

## 585 A Details of Numerical Reasoning

586 In this section, we describe the data processing workflow (including building the testing datasets and  
587 synthesizing training data) and list the prompts used for different methods. For a fair comparison,  
588 we use the same prompts for CoT and the reasoning mode of ToolkenGPT. With the same example  
589 questions, we label the calculation process in place to get the prompts for ReAct<sup>2</sup>. For the tool mode  
590 of ToolkenGPT, we randomly sample 4 examples of the specified tool from the training set, and  
591 transform them into ReAct-style prompts. Because a large number of tools are used, we show the  
592 prompts in the supplementary file instead of listing them here.

### 593 A.1 GSM8K-XL

#### 594 A.1.1 Level-up Strategy

595 To magnify the numbers for the GSM8K to build our dataset GSM8K-XL, we follow the steps  
596 outlined below:

- 597 1. We prompt the gpt-3.5-turbo with two examples to replace the numbers with appropriate  
598 placeholders. The prompt is presented below.
- 599 2. In order to validate the correctness of the number replacements, we develop a solving  
600 function for the GSM8k dataset, which utilizes the formulas presented in the provided  
601 answers and produces the expected solutions. By substituting the original numbers back  
602 into the rewritten question-answer pairs and executing the solving function, we can identify  
603 any erroneous rewrites and exclude them from our dataset. This ensures the reliability and  
604 accuracy of the number replacements.
- 605 3. Once we have confirmed the accuracy of the rewriting process, we can proceed with mag-  
606 nifying the numbers. We systematically replace the original numbers with their respective  
607 cubic values, exponentially amplifying their magnitudes. However, it is essential to exercise  
608 caution in certain scenarios. For instance, when dealing with percentage values (e.g., 10%),  
609 in most cases, a direct cubic magnification to 1000% would be inappropriate. In such cases,  
610 we maintain the original values to preserve their intended meaning.
- 611 4. With the magnified numbers in place, we seamlessly substitute the placeholders within the  
612 rewritten question-answer pairs. By employing the solving function developed earlier, we  
613 can process the magnified formulas and obtain accurate solutions, leading to the GSM8K-XL  
614 dataset.

```
615 {  
616   "role": "system",  
617   "content": "Your task is to rewrite the given question and answer in a  
618     more general form using variables instead of specific numbers."  
619 },  
620 {  
621   "role": "user",  
622   "content": "Question: Janet's ducks lay 16 eggs per day. She eats three  
623     for breakfast every morning and bakes muffins for her friends every  
624     day with four. She sells the remainder at the farmers' market daily  
625     for $2 per fresh duck egg. How much in dollars does she make every  
626     day at the farmers' market\nAnswer: Janet sells 16 - 3 - 4 =  
627     <<16-3-4=9>>9 duck eggs a day.\nShe makes 9 * 2 = $<<9*2=18>>18  
628     every day at the farmer's market.\n#### 18"  
629 },  
630 {  
631   "role": "assistant",  
632   "content": "Question: Janet's ducks lay {v_1} eggs per day. She eats  
633     {v_2} for breakfast every morning and bakes muffins for her friends
```

