 problem statements. This setting has the advantage that most of the required background theory and
104 definition has been formalized in the current libraries of Isabelle, so that formalizations are often
105 possible without introducing additional definitions.

106 We start assessing LLMs’ abilities to do autoformalization with a case study. We manually pick
107 two interesting natural language mathematical statements, and prompt PaLM models of various
108 scales [10] as well as Codex [9] to translate them into a formal statement in Isabelle. Next, we
109 study a dataset in which we have human ground truth formalizations. The dataset is a subset of the
110 miniF2F [24] dataset consisting of 140 algebra problems and 120 number theory problems. Using
111 human formalizations as the reference, we compute the BLEU scores of the formalizations produced
112 by several LLMs. Lastly, we perform human evaluations on failure cases in autoformalization on 150
113 problems.

Note that many mathematical competition statements are often of the form in which one asks to
find the answer to a certain problem, instead of *proving* a given proposition. However, formal
mathematical statements are in the form of propositions, instead of questions. To transform a question
into a proposition, we append the final answer after the question:

\$Problem_Statement The final answer is \$Answer.

114 The format of the prompt we use to do autoformalization is:

Natural language version: \$Natural_Language_Statement.

Translate the natural language version to an Isabelle version:

115 4.1 Mathematical Competition Datasets

116 **MATH** [24] contains in total 12,500 (7,500 training and 5,000 test) middle school and high
117 school mathematical competition problems. Problems are taken from past mathematical com-

petitions, including AMC 10, AMC 12, AIME, and more, and many can be found at http://aops.com/community/c3158_usa_contests. The dataset contains seven categories: algebra, pre-algebra, intermediate algebra, number_theory, precalculus, probability, geometry. Problem statements are written in LaTeX.

MiniF2F [52] is a recently introduced benchmark containing 488 mathematical competition statements manually formalized by humans in three different formal languages. Its goal is to compare and benchmark methods across different theorem provers for machine learning research. Some of these problems come from the valid and test set of MATH algebra and number_theory, and others come from previous International Mathematical Olympiad competitions or AoPS¹. Note that the Isabelle formalizations of the miniF2F benchmark were committed to the repository during March, 2022. According to the public information of the training data, we think it is highly unlikely these formalizations were included in the pre-training corpus.

4.2 Case Studies

Experimental setup For all our experiments, we use the standard greedy decoding (i.e., temperature 0, $p = 1$) to obtain the autoformalizations. We randomly select two mathematical statements for constructing the prompt, which we provide in Appendix A.1. That is, no prompt engineering / tuning is performed when constructing the prompt. The natural language problem statements used in the case studies are taken from the miniF2F dataset. In the case studies below, we highlight the output of language models in red to distinguish it from the prompt.

Case Study 1 (Figure 1) We study the example shown in Figure 1, in which we ask LLMs to autoformalize an International Mathematical Olympiad problem² in natural language. Surprisingly, Codex is able to autoformalize the natural language statement as an Isabelle theorem perfectly, with output given. This is surprising for the following reasons.

First of all, the amount of Isabelle code is very scarce on the internet. The entire AFP library, the largest formal library that contains most of Isabelle proofs, is only 180MB in size. Even assuming that all of this data was included in the training of Codex, this makes at most 0.18% of the pretraining data on which Codex was trained. The fact that the model can write syntactically correct Isabelle code at all is already fascinating.

Second, there is almost zero aligned data from natural language to Isabelle on the web. While some Isabelle files have comments, they typically only give a very high level description of what the theory being formalized is about. So either LLMs are able to transfer knowledge quite successfully between natural language and formal mathematics, or the task was learned mostly via few-shot learning.

Last but not least, the model is capable of understanding and formalizing nontrivial reasoning. First, the model is able to formalize the non-existence statement via proof-by-contradiction. To formalize “there is no function $f...$ ”, it assumes there is such a function, and aims to prove “False”. Second, the model understands what it means by the phrase “to itself”, and correctly infers the domain of function: $f :: \text{nat} \rightarrow \text{nat}$.

On the other hand, PaLM made some syntactic mistakes while getting most of the structure of the proof correctly, with outputs shown in Appendix C.1.

Case Study 2 (Figure 2) In the next example, we ask LLMs to autoformalize a grade school mathematical word problem. Remarkably, PaLM and Codex are both capable of formalizing the statement perfectly. This is surprising because formalizations of grade school math problems in interactive theorem provers are rare (if they exist at all), as this type of mathematics is not of interest to formal mathematicians. Even more, none of the examples in the prompt (see Appendix A.1) that we provide are of this type. It is hence remarkable that the model is capable of extrapolating to this type of statement, showing a great promise of using LLMs for autoformalization.

To study this problem in more depth, we probe PaLM models of various sizes (8B, 64B, 540B) with outputs shown in Appendix C.2, and notice that scale is crucial for the LLMs ability to formalize. We

¹<https://artofproblemsolving.com/>

²A problem from IMO 1987.

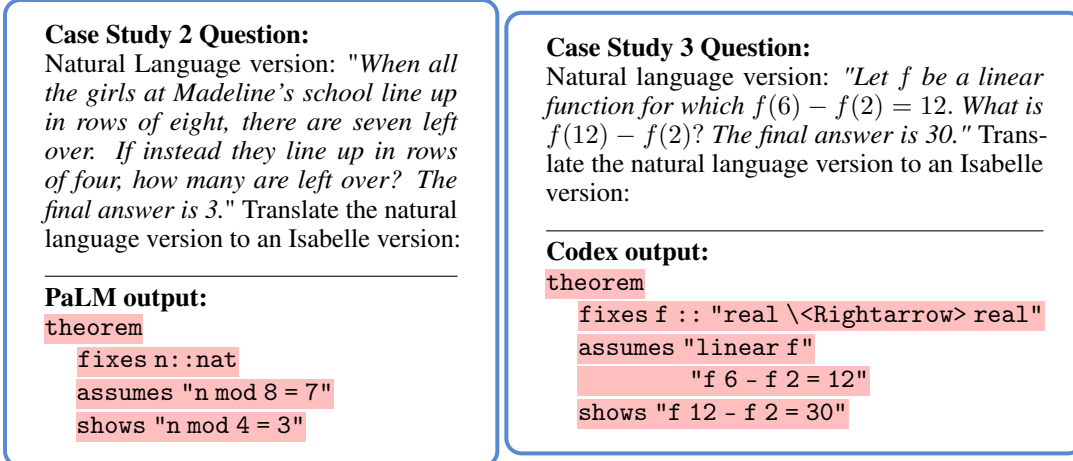


Figure 2: Autoformalizations from natural language to Isabelle code. **Left:** Case study 2 – perfect formalization by PaLM. **Right:** Case study 3 – incorrect formalization by Codex.

observe that the 8B and 64B models are incapable of formalizing this problem, but the largest 540B model is able to produce a correct formalization.

Case Study 3 (Figure 2) In our third case study, Codex gives an incorrect formalization in Isabelle. The mathematical statement involves a concept of “linear function”, which the model fails to formalize correctly. Codex assumes this is already a known concept in Isabelle, and made up a name: `linear f`. Can the model learn to formalize such problems if the prompt contains an example that explains the concept of a line? We explore this and give an affirmative answer to the question (see Appendix C.3). Once seeing a tangentially related problem that explains the concept of a “line”, Codex is able to perfectly formalize a “linear function”. This shows the importance of the few shot examples we include, and also how good a few-shot learners these models are!

Has the model memorized these formalizations? Whilst we do not have access to the training set of Codex, we attempted to find any occurrences of the formalizations produced in the case studies on the internet. We Googled them in different variants and inspected the first page of the search results. We tried variants with and without an “Isabelle” prefix, with and without quotation marks and other special characters, and also individual parts of it, such as “Isabelle `n mod 8 = 7`”, but we did not find any occurrences of related statements. We also tested that we are indeed able to find occurrences of Isabelle formalizations on the web with this methodology, using pieces of formalizations picked from several websites, including the Archive of Formal Proofs. Hence, we are confident that the model has not memorized the formalizations it generated.

4.3 BLEU for Model Comparisons

The miniF2F benchmark contains 140 algebra problems and 120 number theory problems from the MATH dataset. For these problems, we have human ground truth formalizations in Isabelle, which gives us an evaluation set with pairs of natural language statements (from MATH) and their formalizations. We use this dataset to quantitatively compare different LLMs.

Given the observation about few shot learning in Case study 3, we decided to add more relevant examples to each subject to improve the quality of autoformalization. For each subject (i.e., algebra and number_theory), we randomly sample 10 problems to construct the few shot prompt. The rest of the problems are used for evaluation (i.e., 130 for algebra and 110 for number_theory). We provide the prompt used in the Appendix B.1 and B.2.

We use PaLM models of varying sizes and Codex to perform the autoformalization, and compute the BLEU scores of the formalizations, shown in Table 1. Confirming our observation in Case study 2, we see a clear trend that scaling improves translation, as the BLEU scores consistently improve when we scale PaLM models from 8B to 540B, for both subjects. In addition, we see that the Codex model

Table 1: BLEU scores between the autoformalized statements and human formalized ground truth.

Models \ Subject	algebra	number_theory
PaLM 8B	31.49	22.10
PaLM 64B	43.13	31.43
PaLM 540B	50.30	36.16
Codex	57.13	43.33

Table 2: Failure case study of 150 problems formalized by Codex.

Failure cases \ Subjects	algebra	number_theory	inter_alg
Perfect translation	13	17	8
Incomplete/ill-formed/unclear prompt	9	3	14
Fail to align definitions or concepts	10	18	18
Inconsistent/missing assumption	8	9	9
Syntactical/type error	7	2	11
Missing definition in Isabelle	0	12	3
Wrong application of functions	6	13	16
Other	6	2	1

is better at autoformalization measured by BLEU, possibly due to the fact that Codex was trained on more formal data than PaLM.

4.4 Human Evaluation of Failure cases

To better understand LLMs’ ability to do autoformalization, we manually inspect Codex’s autoformalizations of 150 random problems from the MATH dataset [24]. 50 of the problem statements are sampled from the algebra training set, 50 from number_theory and 50 from intermediate_algebra. For algebra and number_theory, we use their corresponding prompt as in the last section, shown in Appendix B.1 and B.2. For intermediate_algebra, we use the prompt we used for algebra (Appendix B.1). We classify the failure modes of these translations, shown in Table 2.

We see that out of 150 problems, Codex is capable of translating 38 problems perfectly – a success rate of 25.3%. The majority of the failures are due to the misalignment of informal and formal definitions. For example, when seeing the phrase “the greatest possible value”, the LLMs often fail to align it with the function Greatest/Max in Isabelle. Another example is the failure to align the factorial of n (i.e., $!n$) to $fact\ n$ in Isabelle. Other common failure modes include the misapplication of functions (e.g., applying a prefix function in an infix way).

5 Autoformalization for Neural Theorem Proving

To demonstrate the usefulness of the formalized statements, we explore if one can improve neural theorem provers by training the neural models on proofs of automatically translated theorems. In this section, we combine autoformalization with expert iteration algorithms [4], and achieve a new state of the art in miniF2F benchmark.

5.1 Expert Iteration with Autoformalization

The basic idea of expert iteration [4] is to iteratively generate a better dataset using the model, and use the data to improve the model quality. This allows the model to generate an even better quality of the dataset and hence a better model, forming a self-improvement cycle.

In neural theorem proving, one way to get better quality data is to use feedback from the proof checker to run many proof searches (or generate multiple proofs) and check the proof attempts for correctness. Newly found correct proofs can then be used as the new training data to improve the neural prover [7, 35, 36]. The main critical ingredient that is needed is a set of problem statements on which the model can perform proof search to obtain new training data. However, unlike in Polu et al.