<sup>2</sup>As Yao et al. [62] doesn't evaluate ReAct on numerical reasoning, we don't follow the exact format (i.e. trajectories marked by [THINK] and [ACT]), but design a more natural way to combine CoT reasoning and tool calling (as shown by our prompts).

```

634     every day with {v_3}. She sells the remainder at the farmers'
635     market daily for ${v_4} per fresh duck egg. How much in dollars
636     does she make every day at the farmers' market?\nAnswer: Janet
637     sells {v_1} - {v_2} - {v_3} = <<{v_1}-{v_2}-{v_3}={r_1}>>{r_1} duck
638     eggs a day.\nShe makes {r_1} * {v_4} = $<<{r_1}*{v_4}={r_2}>>{r_2}
639     every day at the farmer's market.\n#### {r_2}"
640 },
641 {
642     "role": "user",
643     "content": "Question: A wooden bridge can carry no more than 5000
644     pounds. A delivery truck filled with identical boxes, each weighing
645     15 pounds, will pass over the bridge. The combined weight of the
646     driver and the empty truck is 3755 pounds. What is the maximum
647     number of boxes which can be loaded onto the truck while not
648     exceeding the bridge's weight limit?\nAnswer: The boxes can weigh
649     up to 5000 pounds - 3755 pounds = <<5000-3755=1245>>1245 pounds in
650     total.\nThere can be 1245 / 15 = <<1245/15=83>>83 boxes loaded onto
651     the truck without exceeding the bridge's weight limit.\n#### 83"
652 },
653 {
654     "role": "assistant",
655     "content": "Question: A wooden bridge can carry no more than {v_1}
656     pounds. A delivery truck filled with identical boxes, each weighing
657     {v_2} pounds, will pass over the bridge. The combined weight of the
658     driver and the empty truck is {v_3} pounds. What is the maximum
659     number of boxes which can be loaded onto the truck while not
660     exceeding the bridge's weight limit?\nAnswer: The boxes can weigh
661     up to {v_1} pounds - {v_3} pounds = <<{v_1}-{v_3}={r_1}>>{r_1}
662     pounds in total.\nThere can be {r_1} / {v_2} =
663     <<{r_1}/{v_2}={r_2}>>{r_2} boxes loaded onto the truck without
664     exceeding the bridge's weight limit.\n#### {r_2}"
665 },
666 {
667     "role": "user",
668     "content": [INPUT]
669 }

```

## 670 A.1.2 Training Details

671 As mentioned in Section 4.1, the Toolken embeddings are trained with a subset of 5,063 examples. An  
672 additional 1,000 examples are reserved for validation. The embeddings was trained with a learning  
673 rate of  $5e-4$ , performing early stopping based on the development set, with a maximum of 10 epochs.

## 674 A.1.3 Prompt for GSM8K-XL Dataset

675 Prompt for Direct Prompting with ChatGPT:

676 Solve the following math problem step by step, and then provide the final  
677 answer in the format: 'So, the answer is xxx.'

678

679 [QUESTION]

680 Prompt for Chain of Thought (CoT) and ToolkenGPT:

681 Answer the following questions step by step.

682

683 Question: Mark has 3 tanks for pregnant fish. Each tank has 4 pregnant  
684 fish and each fish gives birth to 20 young. How many young fish does he  
685 have at the end?

686 Answer: He has  $4 \times 3 = 12$  pregnant fish They give birth to  $12 \times 20 = 240$  fish ####  
687 240

688

689 Question: The math questions in a contest are divided into three rounds:  
690 easy, average, and hard. There are corresponding points given for each  
691 round. That is 2, 3, and 5 points for every correct answer in the easy,  
692 average, and hard rounds, respectively. Suppose Kim got 6 correct answers  
693 in the easy; 2 correct answers in the average; and 4 correct answers in  
694 the difficult round, what are her total points in the contest?

695 Answer: Kim got 6 points/round x 2 round = 12 points in the easy round.  
696 She got 2 points/round x 3 rounds = 6 points in the average round. She got  
697 4 points/round x 5 rounds = 20 points in the difficult round. So her total  
698 points is 12 points + 6 points + 20 points = 38 points. #### 38

699

700 Question: A clothing store sells 20 shirts and 10 pairs of jeans. A shirt  
701 costs \$10 each and a pair of jeans costs twice as much. How much will the  
702 clothing store earn if all shirts and jeans are sold?

703 Answer: Twenty shirts amount to  $\$10 \times 20 = \$200$ . The cost of each pair of  
704 jeans is  $\$10 \times 2 = \$20$ . So 10 pairs of jeans amount to  $\$20 \times 10 = \$200$ .  
705 Therefore, the store will earn  $\$200 + \$200 = \$400$  if all shirts and jeans  
706 are sold. #### 400

707

708 Question: Arnold's collagen powder has 18 grams of protein for every 2  
709 scoops. His protein powder has 21 grams of protein per scoop. And his  
710 steak has 56 grams of protein. If he has 1 scoop of collagen powder, 1  
711 scoop of protein powder and his steak, how many grams of protein will he  
712 consume?