Table 3: Proof success rates on miniF2F.

Model	<i>valid</i>	<i>test</i>
PACT [22]	23.9%	24.6%
FMSCL [36]	33.6%	29.6%
Base model (M_0)	28.3%	29.9%
After 1 expert iteration (M_1)	36.1%	34.0%
After 2 expert iterations (M_2)	37.3%	35.2%

[36], where one asks humans to manually formalize a set of problems to get formal statements, here we use LLMs to autoformalize the theorems in order to kick off the self-improvement cycle.

More formally, denote a base neural theorem prover as M_0 . Let the set of autoformalized problems be \mathcal{A} . For each iteration $i = 1 \dots N$, we carry out the following procedure: use the language model M_{i-1} with best-first search to prove as many theorems as possible in \mathcal{A} , collect the set of successful proofs S_i , concatenate successful problems from all iterations with the formal mathematics problems to create the set $\mathcal{A}_i = (\bigcup_{j \leq i} S_j) \cup \mathcal{B}$, and fine-tune M_0 on it for exactly one epoch to get a new model M_i . When we take the union of successful proofs from all past iterations, we perform deduplication by problem statements, similar to Polu et al. [36].

5.2 Neural Theorem Provers

To demonstrate the effectiveness of the approach, we start with a recently introduced neural theorem prover for Isabelle, LISA [26]. The LISA agent is fine-tuned on the PISA dataset [26] (extraction and interaction code under a [BSD license](#)), which consists of 2.49 million proof steps from the Isabelle/HOL library (under a [BSD-style license](#)) and the [Archive of Formal Proofs](#) (under various licenses as described [here](#)). The model is trained with the objective to predict the next token in a proof step, given the proof state and the last proof step. [Following the setup of Thor \[3\], which achieves SOTA performance in the no-additional-training-data category on MiniF2F](#), we invoke Sledgehammer with a 30 second timeout when the model generates a step containing any of the keywords `metis`, `meson`, and `smt`.

We use Wang [45]’s implementation (under an [Apache license 2.0](#)) of a GPT-2 [38] style decoder-only transformer [44] model with 700M non-embedding parameters. The model has 24 layers, 24 attention heads, a hidden dimension of 1536, and a vocabulary size of 50400. We pre-train the model on the GitHub + arXiv subsets of The Pile [14] for 200,000 steps, with a context length of 2048 tokens. In pre-training we use a warmup strategy [18], raising the learning rate linearly from 0 to 2×10^{-4} in 3,000 steps. We then use a cosine learning rate scheduler [34] for the rest of the pre-training, with a final learning rate of 1.2×10^{-5} . We use a global batch size of 32 sequences, or 65,536 tokens. For fine-tuning we use the same learning rate schedule, with 10,000 warmup steps, 90,000 annealing steps, maximum learning rate 3×10^{-4} and final learning rate 3×10^{-5} . The global batch size is 144 sequences, or 294,912 tokens. The model’s evaluation loss reaches a minimum after 13,000 steps and we use that checkpoint. [For the fine-tuning in expert iteration, we fix the learning rate at \$3 \times 10^{-4}\$ and the batch size at 144 sequences, and train the model for exactly one epoch.](#)

[Here we give the details regarding the best-first search strategy used in evaluation: we maintain a priority queue of search nodes with queue length 32. The accumulated log probability of the proof steps is used as the queue priority. For each theorem to prove, we first initialize the queue with one node that has the theorem declared and no proof step applied to it. At each time-step, we deque and sample 32 possible proof steps to apply to the node. The nodes corresponding to steps that successfully proceed the proofs then get added to the queue. We repeat this process until the theorem is successfully proven, or we reach our computational budget.](#)

Machine specification For pre-training, fine-tuning, and evaluation, we use a TPUv3 with 8 cores from [Google Cloud Platform](#). The Isabelle process has access to up to 32 CPU cores. We estimate that running all the experiments in this paper requires a total of 780 TPU hours.

Question: *An arbitrary product of compact spaces is compact.*

Codex Output:

```
theorem
  fixes X :: "'a set set"
  assumes "\<forall>x \<in>X. compact x"
  shows "compact (\<Prod> x \<in> X. x)"
```

Figure 3: A formalization for an advanced mathematical statement by Codex.

5.3 Result

We use Codex with greedy decoding to formalize 3908 mathematical problems in algebra, intermediate algebra, and number theory from the training set of MATH [24], with the same few shot prompts used in Section 4.4. Out of them, 3363 of the autoformalized theorems are syntactically correct. We then perform expert iteration on this dataset.

We start with a neural theorem prover (M_0) as described in Section 5.2. In our first iteration, M_0 proves 782 theorems, with a success rate of 23.3% (out of 3363). This gives us a new set of verified proofs to further train the neural theorem prover. We proceed to fine-tune our neural theorem prover in the fashion described in Section 5.1 to get a new prover (M_1). This process is repeated in the second iteration, giving us 1011 successful proofs from the autoformalized theorems (30.1%). We fine-tuned M_0 again on all available deduplicated proofs to obtain M_2 .

After each stage of fine-tuning, we evaluate the neural theorem prover on miniF2F [52]. The results are shown in Table 3. The base model (M_0) has a success rate of 28.3% and 29.9% on the validation and test fractions of miniF2F respectively. It can be observed that the first expert iteration increases the success rate of the neural prover by 7.8% and 4.1% to 36.1% and 34.0% on the valid and test sets. The second iteration further improves them both by 1.2%, to 37.3% and 35.2%. By doing two expert iterations on the autoformalized theorems, the neural prover achieves a success rate that is 5.6% higher than the previous state-of-the-art.

6 An Outlook on Autoformalizing Advanced Mathematics

So far, we focused on mathematical competition problems, in which we achieve significant results using autoformalization. Not only can LLMs autoformalize non-trivial theorems, the autoformalized theorems can also improve neural prover performance. In this section, we take a peek into more advanced mathematics. We hope to identify some of the limitations of our methods when it comes to autoformalizing mathematics in the wild.

6.1 Autoformalization: from Natural Language to Isabelle

Autoformalization is extremely challenging in the sense that the model needs to (1) bridge the logical gaps left in pen-and-paper proofs, (2) assume the implicit contexts and assumptions, and (3) align informal definitions/concepts to formal ones. This task is further complicated considering the context can change from time to time and the same mathematical object can be formalized in subtly different ways. In this paper, we only deal with autoformalization of theorem statements, where the model is mainly challenged by definition alignments. Compared to competition problems in Section 4.4, advanced mathematical statements assume more context information that additionally obscures the alignment.

Figure 3 shows a typical case where the model fails to align the informal definition to the formal one. Based on the notation used, the model does not correctly distinguish between products of numbers, products of sets, and products of topological spaces. In Isabelle, the correct conclusion for the statement should be *compact_space (product_topology X I)*, where I is an index set that ought to have been introduced in the *fixes* and *assumes* sections.

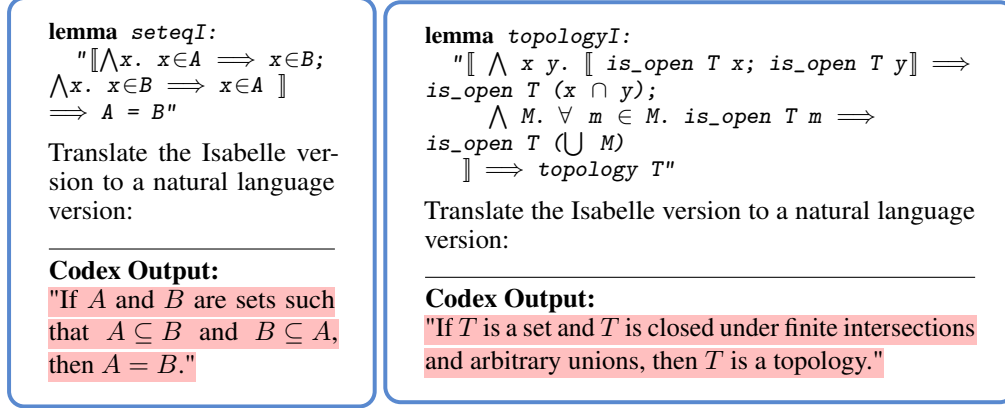


Figure 4: Two perfect translations from Isabelle code to natural language by Codex.

6.2 Informalization: from Isabelle to Natural Language

So far, we explored one direction of translation: from natural language statement to formal statement in Isabelle. The other direction, *informalization*, is also of great importance for two reasons: (1) the informal texts are much easier for humans to comprehend and communicate, and (2) we can align translated informal statements with formal ones to create data, and use the back-translation techniques [41] to potentially boost the translator’s performance further. In this section, we explore Codex’s capability of translating formal Isabelle statement to natural language.

A corpus of 38 formal-language theorems, lemmas, and definitions is selected by an Isabelle expert. These statements are automatically translated to informal mathematics using Codex; to see the prompt we used and the results for all 38 examples, see Appendix B.3 and E.2. We present two examples of informalization in Figure 4. Of the 38 examples, 36 were translated to a reasonably coherent statement, and 29 of these statements (76%) were more-or-less correct, giving a vastly better success rate than the 25% success rate of formalization (Section 4.4). Our main conclusion is that for advanced mathematics, the model is better at informalization than formalization, showing the prospect of backtranslation style algorithms.

Note that the standard is more relaxed here since we assume a human reader will supply the obvious context and correct mistakes when the intended meaning is obvious (intended by the hypothetical human writer of these sentences). To illustrate, an example of a minor “acceptable” error: assuming that “ w, z are in the same connected component of the plane” when, in context, it is clear that w, z should be assumed to be in the same connected component of the complement of a previously specified curve. (The assumption as originally stated is trivial.) For an example of a major error: almost-perfect translation of the Central Limit Theorem that omits the assumption of identical distributions.

7 Discussion

Promise of Autoformalization with LLMs We have seen that automated formalization of informally given natural language statements is generally possible, even with language models not trained for this particular task. Also, automatically formalized statements are useful for training and improving the reasoning capabilities of automated neural provers. Our hope is that improved versions of this methodology will be capable of enabling a positive feedback loop involving formalization and formal reasoning that has the potential of reaching human level capabilities in both respects, as was suggested by [42].

Limitations and future directions We use a static model for the formalization process. For large-scale autoformalization, we will need to formalize larger theories, preferably without fine tuning the model, as training it could be cumbersome and resource consuming. However, in order to utilize the newly added notions, the model would need to keep whole large theories in the current context window, which exceeds those of the current LLMs. This limits our approach to the generation of fairly small pieces of formal mathematics and the automatic formalization of entire theories including

their definitions will require new research ideas. One path towards this goal might be the use of continuous training or expert iteration, cycle-consistency-based training [30, 46], or novel uses of in-context learning. To generate larger theories we will also need neural networks that can recall longer sequences (current LLMs are typically limited to a few thousand words). Retrieval-augmented language models, such as the memorizing transformer [51] offer one path to overcome this limitation.

Societal Impact While the potential of creating negative societal impact through formalizations is small, the use of LLMs always comes with risks. For example, for deploying an autoformalization tool using LLMs we would need to consider the inclusivity of variable and lemma names, and of the attribution of scientific ideas.

References

- [1] Alexander A. Alemi, François Chollet, Niklas Eén, Geoffrey Irving, Christian Szegedy, and Josef Urban. Deepmath - deep sequence models for premise selection. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 2235–2243, 2016. URL <http://papers.nips.cc/paper/6280-deepmath-deep-sequence-models-for-premise-selection>.
- [2] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *CoRR*, abs/1707.01495, 2017. URL <http://arxiv.org/abs/1707.01495>.
- [3] Anon Anonymous. Thor: Wielding hammers to integrate language models and automated theorem provers. 2022.
- [4] Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5360–5370, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/d8e1344e27a5b08cdfd5d027d9b8d6de-Abstract.html>.
- [5] Eser Aygün, Laurent Orseau, Ankit Anand, Xavier Glorot, Vlad Firoiu, Lei M. Zhang, Doina Precup, and Shibl Mourad. Proving theorems using incremental learning and hindsight experience replay. *CoRR*, abs/2112.10664, 2021. URL <https://arxiv.org/abs/2112.10664>.
- [6] Kshitij Bansal, Sarah M. Loos, Markus N. Rabe, and Christian Szegedy. Learning to reason in large theories without imitation. *CoRR*, abs/1905.10501, 2019. URL <http://arxiv.org/abs/1905.10501>.
- [7] Kshitij Bansal, Sarah M. Loos, Markus N. Rabe, Christian Szegedy, and Stewart Wilcox. Holist: An environment for machine learning of higher order logic theorem proving. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 454–463. PMLR, 2019. URL <http://proceedings.mlr.press/v97/bansal19a.html>.
- [8] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul

- 394 Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke
395 Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad
396 Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias
397 Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex
398 Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain,
399 William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant
400 Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie
401 Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and
402 Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374,
403 2021. URL <https://arxiv.org/abs/2107.03374>.
- 404 [10] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
405 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,
406 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay,
407 Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope,
408 James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke,
409 Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant
410 Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek
411 Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal,
412 Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor
413 Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou,
414 Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy
415 Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language
416 modeling with pathways. *CoRR*, abs/2204.02311, 2022. doi: 10.48550/arXiv.2204.02311. URL
417 <https://doi.org/10.48550/arXiv.2204.02311>.
- 418 [11] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Jacob Hilton, Reiichiro Nakano, Christo-
419 pher Hesse, and John Schulman. Training verifiers to solve math word problems. *CoRR*,
420 abs/2110.14168, 2021. URL <https://arxiv.org/abs/2110.14168>.
- 421 [12] Coq. The Coq Proof Assistant. <http://coq.inria.fr>. URL <http://coq.inria.fr>.
- 422 [13] Leonardo Mendonça de Moura, Soonho Kong, Jeremy Avigad, Floris van Doorn, and Jakob
423 von Raumer. The lean theorem prover (system description). In Amy P. Felty and Aart Middeldorp
424 (eds.), *Automated Deduction - CADE-25 - 25th International Conference on Automated*
425 *Deduction, Berlin, Germany, August 1-7, 2015, Proceedings*, volume 9195 of *Lecture Notes in*
426 *Computer Science*, pp. 378–388. Springer, 2015. doi: 10.1007/978-3-319-21401-6_26. URL
427 https://doi.org/10.1007/978-3-319-21401-6_26.
- 428 [14] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
429 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The Pile:
430 An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*,
431 2020.
- 432 [15] Thibault Gauthier, Cezary Kaliszyk, Josef Urban, Ramana Kumar, and Michael Norrish.
433 Tactictoe: Learning to prove with tactics. *J. Autom. Reason.*, 65(2):257–286, 2021. doi:
434 10.1007/s10817-020-09580-x. URL <https://doi.org/10.1007/s10817-020-09580-x>.
- 435 [16] Georges Gonthier. The four colour theorem: Engineering of a formal proof. In Deepak Kapur
436 (ed.), *Computer Mathematics, 8th Asian Symposium, ASCM 2007, Singapore, December 15-17,*
437 *2007. Revised and Invited Papers*, volume 5081 of *Lecture Notes in Computer Science*, pp. 333.
438 Springer, 2007. doi: 10.1007/978-3-540-87827-8_28.
- 439 [17] Georges Gonthier, Andrea Asperti, Jeremy Avigad, Yves Bertot, Cyril Cohen, François Garillot,
440 Stéphane Le Roux, Assia Mahboubi, Russell O’Connor, Sidi Ould Biha, Ioana Pasca, Laurence
441 Rideau, Alexey Solovyev, Enrico Tassi, and Laurent Théry. A machine-checked proof of the
442 odd order theorem. In Sandrine Blazy, Christine Paulin-Mohring, and David Pichardie (eds.),
443 *Interactive Theorem Proving - 4th International Conference, ITP 2013, Rennes, France, July*
444 *22-26, 2013. Proceedings*, volume 7998 of *Lecture Notes in Computer Science*, pp. 163–179.
445 Springer, 2013. doi: 10.1007/978-3-642-39634-2_14.

- [18] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- [19] Christopher Hahn, Frederik Schmitt, Jens U. Kreber, Markus N. Rabe, and Bernd Finkbeiner. Teaching temporal logics to neural networks. In *ICLR*, 2021.
- [20] Thomas Hales, Mark Adams, Gertrud Bauer, Tat Dat Dang, John Harrison, Hoang Le Truong, Cezary Kaliszyk, Victor Magron, Sean McLaughlin, Tat Thang Nguyen, et al. A formal proof of the Kepler conjecture. In *Forum of Mathematics, Pi*, volume 5, pp. e2. Cambridge University Press, 2017.
- [21] Jesse Michael Han, Igor Babuschkin, Harrison Edwards, Arvind Neelakantan, Tao Xu, Stanislas Polu, Alex Ray, Pranav Shyam, Aditya Ramesh, Alec Radford, and Ilya Sutskever. Unsupervised neural machine translation with generative language models only. *CoRR*, abs/2110.05448, 2021. URL <https://arxiv.org/abs/2110.05448>.
- [22] Jesse Michael Han, Jason Rute, Yuhuai Wu, Edward Ayers, and Stanislas Polu. Proof artifact co-training for theorem proving with language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=rpxJc9j04U>.
- [23] John Harrison. HOL Light: A tutorial introduction. In Mandayam K. Srivas and Albert John Camilleri (eds.), *Formal Methods in Computer-Aided Design, First International Conference, FMCAD '96, Palo Alto, California, USA, November 6-8, 1996, Proceedings*, volume 1166 of *Lecture Notes in Computer Science*, pp. 265–269. Springer, 1996.
- [24] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. *CoRR*, abs/2103.03874, 2021. URL <https://arxiv.org/abs/2103.03874>.
- [25] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022. doi: 10.48550/arXiv.2203.15556. URL <https://doi.org/10.48550/arXiv.2203.15556>.
- [26] Albert Q. Jiang, Wenda Li, Jesse Michael Han, and Yuhuai Wu. Lisa: Language models of isabelle proofs. *6th Conference on Artificial Intelligence and Theorem Proving*, 2021.
- [27] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. seL4: formal verification of an OS kernel. In Jeanna Neeffe Matthews and Thomas E. Anderson (eds.), *Proceedings of the 22nd ACM Symposium on Operating Systems Principles 2009, SOSP 2009, Big Sky, Montana, USA, October 11-14, 2009*, pp. 207–220. ACM, 2009. doi: 10.1145/1629575.1629596.
- [28] Gerwin Klein, June Andronick, Matthew Fernandez, Ihor Kuz, Toby C. Murray, and Gernot Heiser. Formally verified software in the real world. *Commun. ACM*, 61(10):68–77, 2018. doi: 10.1145/3230627. URL <https://doi.org/10.1145/3230627>.
- [29] Jens U. Kreber and Christopher Hahn. Generating symbolic reasoning problems with transformer GANs. *CoRR*, abs/2110.10054, 2021. URL <https://arxiv.org/abs/2110.10054>.
- [30] Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. Unsupervised machine translation using monolingual corpora only. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=rkYTTf-AZ>.
- [31] Gil Lederman, Markus N. Rabe, Sanjit Seshia, and Edward A. Lee. Learning heuristics for quantified boolean formulas through reinforcement learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=BJluxREKDB>.

- [32] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *DeepMind*, 2022.
- [33] Sarah M. Loos, Geoffrey Irving, Christian Szegedy, and Cezary Kaliszyk. Deep network guided proof search. In Thomas Eiter and David Sands (eds.), *LPAR-21, 21st International Conference on Logic for Programming, Artificial Intelligence and Reasoning, Maun, Botswana, May 7-12, 2017*, volume 46 of *EPiC Series in Computing*, pp. 85–105. EasyChair, 2017. doi: 10.29007/8mwc. URL <https://doi.org/10.29007/8mwc>.
- [34] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [35] Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. *CoRR*, abs/2009.03393, 2020. URL <https://arxiv.org/abs/2009.03393>.
- [36] Stanislas Polu, Jesse Michael Han, Kunhao Zheng, Mantas Baksys, Igor Babuschkin, and Ilya Sutskever. Formal mathematics statement curriculum learning, 2022. URL <https://arxiv.org/abs/2202.01344>.
- [37] Markus N. Rabe, Dennis Lee, Kshitij Bansal, and Christian Szegedy. Mathematical reasoning via self-supervised skip-tree training. In *ICLR*, 2021.
- [38] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [39] Frederik Schmitt, Christopher Hahn, Markus N. Rabe, and Bernd Finkbeiner. Neural circuit synthesis from specification patterns. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 15408–15420, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/8230bea7d54bcdf99cdf99cfe85cb07313d5-Abstract.html>.
- [40] Stephan Schulz. Learning search control knowledge for equational theorem proving. In Franz Baader, Gerhard Brewka, and Thomas Eiter (eds.), *KI 2001: Advances in Artificial Intelligence, Joint German/Austrian Conference on AI, Vienna, Austria, September 19-21, 2001, Proceedings*, volume 2174 of *Lecture Notes in Computer Science*, pp. 320–334. Springer, 2001. doi: 10.1007/3-540-45422-5_23.
- [41] Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics, 2016. doi: 10.18653/v1/p16-1009. URL <https://doi.org/10.18653/v1/p16-1009>.
- [42] Christian Szegedy. A promising path towards autoformalization and general artificial intelligence. In Christoph Benzmüller and Bruce R. Miller (eds.), *Intelligent Computer Mathematics - 13th International Conference, CICM 2020, Bertinoro, Italy, July 26-31, 2020, Proceedings*, volume 12236 of *Lecture Notes in Computer Science*, pp. 3–20. Springer, 2020. doi: 10.1007/978-3-030-53518-6_1. URL https://doi.org/10.1007/978-3-030-53518-6_1.
- [43] Josef Urban. MPTP - motivation, implementation, first experiments. *J. Autom. Reason.*, 33(3-4):319–339, 2004. doi: 10.1007/s10817-004-6245-1.
- [44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [45] Ben Wang. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.

- 548 [46] Qingxiang Wang, Cezary Kaliszyk, and Josef Urban. First experiments with neural translation
549 of informal to formal mathematics. In Florian Rabe, William M. Farmer, Grant O. Passmore,
550 and Abdou Youssef (eds.), *Intelligent Computer Mathematics - 11th International Conference,*
551 *CICM 2018, Hagenberg, Austria, August 13-17, 2018, Proceedings*, volume 11006 of *Lecture*
552 *Notes in Computer Science*, pp. 255–270. Springer, 2018. doi: 10.1007/978-3-319-96812-4_22.
553 URL https://doi.org/10.1007/978-3-319-96812-4_22.
- 554 [47] Qingxiang Wang, Chad Brown, Cezary Kaliszyk, and Josef Urban. Exploration of neural
555 machine translation in autoformalization of mathematics in mizar. In *International Conference*
556 *on Certified Programs and Proofs*, 2020.
- 557 [48] Makarius Wenzel, Lawrence C. Paulson, and Tobias Nipkow. The Isabelle framework. In
558 Otmane Aït Mohamed, César A. Muñoz, and Sofiène Tahar (eds.), *Theorem Proving in Higher*
559 *Order Logics, 21st International Conference, TPHOLs 2008, Montreal, Canada, August 18-21,*
560 *2008. Proceedings*, volume 5170 of *Lecture Notes in Computer Science*, pp. 33–38. Springer,
561 2008. doi: 10.1007/978-3-540-71067-7_7.
- 562 [49] Minchao Wu, Michael Norrish, Christian Walder, and Amir Dezfouli. TacticZero: Learning to
563 prove theorems from scratch with deep reinforcement learning. *CoRR*, abs/2102.09756, 2021.
564 URL <https://arxiv.org/abs/2102.09756>.
- 565 [50] Yuhuai Wu, Albert Jiang, Jimmy Ba, and Roger Baker Grosse. INT: an inequality benchmark
566 for evaluating generalization in theorem proving. In *9th International Conference on Learning*
567 *Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
568 URL <https://openreview.net/forum?id=06LPudownQm>.
- 569 [51] Yuhuai Wu, Markus N Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing trans-
570 formers. *arXiv preprint arXiv:2203.08913*, 2022.
- 571 [52] Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. MiniF2F: a cross-system benchmark
572 for formal olympiad-level mathematics. *arXiv preprint arXiv:2109.00110*, 2021.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section ??.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
- (b) Did you describe the limitations of your work? **[Yes]** See discussion section.
- (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See discussion section.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- (b) Did you include complete proofs of all theoretical results? **[N/A]**

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]**
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]**
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]** It is expensive to run multiple times, and we believe the results are significant.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** See Section 5.2.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cited the creators of the assets when we first mentioned them.
- (b) Did you mention the license of the assets? **[Yes]** We included the links to the licenses of the assets when we first mentioned them in Section 5.
- (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]** We include the code we used to train the models in the supplemental material. The data we used are all under open-source licenses.
- (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[Yes]** The assets we used mostly have open-sourced licenses as mentioned previously.
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[No]** The data we used are mathematical proofs so we think it is apparent that they do not contain personally identifiable information or any offensive content.

5. If you used crowdsourcing or conducted research with human subjects...

- 622 (a) Did you include the full text of instructions given to participants and screenshots, if
623 applicable? [N/A]
- 624 (b) Did you describe any potential participant risks, with links to Institutional Review
625 Board (IRB) approvals, if applicable? [N/A]
- 626 (c) Did you include the estimated hourly wage paid to participants and the total amount
627 spent on participant compensation? [N/A]

Appendix

A Prompts

A.1 Prompt Used in the Case Studies in Section 4.2:

Natural language version: "Let $z = \frac{1+i}{\sqrt{2}}$, find $(\sum_{i=1}^1 2(z^{i^2})) \cdot (\sum_{i=1}^1 2(\frac{1}{z^{i^2}}))$. The final answer is 36."

Translate the natural language version to an Isabelle version:

```
theorem:
  fixes z::complex
  assumes h0: "z = (Complex (1/sqrt 2) (1/sqrt 2))"
  shows "(\\<Sum>k::nat=1..12. (z^(k^2)))
    * (\\<Sum> k::nat=1..12. 1/(z^(k^2))) =36"
```

Natural language version: "Determine the value of ab if $\log_8 a + \log_4 b^2 = 5$ and $\log_8 b + \log_4 a^2 = 7$. The final answer is 512". Translate the natural language version to an Isabelle version:

```
theorem:
  fixes a b ::real
  assumes "(ln a) / (ln 8) + (ln (b^2)) / (ln 4) = 5"
    "(ln b) / (ln 8) + (ln (a^2)) / (ln 4) = 7"
  shows "a * b = 512"
```

632 B Few-shot Prompts

633 B.1 Prompt used to formalize algebra problems

Natural language version: *"Simplify $(\frac{4}{x})^{-1} \left(\frac{3x^3}{x}\right)^2 \left(\frac{1}{2x}\right)^{-3}$. The final answer is $18x^8$ ".*
 Translate the natural language version to an Isabelle version:

```
theorem
  fixes x :: real
  assumes h0 : "x \<noteq> 0"
  shows "1/(4/x) * ((3*x^3)/x)^2 * (1/(1 / (2 * x)))^3 = 18 * x^8"
```

Natural language version: *"For integers n , let*

$$f(n) = \begin{cases} n^2 & \text{if } n \text{ is odd,} \\ n^2 - 4n - 1 & \text{if } n \text{ is even.} \end{cases}$$

Find $f(f(f(f(f(4)))))$. The final answer is 1". Translate the natural language version to an Isabelle version:

```
theorem
  fixes f :: "int \<Rightarrow> int"
  assumes "\<forall>n. odd n \<longrightarrow> f n = n^2"
    and "\<forall> n. even n \<longrightarrow> f n = n^2 - 4*n -1"
  shows "f 4 = -1"
```

Natural language version: *"The volume of a cone is given by the formula $V = \frac{1}{3}Bh$, where B is the area of the base and h is the height. The area of the base of a cone is 30 square units, and its height is 6.5 units. What is the number of cubic units in its volume? The final answer is 65".* Translate the natural language version to an Isabelle version:

```
theorem
  fixes b h v :: real
  assumes "0 < b \<and> 0 < h \<and> 0 < v"
    and "v = 1 / 3 * (b * h)"
    and "b = 30"
    and "h = 13 / 2"
  shows "v = 65"
```

Natural language version: *"If $3a + b + c = -3$, $a + 3b + c = 9$, $a + b + 3c = 19$, then find abc . The final answer is -56".*

Translate the natural language version to an Isabelle version:

```
theorem
  fixes a b c :: real
  assumes "3 * a + b + c = -3"
    and "a + 3 * b + c = 9"
    and "a + b + 3 * c = 19"
  shows "a * b * c = -56"
```

Natural language version: "If $f(x) = 5x - 12$, find a value for x so that $f^{-1}(x) = f(x + 1)$. The final answer is $\frac{47}{24}$ ". Translate the natural language version to an Isabelle version:

```
theorem
  fixes x :: real and \<sigma>::"real \<Rightarrow> real"
  assumes "bij \<sigma>"
    and \<sigma>:"\<forall> x. \<sigma> x = 5 * x - 12"
    and "\<sigma> (x + 1) = (inv \<sigma>) x"
  shows "x = 47 / 24"
```

Natural language version: "What is the x -coordinate for the x -intercept of the line containing the points $(7, 4)$ and $(6, 3)$? The final answer is 3". Translate the natural language version to an Isabelle version:

```
theorem
  fixes a b :: real
  and f :: "real \<Rightarrow> real"
  assumes h0 : "\<And>x. f x = a * x + b"
    and h1 : "f 7 = 4"
    and h2 : "f 6 = 3"
  shows "f 3 = 0"
```

Natural language version: "Given $2^a = 32$ and $a^b = 125$ find b^a . The final answer is 243". Translate the natural language version to an Isabelle version:

```
theorem
  fixes a b :: real
  assumes "2 powr a = 32"
    and "a powr b = 125"
  shows "b powr a = 243"
```

Natural language version: "Let

$$f(x) = \begin{cases} x^2 + 9 & \text{if } x < -5, \\ 3x - 8 & \text{if } x \geq -5. \end{cases}$$

If $f(x) = 10$, find the sum of all possible values of x . The final answer is 6". Translate the natural language version to an Isabelle version:

```
theorem
  fixes f :: "real \<Rightarrow> real"
  assumes "\<forall> x < -5. f x = x^2 + 9"
    and "\<forall> x \<ge> -5. f x = 3 * x - 8"
  shows "(\<Sum> k \<in> (f -> {10}). k) = 6"
```

Natural language version: "Simplify $(9 - 4i) - (-3 - 4i)$. The final answer is 12". Translate the natural language version to an Isabelle version:

```
theorem
  fixes q e :: complex
  assumes h0 : "q = Complex (Re 9) (Im (-4))"
    and h1 : "e = Complex (Re (-3)) (Im (-4))"
  shows "q - e = 12"
```

Natural language version: "What is the minimum possible value for y in the equation $y = x^2 - 6x + 13$? The final answer is 4".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes x y :: real
  assumes h0 : "y = x^2 - 6 * x + 13"
  shows "4 \<le> y"
```

Natural language version: "If n is a positive integer such that $2n$ has 28 positive divisors and $3n$ has 30 positive divisors, then how many positive divisors does $6n$ have? The final answer is 35". Translate the natural language version to an Isabelle version:

```
theorem
  fixes n :: nat
  assumes "n>0"
    and "card ({k. k dvd (2*n)}) = 28"
    and "card ({k. k dvd (3*n)}) = 30"
  shows "card ({k. k dvd (6*n)}) = 35"
```

Natural language version: "Let n be the number of integers m in the range $1 \leq m \leq 8$ such that $\gcd(m, 8) = 1$. What is the remainder when 3^n is divided by 8? The final answer is 1". Translate the natural language version to an Isabelle version:

```
theorem
  fixes n :: nat
  assumes "n = card {k::nat. gcd k 8 = 1 \<and> 1<le>k \<and> k < 8}"
  shows "(3^n) mod 8 = (1::nat)"
```

Natural language version: "What is the remainder when $1 + 2 + 3 + 4 + \dots + 9 + 10$ is divided by 9? The final answer is 1".

Translate the natural language version to an Isabelle version:

```
theorem
  "(\<Sum> k< 11. k) mod 9 = (1::nat)"
```

Natural language version: "Cards are numbered from 1 to 100. One card is removed and the values on the other 99 are added. The resulting sum is a multiple of 77. What number was on the card that was removed? The final answer is 45".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes x :: nat
  assumes h0 : "1 \<le> x \<and> x \<le> 100"
    and h1 : "77 dvd ((\<Sum>k::nat=0..100. k)-x)"
  shows "x=45"
```

Natural language version: "Find $9^{-1} \pmod{100}$, as a residue modulo 100. (Give an answer between 0 and 99, inclusive.) The final answer is 89 $\pmod{100}$ ".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes x::nat
  assumes "x < 100"
    and "x*9 mod 100 = 1"
  shows "x = 89"
```


Natural language version: "Suppose m is a two-digit positive integer such that $6^{-1} \pmod{m}$ exists and $6^{-1} \equiv 6^2 \pmod{m}$. What is m ? The final answer is 43".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes m x :: nat
  assumes h0 : "10 \<le> m"
    and h1 : "m \<le> 99"
    and h2 : "(6 * x) mod m = 1"
    and h3 : "(x - 6^2) mod m = 0"
  shows "m = 43"
```

Natural language version: "Find $24^{-1} \pmod{11^2}$. That is, find the residue b for which $24b \equiv 1 \pmod{11^2}$. Express your answer as an integer from 0 to $11^2 - 1$, inclusive. The final answer is 116".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes b::int
  assumes "\<forall>b::int. 0\<le>b \<and> b\<le>11^2 \<and> [b * 24 = 1]
    (mod (11^2))"
  shows "b = 116"
```

Natural language version: "Given that $p \geq 7$ is a prime number, evaluate

$$2^{-1} \cdot 2^{-1} + 2^{-1} \cdot 3^{-1} + 3^{-1} \cdot 4^{-1} + \cdots + (p-2)^{-1} \cdot (p-1)^{-1} \pmod{p}.$$

The final answer is $2 \pmod{p}$ ".

Translate the natural language version to an Isabelle version:

```
theorem
  fixes p :: nat
  assumes "prime p"
    and "7 \<le> p"
  shows "(\<Sum> k \<in> {1..<p-1>}. (inv_mod k p) * (inv_mod (k+1) p)) = 2"
```

Natural language version: "What is the remainder when $2000 + 2001 + 2002 + 2003 + 2004 + 2005 + 2006$ is divided by 7? The final answer is 0". Translate the natural language version to an Isabelle version:

```
theorem
  "(2000 + 2001 + 2002 + 2003 + 2004 + 2005 + 2006) mod 7 = (0::nat)"
```

Natural language version: "One morning each member of Angela's family drank an 8-ounce mixture of coffee with milk. The amounts of coffee and milk varied from cup to cup, but were never zero. Angela drank a quarter of the total amount of milk and a sixth of the total amount of coffee. How many people are in the family? The final answer is 5". Translate the natural language version to an Isabelle version:

```
theorem
  fixes x y n :: nat
  assumes "x / 4 + y / 6 = (x + y) / n"
    and "n\<noteq>0"
    "x\<noteq>0" "y\<noteq>0"
  shows "n = 5"
```

639 B.3 Prompt used for informalization

Isabelle version:

```
theorem
  fixes z::complex
  assumes h0: "z = (Complex (1/sqrt 2) (1/sqrt 2))"
  shows "(\\<Sum>k::nat=1..12. (z^(k^2)))
    * (\\<Sum> k::nat=1..12. 1/(z^(k^2))) =36"
```

Translate the Isabelle version to a natural language version: "Let $z = \frac{1+i}{\sqrt{2}}$, find $(\sum_{i=1}^1 2(z^{i^2})) \cdot (\sum_{i=1}^1 2(\frac{1}{z^{i^2}}))$. The final answer is 36."

Isabelle version:

```
theorem
  fixes a b ::real
  assumes "(ln a) / (ln 8) + (ln (b^2)) / (ln 4) = 5"
    "(ln b) / (ln 8) + (ln (a^2)) / (ln 4) = 7"
  shows "a * b = 512"
```

Translate the Isabelle version to a natural language version: "Determine the value of ab if $\log_8 a + \log_4 b^2 = 5$ and $\log_8 b + \log_4 a^2 = 7$. The final answer is 512".

640

641 C More on case Studies

642 C.1 Case Study 1

643 Here we show the output given by PaLM model 540B. The model got the most structure correctly,
644 but made some syntactical errors.

Question:

Natural language version: "Prove that there is no function f from the set of non-negative integers into itself such that $f(f(n)) = n + 1987$ for every n ." Translate the natural language version to an Isabelle version:

PaLM 540B Output:

```
theorem
  fixes f :: "nat -> nat"
  assumes "forall n::nat. f(f(n)) = n + 1987"
  shows "false"
```

645

646 **C.2 Case Study 2**

647 We show below the output given by PaLM models of different sizes.

Question:

Natural Language version: *"When all the girls at Madeline's school line up in rows of eight, there are seven left over. If instead they line up in rows of four, how many are left over? The final answer is 3."*

Translate the natural language version to an Isabelle version:

PaLM 540B output:

```
theorem
  fixes n::nat
  assumes "n mod 8 = 7"
  shows "n mod 4 = 3"
```

648

Question:

Natural Language version: *"When all the girls at Madeline's school line up in rows of eight, there are seven left over. If instead they line up in rows of four, how many are left over? The final answer is 3."*

Translate the natural language version to an Isabelle version:

PaLM 64B output:

```
theorem
  fixes girls :: nat
  assumes "(girls / 8) + 7 = girls"
  shows "girls = 3"
```

649

Question:

Natural Language version: *"When all the girls at Madeline's school line up in rows of eight, there are seven left over. If instead they line up in rows of four, how many are left over? The final answer is 3."*

Translate the natural language version to an Isabelle version:

PaLM 8B output:

```
theorem
  fixes a b ::real
  assumes "(ln a) / (ln 8) + (ln (b^2)) / (ln 4) = 7"
  "(ln b) / (ln 8) + (ln (a^2)) / (ln 4) = 3"
  shows "a * b = 3"
```

650

651 C.3 Case Study 3

652 If we include the following example in the prompt in addition to the two examples shown in
653 Appendix B.3:

Natural language version: "What is the x -coordinate for the x -intercept of the line containing the points $(7, 4)$ and $(6, 3)$? The final answer is 3". Translate the natural language version to an Isabelle version:

```
theorem
  fixes a b :: real
    and f :: "real \ $\rightarrow$  real"
  assumes h0 : "\ $\wedge$ x. f x = a * x + b"
    and h1 : "f 7 = 4"
    and h2 : "f 6 = 3"
  shows "f 3 = 0"
```

654

655 Then Codex was able to formalize the problem perfectly:

Question:

Natural language version: "Let f be a linear function for which $f(6) - f(2) = 12$. What is $f(12) - f(2)$? The final answer is 30."

Translate the natural language version to an Isabelle version:

Codex Output:

```
theorem
  fixes a b :: real
    and f :: "real \ $\rightarrow$  real"
  assumes h0 : "\ $\wedge$ x. f x = a * x + b"
    and h1 : "f 6 - f 2 = 12"
  shows "f 12 - f 2 = 30"
```

656

657 Note that this additional example in the prompt is only tangentially related. The phrase "linear
658 function" is not even mentioned. Instead, it shows an example how to formalize a *line*. But that is
659 sufficient for the model to understand how to formalize a "linear function".

660 D More Examples on Mathematical Competition Problems

661 D.1 Perfect Translations

Question:

Natural language version: "At what point does the line containing the points $(1, 7)$ and $(3, 11)$ intersect the y -axis? Express your answer as an ordered pair. The final answer is $(0, 5)$ ".

Translate the natural language version to an Isabelle version:

Codex Output:

```
theorem
  fixes a b :: real
    and f :: "real \ $\rightarrow$  real"
  assumes h0 : "\ $\wedge$ x. f x = a * x + b"
    and h1 : "f 1 = 7"
    and h2 : "f 3 = 11"
  shows "f 0 = 5"
```

662

Question:

Natural language version: "A map of the town that Annie, Barbara, and Charlie live in can be represented by the Cartesian plane. Annie is located at $(6, -20)$ and Barbara is located at $(1, 14)$. They agree to meet at the closest point that is equidistant from their current locations and walk upwards together to get to Charlie's location at $(\frac{7}{2}, 2)$. How many units upward do Annie and Barbara walk together to get to Charlie? The final answer is 5." Translate the natural language version to an Isabelle version:

Codex Output:

```
theorem
  fixes a b c :: "real \ real"
  assumes h0 : "a = (6, -20)"
    and h1 : "b = (1, 14)"
    and h2 : "c = (7 / 2, 2)"
  shows "dist a c + dist b c = 5"
```

663

Question:

Natural language version: "Let x be a positive integer, and define the integers $n = x^2 + 2x + 17$ and $d = 2x + 5$. When dividing n by d , the quotient is x , and the remainder is 7. Find x . The final answer is 2".

Codex Output:

```
theorem
  fixes x n d :: int
  assumes h0 : "x > 0"
    and h1 : "n = x^2 + 2 * x + 17"
    and h2 : "d = 2 * x + 5"
    and h3 : "n div d = x"
    and h4 : "n mod d = 7"
  shows "x = 2"
```

664

Question:

Natural language version: "A line segment begins at $(1, 3)$. It is 5 units long and ends at the point $(4, x)$ and $x > 0$. What is the value of x ? The final answer is 7".

Codex Output:

```
theorem
  fixes x :: real
  assumes h0 : "x > 0"
    and h1 : "5 = sqrt ((4 - 1)^2 + (x - 3)^2)"
  shows "x = 7"
```

665

Natural language version: "A line segment begins at $(1, 3)$. It is 5 units long and ends at the point $(4, x)$ and $x > 0$. What is the value of x ? The final answer is 7".

Codex Output:

```
theorem
  fixes x :: real
  assumes h0 : "x > 0"
    and h1 : "5 = sqrt ((4 - 1)^2 + (x - 3)^2)"
  shows "x = 7"
```

666

Question:

Natural language version: "Let $f(x) = x^2 - 2x$. How many distinct real numbers c satisfy $f(f(f(f(c)))) = 3$? The final answer is 9".

Codex Output:

```
theorem
  fixes f :: "real \ $\rightarrow$  real"
  assumes h0 : "\forall x. f x = x^2 - 2 * x"
  shows "card {c. f (f (f (f c))) = 3} = 9"
```

667

Question:

Natural language version: "Find the solutions to $z^4 = -4$. Enter the solutions, separated by commas. The final answer is $1 + i, 1 - i, -1 + i, -1 - i$ ".

Codex Output:

```
theorem
  fixes z :: complex
  assumes "z^4 = -4"
  shows "z = 1 + \ $i$  \ $\vee$  z = 1 - \ $i$  \ $\vee$  z = -1 + \ $i$  \ $\vee$  z = -1 - \ $i$ "
```

668

Question:

Natural language version: "Let z be a complex number satisfying $z^2 = 4z - 19 + 8i$. Given that $|z|$ is an integer, find z . The final answer is $3 + 4i$ ".

Codex Output:

```
theorem
  fixes z :: complex
  assumes h0 : "z^2 = 4 * z - 19 + 8 * \ $i$ "
    and h1 : "norm z \ $\in$  \ $\mathbb{Z}$ "
  shows "z = 3 + 4 * \ $i$ "
```

669

Question:

Natural language version: "The sequence (a_n) satisfies

$$a_1 + a_2 + a_3 + \cdots + a_n = n^2 a_n$$

for all $n \geq 2$. If $a_{63} = 1$, find a_1 . The final answer is 2016".

Codex Output:

```
theorem
  fixes a :: "nat \ $\rightarrow$  real"
  assumes h0 : "\forall n \ge 2. (\sum k = 1..n. a k) = n^2 * a n"
    and h1 : "a 63 = 1"
  shows "a 1 = 2016"
```

670

Question:

Natural language version: "Find the units digit of the following within the indicated number base: $52_7 + 62_7$. The final answer is 4".

Codex Output:

```
theorem
  fixes x y :: nat
  assumes "x < 7"
    and "y < 7"
  shows "((x*7+2) + (y*7+2)) mod 7 = 4"
```

671

Question:

Natural language version: "How many factors of 8000 are perfect squares? The final answer is 8".

Codex Output:

```
theorem
  fixes n :: nat
  assumes "n = card {k::nat. k^2 dvd 8000}"
  shows "n = 8"
```

672

Question:

Natural language version: "Solve the congruence $11n \equiv 7 \pmod{43}$, as a residue modulo 43. (Give an answer between 0 and 42.) The final answer is 28".

Codex Output:

```
theorem
  fixes n :: nat
  assumes "n < 43"
    and "(11 * n) mod 43 = 7"
  shows "n = 28"
```

673

Question:

Natural language version: "What is the 100th digit after the decimal point of the decimal representation of $1/7$? The final answer is 8".

Codex Output:

```
theorem
  fixes x :: nat
  assumes "x = (1/7) * 10^100"
  shows "x mod 10 = 8"
```

674

Question:

Natural language version: "Two positive integers m and n are chosen such that m is the smallest positive integer with only two positive divisors and n is the largest integer less than 100 with exactly three positive divisors. What is $m + n$? The final answer is 51".

Codex Output:

```
theorem
  fixes m n :: nat
  assumes "m = (LEAST k. card {d. d dvd k \<and> d > 0} = 2)"
    and "n = (GREATEST k. card {d. d dvd k \<and> d > 0} = 3 \<and> k < 100)"
  shows "m + n = 51"
```

675

676 D.2 Incorrect Translations

677 In this section, we present failure cases for each of the categories in Table 2 in Section. 4.4.

Error category: incomplete/ill-formed/unclear prompt

Question:

Natural language version: "Find the largest value of t such that

$$\frac{13t^2 - 34t + 12}{3t - 2} + 5t = 6t - 1.$$

The final answer is $\frac{5}{2}$.

Codex Output:

```
theorem
  fixes t :: real
  assumes h0 : "13 * t^2 - 34 * t + 12 = 3 * t - 2"
    and h1 : "5 * t = 6 * t - 1"
  shows "t = 5 / 2"
```

678

Error category: fail to align definitions or concepts

Question:

Natural language version: "What is the smallest four-digit positive integer that is divisible by 47? The final answer is 1034".

Codex Output:

```
theorem
  fixes x :: nat
  assumes "x > 999"
    and "x < 10000"
    and "47 dvd x"
  shows "x = 1034"
```

679

680 Here, the translator correctly captured the meaning of "four-digit positive integer" (i.e., " $x > 999$ "
681 and " $x < 10000$ ") and that of "divisible by 47" (i.e., " $47 \text{ dvd } x$ "). However, it fails to align the
682 "smallest" concept to the one in Isabelle.

Error category: inconsistent/missing assumption

Question:

Natural language version: "Jack walked up a hill at a speed of $(x^2 - 11x - 22)$ miles per hour. Meanwhile, Jill walked a total distance of $(x^2 - 3x - 54)$ miles in $(x + 6)$ hours. If Jack and Jill walked at the same speed, what is that speed, in miles per hour? The final answer is 4".

Codex Output:

```
theorem
  fixes x :: real
  assumes h0 : "x^2 - 11 * x - 22 = x^2 - 3 * x - 54"
  shows "x = 4"
```

683

684 Here, the translator missed the " $(x + 6)$ " part in the natural language—it should be translated to " x^2
685 $- 11 * x - 22 = (x^2 - 3 * x - 54) / (x + 6)$ " instead. Also, the conclusion is incorrect either: it
686 should be " $x^2 - 11 * x - 22 = 4$ " or " $(x^2 - 3 * x - 54) / (x + 6) = 4$ ".

Error category: syntactical/type error)

Question:

Natural language version: "Find

$$\sum_{N=1}^{1024} \lfloor \log_2 N \rfloor.$$

The final answer is 8204".

Codex Output:

```
theorem
  fixes N :: nat
  shows "(\<Sum> N = 1..1024. nat \<lfloor>log 2 N\<rfloor>) = 8204"
```

687

688 In the current version, N will be inferred as a real number, so that the interval "1..1024" will be a
689 real interval from 1 to 1024. This is incorrect, and sum over such infinite interval will result in 0 (i.e.,
690 "(\<Sum> N = 1..1024. nat \<lfloor>log 2 N\<rfloor>)= 0"). The correct version requires
691 an explicit type conversion: "(\<Sum> N = 1..1024. \<lfloor>log 2 (of_int N)\<rfloor>)=
692 8204".

Error category: missing definition in Isabelle

Question:

Natural language version: "The product of two positive three-digit palindromes is 436,995.
What is their sum? The final answer is 1332".

Codex Output:

```
theorem
  fixes x y :: nat
  assumes "100 \<le> x \<and> x \<le> 999"
    and "100 \<le> y \<and> y \<le> 999"
    and "x * y = 436995"
  shows "x + y = 1332"
```

693

694 The only missing part in this example is the definition of palindromes. It appears that the translator
695 does not know what a palindrome is so that it chooses to drop it from the assumptions.

Error category: wrong application of functions

Question:

Natural language version: "If the operation Z is defined as $aZb = b + 10a - a^2$, what is the
value of $2Z6$? The final answer is 22".

Codex Output:

```
theorem
  fixes a b :: real
  assumes h0 : "a Z b = b + 10 * a - a^2"
  shows "2 Z 6 = 22"
```

696

697 The function Z should be applied in a prefix way in Isabelle. Also, the translator missed the "for all"
698 quantifier in the assumption. The correct assumption should be "\<forall> a b. Z a b = b + 10
699 * a - a^2".

Error category: other

Question:

Natural language version: "Simplify: $|-3^2 + 4|$ The final answer is 5".

Codex Output:

```
theorem
  fixes x :: real
  assumes h0 : "x = - 3 ^2 + 4"
  shows "abs x = 5"
```

700

701 Here, a parenthesis should be placed around 3^2 (i.e., " $x = - (3^2) + 4$ "). This is a subtle issue to
702 operator precedence in Isabelle.

Error category: other

Question:

Natural language version: "How many numbers from 1 to 150 are not perfect squares or perfect cubes? The final answer is 135".

Codex Output:

```
theorem
  fixes n :: nat
  assumes "n <= 150"
  shows "card {k::nat. k <= n & \<not> (\<exists> m. m^2 = k \<or>
    m^3 = k) } = 135"
```

703

704 The translator somehow introduces an unnecessary variable n and assumed " $n \leq 150$ ". The
705 correct translation should have n replaced with 150 in the conclusion: " $\text{card } \{k::\text{nat. } k \leq$
706 $150 \text{ \& \<not> (\<exists> m. m^2 = k \text{ \<or> m^3 = k)}\} = 135$ ".

707 E More Examples on Advanced Mathematics

708 E.1 Formalization

Question:

Natural language version: "Let f be a continuous function from a compact convex set S to itself. Then there exists $x \in S$ such that $f(x) = x$." Translate the natural language version to an Isabelle version:

Codex Output:

```
theorem
  fixes f :: "'a::real_normed_vector \ $\rightarrow$  'a"
  assumes "continuous_on S f"
        "compact S"
        "convex S"
  shows "\exists x \in S. f x = x"
```

709

710 Here, we asked the model to formalize the statement of the famous Brouwer's fixed-point theorem.
711 The model made two mistakes. The first one is that the model erroneously assume the function is on a
712 real normed vector space rather than a euclidean space or a banach space. This is quite understandable
713 as we left this information in the source. The second and more critical error is that the absence of the
714 assumption of f mapping from S to itself (i.e., $f(S) \subseteq S$ or $f(S) = S$).

Question: Let γ be a closed path in the complex plane, and z be a point not on γ . Then the winding number of γ around z is an integer.

Codex Output:

```
theorem
  fixes \gamma :: "real \ $\rightarrow$  complex" and z :: complex
  assumes "continuous_on {0..1} \gamma"
        "\gamma 0 = \gamma 1"
        "z \notin \gamma ` {0..1}"
  shows "winding_number \gamma z \in \int"
```

715

716 Here, the model has made a perfect translation—not only the model correctly infers that a path in the
717 complex plane should be a function from the real interval $[0, 1]$ to complex, it also manages to align a
718 path being closed to the assumption $\gamma 0 = \gamma 1$.

720 The following lemma statements are from the standard library of Isabelle (i.e., HOL-
721 Complex_Analysis and HOL-Probability).

```

proposition openin_clauses:
  fixes U :: "'a topology"
  shows
    "openin U {}"
    " $\bigwedge S T. \text{openin } U S \implies \text{openin } U T \implies \text{openin } U (S \cap T)$ "
    " $\bigwedge K. (\forall S \in K. \text{openin } U S) \implies \text{openin } U (\bigcup K)$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

Let U be a topology, then U is open if and only if it satisfies the following conditions:

1. U is open.
2. If S and T are open, then $S \cap T$ is open.
3. If K is a collection of open sets, then $\bigcup K$ is open.

722

```

lemma openin_Int[intro]:
  "openin U S  $\implies$  openin U T
   $\implies$  openin U (S  $\cap$  T)"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let U be a topology, then \emptyset is open in U , if S and T are open in U , then $S \cap T$ is open in U , if K is a collection of sets, and each set in K is open in U , then $\bigcup K$ is open in U ".

723

```

lemma openin_Int[intro]:
  "openin U S  $\implies$  openin U T
   $\implies$  openin U (S  $\cap$  T)"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If S and T are open in U , then $S \cap T$ is open in U ".

724

```

lemma openin_Inter [intro]:
  assumes "finite  $\mathcal{F}$ " "  $\mathcal{F} \neq \{\}$ " " $\bigwedge X. X \in \mathcal{F} \implies \text{openin } T X$ "
  shows "openin T ( $\bigcap \mathcal{F}$ )"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If T is a topological space, and \mathcal{F} is a finite non-empty collection of open subsets of T , then $\bigcap \mathcal{F}$ is open in T ".

725

lemma *openin_subtopology*:
 "openin (subtopology U V) S \longleftrightarrow ($\exists T$. openin U T \wedge S = T \cap V)"

Translate the Isabelle version to a natural language version:

Codex Output:

"S is open in the subtopology of U induced by V if and only if there exists an open set T in U such that S = T \cap V".

726

lemma *closedin_subtopology_Un*:
 "[closedin (subtopology X T) S; closedin (subtopology X U) S]
 \implies closedin (subtopology X (T \cup U)) S"

Translate the Isabelle version to a natural language version:

Codex Output:

"If S is closed in T and S is closed in U, then S is closed in T \cup U".

727

lemma *connected_closedin*:
 "connected S \longleftrightarrow
 (E1 E2.
 closedin (top_of_set S) E1 \wedge
 closedin (top_of_set S) E2 \wedge
 S \subseteq E1 \cup E2 \wedge E1 \cap E2 = {} \wedge E1 \neq {} \wedge E2 \neq {})"

Translate the Isabelle version to a natural language version:

Codex Output:

"A set S is connected if and only if there does not exist two closed sets E₁ and E₂ such that S \subseteq E₁ \cup E₂, E₁ \cap E₂ = \emptyset , E₁ $\neq \emptyset$ and E₂ $\neq \emptyset$ ".

728

definition *continuous_map* **where**
 "continuous_map X Y f \equiv
 ($\forall x \in \text{topspace } X$. f x \in topspace Y) \wedge
 ($\forall U$. openin Y U \longrightarrow openin X {x \in topspace X. f x \in U})"

Translate the Isabelle version to a natural language version:

Codex Output:

"A function f from a topological space X to a topological space Y is continuous if for every open set U in Y, the preimage f⁻¹(U) is open in X".

729

definition *continuous_map* **where**
 "continuous_map X Y f \equiv
 ($\forall x \in \text{topspace } X$. f x \in topspace Y) \wedge
 ($\forall U$. openin Y U \longrightarrow openin X {x \in topspace X. f x \in U})"

Translate the Isabelle version to a natural language version:

Codex Output:

"A continuous map f : X \rightarrow Y is a map such that for every open set U \subset Y, the preimage f⁻¹(U) is open in X".

730

lemma *connectedin_separation*:
 "connectedin X $S \longleftrightarrow$
 $S \subseteq \text{topspace } X \wedge$
 $(C_1 \ C_2. \ C_1 \cup C_2 = S \wedge C_1 \neq \{\} \wedge C_2 \neq \{\}$
 $\wedge C_1 \cap X \text{ closure_of } C_2 = \{\} \wedge C_2$
 $\cap X \text{ closure_of } C_1 = \{\})$ " (is "?lhs = ?rhs")

Translate the Isabelle version to a natural language version:

Codex Output:

" S is connected in X if and only if $S \subseteq X$ and there do not exist two non-empty sets C_1 and C_2 such that $C_1 \cup C_2 = S$ and $C_1 \cap X\overline{C_2} = \emptyset$ and $C_2 \cap X\overline{C_1} = \emptyset$ ".

731

lemma *connectedin_Union*:
 assumes \mathcal{U} : " $\bigwedge S. S \in \mathcal{U} \implies \text{connectedin } X \ S$ " and ne: " $\bigcap \mathcal{U} \neq \{\}$ "
 shows " $\text{connectedin } X \ (\bigcup \mathcal{U})$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If S is connected in X for all $S \in \mathcal{U}$ and $\bigcap_{S \in \mathcal{U}} S \neq \emptyset$, then $\bigcup_{S \in \mathcal{U}} S$ is connected in X ".

732

lemma *compact_space_alt*:
 " $\text{compact_space } X \longleftrightarrow$
 $(\forall \mathcal{U}. (\forall U \in \mathcal{U}. \text{openin } X \ U) \wedge \text{topspace } X \subseteq \bigcup \mathcal{U}$
 $\longrightarrow (\exists \mathcal{F}. \text{finite } \mathcal{F} \wedge \mathcal{F} \subseteq \mathcal{U} \wedge \text{topspace } X \subseteq \bigcup \mathcal{F}))$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"A topological space X is compact if and only if for every open cover \mathcal{U} of X , there exists a finite subcover \mathcal{F} of \mathcal{U} ".

733

lemma *compactin_subtopology*:
 " $\text{compactin } (\text{subtopology } X \ S) \ T \longleftrightarrow \text{compactin } X \ T \wedge T \subseteq S$ "

Translate the Isabelle version to a natural language version:

Codex Output:

" T is compact in S if and only if T is compact in X and $T \subseteq S$ ".

734

lemma *closedin_compact_space*:
 " $\llbracket \text{compact_space } X; \text{closedin } X \ S \rrbracket \implies \text{compactin } X \ S$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If X is a compact space and S is a closed subset of X , then S is a compact space".

735

lemma *compactin_discrete_topology*:
 "compactin (discrete_topology X) S \longleftrightarrow S \subseteq X \wedge finite S"

Translate the Isabelle version to a natural language version:

Codex Output:

"S is compact in the discrete topology on X if and only if S is a finite subset of X".

736

lemma *image_compactin*:
 assumes *cpt*: "compactin X S" and *cont*: "continuous_map X Y f"
 shows "compactin Y (f ` S)"

Translate the Isabelle version to a natural language version:

Codex Output:

"If S is compact in X and f is continuous from X to Y, then f(S) is compact in Y".

737

lemma *continuous_on_open*:
 "continuous_on S f \longleftrightarrow
 ($\forall T$. openin (top_of_set (f ` S)) T \longrightarrow
 openin (top_of_set S) (S \cap f $^{-1}$ T))"

Translate the Isabelle version to a natural language version:

Codex Output:

"f is continuous on S if and only if for any open set T in f(S), f⁻¹(T) is open in S".

738

lemma *continuous_closedin_preimage_gen*:
 assumes "continuous_on S f" "f ` S \subseteq T"
 "closedin (top_of_set T) U"
 shows "closedin (top_of_set S) (S \cap f $^{-1}$ U)"

Translate the Isabelle version to a natural language version:

Codex Output:

"If f is continuous on S and f(S) \subseteq T, and U is closed in T, then S \cap f⁻¹(U) is closed in S".

739

theorem *Casorati_Weierstrass*:
 assumes "open M" "z \in M" "f holomorphic_on (M - {z})"
 and " $\bigwedge l$. \neg (f l) (at z)"
 " $\bigwedge l$. \neg ((inverse \circ f) l) (at z)"
 shows "closure(f ` (M - {z})) = UNIV"

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a holomorphic function on an open set M containing z. If f is not constant in any neighborhood of z and f⁻¹ is not constant in any neighborhood of z, then f(M - {z}) is dense in \mathbb{C} ".

740

theorem

```

assumes "open M" "z ∈ M" "f holomorphic_on (M - {z})"
and "⋀ l. ¬ (f l) (at z)"
      "⋀ l. ¬ ((inverse ∘ f) l) (at z)"
shows "closure(f ` (M - {z})) = UNIV"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let M be an open set and $z \in M$. Let f be a holomorphic function on $M - \{z\}$. If f is not bounded in any neighborhood of z and $\frac{1}{f}$ is not bounded in any neighborhood of z , then $f(M - \{z\})$ is dense in \mathbb{C} .

741

theorem great_Picard:

```

assumes "open M" "z ∈ M" "a ≠ b"
      and holf: "f holomorphic_on (M - {z})"
and fab: "⋀ w. w ∈ M - {z} ⇒ f w ≠ a ∧ f w ≠ b"
obtains l where "(f l) (at z)"
      ∨ ((inverse ∘ f) l) (at z)"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a holomorphic function on $M - \{z\}$, where M is an open set and $z \in M$. If $f(z) \neq a$ and $f(z) \neq b$, then f or $\frac{1}{f}$ has a limit at z ".

742

theorem brouwer:

```

fixes f :: "'a::euclidean_space ⇒ 'a"
assumes S: "compact S" "convex S" "S ≠ {}"
      and conf: "continuous_on S f"
      and fim: "f ` S ⊆ S"
obtains x where "x ∈ S" and "f x = x"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a continuous function from a compact convex set S to itself. Then there exists $x \in S$ such that $f(x) = x$ ".

743

theorem Jordan_curve:

```

fixes c :: "real ⇒ complex"
assumes "simple_path c" and loop: "pathfinish c = pathstart c"
obtains inner outer where
  "inner ≠ {}" "open inner" "connected inner"
  "outer ≠ {}" "open outer" "connected outer"
  "bounded inner" "¬ bounded outer" "inner ∩ outer = {}"
  "inner ∪ outer = - path_image c"
  "frontier inner = path_image c"
  "frontier outer = path_image c"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let c be a simple closed curve. Then there exist two disjoint open sets U and V such that U is bounded, V is unbounded, $U \cup V = \mathbb{C} \setminus c$, $U \cap V = \emptyset$, $\partial U = c$, and $\partial V = c$ ".

744

corollary *Jordan_disconnected*:
 fixes $c :: \text{"real"} \Rightarrow \text{"complex"}$
 assumes $\text{"simple_path } c"$ $\text{"pathfinish } c = \text{pathstart } c"$
 shows $\neg \text{connected}(\text{path_image } c)$

Translate the Isabelle version to a natural language version:

Codex Output:

"If c is a simple closed curve, then $R^2 - c$ is not connected".

745

theorem *inverse_squares_sums*: $(\lambda n. 1 / (n + 1)^2) \text{ sums } (\pi^2 / 6)$

Translate the Isabelle version to a natural language version:

Codex Output:

"Find the sum of the series $\sum_{i=1}^{\infty} \frac{1}{(i+1)^2}$. The final answer is $\frac{\pi^2}{6}$ ".

746

proposition *connected_space_product_topology*:
 $\text{"connected_space}(\text{product_topology } X \ I) \longleftrightarrow$
 $(\prod_{i \in I} \text{topspace } (X \ i)) = \{\} \vee (\forall i \in I. \text{connected_space}(X \ i))"$

Translate the Isabelle version to a natural language version:

Codex Output:

"The product topology of a family of topological spaces is connected if and only if the product of the topological spaces is empty or each topological space is connected".

747

theorem (in *ring_of_sets*) *caratheodory'*:
 assumes $\text{posf: "positive } M \ f"$ and $\text{ca: "countably_additive } M \ f"$
 shows $\exists \mu :: 'a \text{ set} \Rightarrow \text{ennreal.}$
 $(\forall s \in M. \mu \ s = f \ s) \wedge \text{measure_space } \Omega \ (\text{sigma_sets } \Omega \ M) \ \mu"$

Translate the Isabelle version to a natural language version:

Codex Output:

"Let M be a σ -algebra on Ω , and f be a positive countably additive function on M . Then there exists a measure μ on M such that $\mu(s) = f(s)$ for all $s \in M$ and (Ω, M, μ) is a measure space".

748

lemma *sigma_finite_density_unique*:
 assumes $\text{borel: "f} \in \text{borel_measurable } M"$ $\text{"g} \in \text{borel_measurable } M"$
 and $\text{fin: "sigma_finite_measure (density } M \ f)"}$
 shows $\text{"density } M \ f = \text{density } M \ g \longleftrightarrow (\text{AE } x \text{ in } M. f \ x = g \ x)"}$

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f and g be two measurable functions on a sigma-finite measure space M . Then f and g are equal almost everywhere if and only if the measures μ_f and μ_g are equal".

749

```

lemma mvt_simple:
  fixes f :: "real  $\Rightarrow$  real"
  assumes "a < b"
  and derf: " $\bigwedge x. [a \leq x; x \leq b] \implies (f \text{ has\_derivative } f' \ x) \text{ (at } x \text{ within } \{a..b\})$ "
  shows " $\exists x \in \{a <..<b\}. f \ b - f \ a = f' \ x \ (b - a)$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a function defined on $[a, b]$ and differentiable on (a, b) . Then there exists a point $c \in (a, b)$ such that $f(b) - f(a) = f'(c)(b - a)$ ".

750

```

theorem second_countable_borel_measurable:
  fixes X :: "'a::second_countable_topology set set"
  assumes eq: "open = generate_topology X"
  shows "borel = sigma UNIV X"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let X be a set of open sets in a second countable topology. Then the Borel σ -algebra is the σ -algebra generated by X ".

751

```

theorem borel_measurable_ereal_iff_real:
  fixes f :: "'a  $\Rightarrow$  ereal"
  shows "f  $\in$  borel_measurable M  $\longleftrightarrow$ 
    (( $\lambda x. \text{real\_of\_ereal } (f \ x)$ )  $\in$  borel_measurable M  $\wedge$  f -' { $\infty$ }
       $\cap$  space M  $\in$  sets M  $\wedge$  f -' {- $\infty$ }  $\cap$  space M  $\in$  sets M)"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a function from M to \mathbb{R} , then f is measurable if and only if f is measurable and $f^{-1}(\infty)$ and $f^{-1}(-\infty)$ are measurable sets".

752

```

theorem condensation_test:
  assumes mono: " $\bigwedge m. 0 < m \implies f \ (Suc \ m) \leq f \ m$ "
  assumes nonneg: " $\bigwedge n. f \ n \geq 0$ "
  shows "summable f  $\longleftrightarrow$  summable ( $\lambda n. 2^n * f \ (2^n)$ )"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a non-negative function such that $f(n + 1) \leq f(n)$ for all n . Prove that f is summable if and only if $2^n f(2^n)$ is summable".

753

```

theorem ratio_test_convergence:
  fixes f :: "nat  $\Rightarrow$  real"
  assumes pos_f: "eventually ( $\lambda n. f\ n > 0$ ) sequentially"
  defines "l  $\equiv$  liminf ( $\lambda n. ereal (f\ n / f\ (Suc\ n))$ )"
  assumes l: "l > 1"
  shows "summable f"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a positive function, and let l be the limit inferior of the sequence $\frac{f(n)}{f(n+1)}$. If $l > 1$, then f is summable".

754

```

theorem Stone_Weierstrass_polynomial_function:
  fixes f :: "'a::euclidean_space  $\Rightarrow$  'b::euclidean_space"
  assumes S: "compact S"
    and f: "continuous_on S f"
    and e: "0 < e"
  shows " $\exists g. \text{polynomial\_function } g \wedge (\forall x \in S. \text{norm}(f\ x - g\ x) < e)$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a continuous function on a compact set S . Then there exists a polynomial function g such that $\forall x \in S, |f(x) - g(x)| < \epsilon$."

755

```

theorem Riemann_mapping_theorem:
  "open S  $\wedge$  simply_connected S  $\longleftrightarrow$ 
  S = {}  $\vee$  S = UNIV  $\vee$ 
  ( $\exists f\ g. f \text{ holomorphic\_on } S \wedge g \text{ holomorphic\_on ball } 0\ 1 \wedge$ 
    ( $\forall z \in S. f\ z \in \text{ball } 0\ 1 \wedge g(f\ z) = z$ )  $\wedge$ 
    ( $\forall z \in \text{ball } 0\ 1. g\ z \in S \wedge f(g\ z) = z$ ))"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let S be a simply connected open set. Then S is either empty, the whole plane, or there exists a holomorphic function f from S to the unit disk and a holomorphic function g from the unit disk to S such that $g \circ f$ is the identity function on S and $f \circ g$ is the identity function on the unit disk."

756


```

theorem Residue_theorem:
  fixes s pts::"complex set" and f::"complex  $\Rightarrow$  complex"
  and g::"real  $\Rightarrow$  complex"
  assumes "open s" "connected s" "finite pts" and
    holo:"f holomorphic_on s-pts" and
    "valid_path g" and
    loop:"pathfinish g = pathstart g" and
    "path_image g  $\subseteq$  s-pts" and
    homo:" $\forall z. (z \notin s) \longrightarrow \text{winding\_number } g \ z = 0$ "
  shows "contour_integral g f =
    2 * pi * i * ( $\sum_{p \in \text{pts.}} \text{winding\_number } g \ p * \text{residue } f \ p$ )"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let f be a holomorphic function on a connected open set S and let g be a closed path in S that does not pass through any of the points in the finite set P . Then $\int_g f(z)dz = 2\pi i \sum_{p \in P} \text{Res}(f, p)$ ".

757

```

theorem integer_winding_number:
  "[[path  $\gamma$ ; pathfinish  $\gamma$  = pathstart  $\gamma$ ;
   $z \notin \text{path\_image } \gamma$ ]  $\implies \text{winding\_number } \gamma \ z \in \mathbb{Z}$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let γ be a closed path in the complex plane, and z be a point not on γ . Then the winding number of γ around z is an integer".

758

```

lemma winding_number_eq:
  "[[path  $\gamma$ ; pathfinish  $\gamma$ 
  = pathstart  $\gamma$ ;  $w \in S$ ;  $z \in S$ ; connected  $S$ ;
   $S \cap \text{path\_image } \gamma = \{z\}$ ]
 $\implies \text{winding\_number } \gamma \ w = \text{winding\_number } \gamma \ z$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let γ be a closed path in the plane, and let w, z be two points in the plane. If γ does not pass through w or z , and w, z are in the same connected component of the plane, then the winding number of γ about w is the same as the winding number of γ about z ".

759

```

corollary winding_number_zero_outside:
  "[[path  $\gamma$ ; convex  $s$ ; pathfinish  $\gamma$ 
  = pathstart  $\gamma$ ;  $z \notin s$ ; path_image  $\gamma$ 
 $\subseteq s$ ]  $\implies \text{winding\_number } \gamma \ z = 0$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If γ is a path, s is a convex set, γ is a closed path, z is not in s , and γ is in s , then the winding number of γ around z is 0".

760

```

lemma winding_number_zero_at_infinity:
  assumes  $\gamma$ : "path  $\gamma$ " and loop: "pathfinish  $\gamma$  = pathstart  $\gamma$ "
  shows " $\exists B. \forall z. B \leq \text{norm } z \longrightarrow \text{winding\_number } \gamma \ z = 0$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let γ be a closed path, then there exists a constant B such that for all z , if $|z| > B$, then the winding number of γ at z is 0".

761

```

lemma winding_number_homotopic_paths:
  assumes "homotopic_paths  $(-\{z\})$   $g$   $h$ "
  shows "winding_number  $g$   $z$  = winding_number  $h$   $z$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If g and h are homotopic paths in $\mathbb{C} - \{z\}$, then the winding number of g around z is equal to the winding number of h around z ".

762

```

lemma simple_closed_path_winding_number_cases:
  assumes "simple_path  $\gamma$ " "pathfinish  $\gamma$  = pathstart  $\gamma$ "
  "z  $\notin$  path_image  $\gamma$ "
  shows "winding_number  $\gamma$   $z \in \{-1, 0, 1\}$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If γ is a simple closed path, z is not in the image of γ , then the winding number of γ at z is either -1, 0 or 1".

763

```

corollary Cauchy_theorem_primitive:
  assumes " $\bigwedge x. x \in S \implies$ "
    "(f has_field_derivative f' x) (at x within S)"
  and "valid_path  $g$ " "path_image  $g \subseteq S$ "
    "pathfinish  $g$  = pathstart  $g$ "
  shows "(f' has_contour_integral 0)  $g$ "

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If f is a function with a derivative f' on a set S , and g is a closed path in S , then $\int_g f'(z)dz = 0$ ".

764

```

theorem (in prob_space) central_limit_theorem_zero_mean:
  fixes X :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
  and  $\mu$  :: "real measure"
  and  $\sigma$  :: real
  and S :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
  assumes X_indep: "indep_vars ( $\lambda$ i. borel) X UNIV"
  and X_mean_0: " $\bigwedge$ n. expectation (X n) = 0"
  and  $\sigma$ _pos: " $\sigma > 0$ "
  and X_square_integrable: " $\bigwedge$ n. integrable M ( $\lambda$ x. (X n x)2)"
  and X_variance: " $\bigwedge$ n. variance (X n) =  $\sigma^2$ "
  and X_distrib: " $\bigwedge$ n. distr M borel (X n) =  $\mu$ "
  defines "S n  $\equiv$   $\lambda$ x.  $\sum$  i<n. X i x"
  shows "weak_conv_m ( $\lambda$ n. distr M borel
    ( $\lambda$ x. S n x / sqrt (n *  $\sigma^2$ ))) std_normal_distribution"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"Let X_1, X_2, \dots be a sequence of independent random variables with mean 0 and variance σ^2 . Let $S_n = \sum_{i=1}^n X_i$. Then $S_n/\sqrt{n\sigma^2}$ converges in distribution to a standard normal distribution".

765

```

proposition compact_space_product_topology:
  "compact_space(product_topology X I)  $\longleftrightarrow$ 
   topspace(product_topology X I) = { }  $\vee$  ( $\forall$  i  $\in$  I. compact_space(X i))"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If X_i is a compact space for all $i \in I$, then $X_1 \times X_2 \times \dots \times X_n$ is a compact space".

766

767 Lemma statements in the remainder of this section are from the *Topology* entry in the AFP.³ All the
 768 natural language version is generated by the Codex model.

```

lemma seteqI:
  "[ $\bigwedge$  x. x  $\in$  A  $\implies$  x  $\in$  B;  $\bigwedge$  x. x  $\in$  B  $\implies$  x  $\in$  A ]  $\implies$  A = B"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If A and B are sets such that $A \subseteq B$ and $B \subseteq A$, then $A = B$."

769

```

lemma (in carrier) openI: "m  $\in$  T  $\implies$  m open"

```

Translate the Isabelle version to a natural language version:

Codex Output:

"If m is a topology on a set T , then m is open."

770

³https://www.isa-afp.org/browser_info/current/AFP/Topology/Topology.html

lemma (in carrier) openE: " $\llbracket m \text{ open}; m \in T \implies R \rrbracket \implies R$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If m is an open set in a topological space T , then $m \in T$."

771

lemma (in topology) open_kriterion:
 assumes t_contains_open: " $\bigwedge x. x \in t \implies \exists t'. t' \text{ open} \wedge x \in t' \wedge t' \subseteq t$ "
 shows " $t \text{ open}$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If t is a subset of a topological space T such that for every $x \in t$, there exists an open set t' such that $x \in t'$ and $t' \subseteq t$, then t is open."

772

lemma (in carrier) carrierI [intro]:
 " $\llbracket t \text{ open}; x \in t \rrbracket \implies x \in \text{carrier}$ "

Translate the Isabelle version to a natural language version: "If t is an open set in a topological space T and $x \in t$, then $x \in T$."

773

lemma (in topology) Un_open [intro!]:
 assumes abopen: " $A \text{ open}$ " " $B \text{ open}$ "
 shows " $A \cup B \text{ open}$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If A and B are open sets in a topological space T , then $A \cup B$ is open."

774

lemma topologyI:
 " $\llbracket \bigwedge x y. \llbracket \text{is_open } T x; \text{is_open } T y \rrbracket \implies \text{is_open } T (x \cap y);$
 $\bigwedge M. \forall m \in M. \text{is_open } T m \implies \text{is_open } T (\bigcup M)$
 $\rrbracket \implies \text{topology } T$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If T is a set and T is closed under finite intersections and arbitrary unions, then T is a topology."

775

lemma (in topobase) topo_open:
 " $t \text{ open} = (t \in \text{topo } B)$ "

Translate the Isabelle version to a natural language version:

Codex Output:

"If B is a topological base of a topological space T , then t is open if and only if $t \in B$."

776

```

lemma (in topobase) topo_induct
  [case_names basic inter union, induct set: topo, consumes 1]:
  assumes opn: "x open"
    and bas: " $\bigwedge x. x \in B \implies P\ x$ "
    and int: " $\bigwedge x\ y. [x\ \text{open}; P\ x; y\ \text{open}; P\ y] \implies P\ (x \cap y)$ "
    and uni: " $\bigwedge M. (\forall t \in M. t\ \text{open} \wedge P\ t) \implies P\ (\bigcup M)$ "
  shows "P x"

```

Translate the Isabelle version to an natural language version:

Codex Output:

"If B is a topological base of a topological space T , then $P(x)$ holds for all open sets x if $P(x)$ holds for all $x \in B$ and $P(x)$ is closed under finite intersections and arbitrary unions."

777

```

lemma subtopologyI:
  fixes S (structure)
  fixes T (structure)
  assumes H1: " $\bigwedge s. s\ \text{open} \implies \exists t. t\ \text{open}_T \wedge s = t \cap \text{carrier}$ "
    and H2: " $\bigwedge t. t\ \text{open}_T \implies t \cap \text{carrier}\ \text{open}$ "
  shows "subtopology S T"

```

Translate the Isabelle version to an natural language version:

Codex Output:

"If S is a subset of a topological space T such that for every open set s in S , there exists an open set t in T such that $s = t \cap S$ and for every open set t in T , $t \cap S$ is open in S , then S is a subspace of T ."

778

```

lemma (in subtopology) subtopologyE [elim]:
  assumes major: "s open"
    and minor: " $\bigwedge t. [t\ \text{open}_T; s = t \cap \text{carrier}] \implies R$ "
  shows "R"

```

Translate the Isabelle version to an natural language version:

Codex Output:

"If s is an open set in a topological space T , then s is an open set in a subspace S of T if and only if $s = t \cap S$ for some open set t in T ."

779

```

lemma continuousI:
  fixes S (structure)
  fixes T (structure)
  assumes "f :  $\text{carrier}_S \rightarrow \text{carrier}_T$ "
    " $\bigwedge m. m\ \text{open}_T \implies \text{carrier} \cap (f^{-1} m)\ \text{open}$ "
  shows "continuous f S T"

```

Translate the Isabelle version to an natural language version:

Codex Output:

"If f is a function from a topological space S to a topological space T such that the preimage of every open set in T is open in S , then f is continuous."

780

```

lemma continuousE:
  fixes S (structure)
  fixes T (structure)
  shows
    "[continuous f S T;
      [f : carrierS → carrierT;
        ∀ m. m openT → carrierS ∩ (f -' m) open ] ⇒ P
    ] ⇒ P"

```

Translate the Isabelle version to an natural language version:

Codex Output:

"If x is a limit point of a set A in a topological space T , then there exists a filter F such that $x \in F$ and $A \in F$."

781