713 Answer: 2 scoops of collagen powder have 18 grams of protein and he only  
714 has 1 scoop so he consumes  $18/2 = 9$  grams of protein He has 9 grams  
715 collagen powder, 21 grams of protein powder and 56 grams in his steak for  
716 a total of  $9 + 21 + 56 = 86$  grams of protein #### 86

717

718 Question: [QUESTION]

719 Answer:

720 Prompt for ReAct:

721 Answer the following questions with <add>, <subtract>, <multiply>,  
722 <divide> operators

723

724 Question: Mark has 3 tanks for pregnant fish. Each tank has 4 pregnant  
725 fish and each fish gives birth to 20 young. How many young fish does he  
726 have at the end?

727 Answer: He has  $4 \times 3 = \langle \text{multiply} \rangle (4, 3) = 12$  pregnant fish They give birth to  
728  $12 \times 20 = \langle \text{multiply} \rangle (12, 20) = 240$  fish #### 240

729

730 Question: The math questions in a contest are divided into three rounds:  
731 easy, average, and hard. There are corresponding points given for each  
732 round. That is 2, 3, and 5 points for every correct answer in the easy,  
733 average, and hard rounds, respectively. Suppose Kim got 6 correct answers  
734 in the easy; 2 correct answers in the average; and 4 correct answers in  
735 the difficult round, what are her total points in the contest?

736 Answer: Kim got 6 points/round x 2 round =  $\langle \text{multiply} \rangle (6, 2) = 12$  points in  
737 the easy round. She got 2 points/round x 3 rounds =  $\langle \text{multiply} \rangle (2, 3) = 6$   
738 points in the average round. She got 4 points/round x 5 rounds =  
739  $\langle \text{multiply} \rangle (4, 5) = 20$  points in the difficult round. So her total points is  
740  $12 \text{ points} + 6 \text{ points} + 20 \text{ points} = \langle \text{add} \rangle (12, 6, 20) = 38$  points. #### 38

741

742 Question: A clothing store sells 20 shirts and 10 pairs of jeans. A shirt  
743 costs \$10 each and a pair of jeans costs twice as much. How much will the  
744 clothing store earn if all shirts and jeans are sold?

745 Answer: Twenty shirts amount to  $\$10 \times 20 = \$\langle\text{multiply}\rangle(10, 20)=200$ . The  
746 cost of each pair of jeans is  $\$10 \times 2 = \$\langle\text{multiply}\rangle(10, 2)=20$ . So 10 pairs  
747 of jeans amount to  $\$20 \times 10 = \$\langle\text{multiply}\rangle(20, 10)=200$ . Therefore, the  
748 store will earn  $\$200 + \$200 = \$\langle\text{add}\rangle(200, 200)=400$  if all shirts and jeans  
749 are sold. ##### 400

750

751 Question: Arnold's collagen powder has 18 grams of protein for every 2  
752 scoops. His protein powder has 21 grams of protein per scoop. And his  
753 steak has 56 grams of protein. If he has 1 scoop of collagen powder, 1  
754 scoop of protein powder and his steak, how many grams of protein will he  
755 consume?

756 Answer: 2 scoops of collagen powder have 18 grams of protein and he only  
757 has 1 scoop so he consumes  $18/2 = \langle\text{divide}\rangle(18, 2)=9$  grams of protein He  
758 has 9 grams collagen powder, 21 grams of protein powder and 56 grams in  
759 his steak for a total of  $9+21+56 = \langle\text{add}\rangle(9, 21, 56)=86$  grams of protein  
760 ##### 86

761

762 Question: [QUESTION]

763 Answer:

## 764 A.2 FuncQA

### 765 A.2.1 Training Details

766 As mentioned in Section 4.1, Toolken embeddings are trained using a subset of 611 examples, with  
767 an additional 39 examples reserved for development/validation purposes. The learning rate we use is  
768  $1e-4$ , and we perform early stopping based on the development set, with the maximal training epochs  
769 to be 20.

### 770 A.2.2 Prompt for Synthetic Training Data

771 In Section 4.1, we discussed the utilization of ChatGPT for synthesizing the training set. To create the  
772 training data, we begin by manually crafting two examples that adhere to the desired format, and then  
773 use the follow specific prompt to generate more examples. However, it is important to acknowledge  
774 that the prompt does not guarantee the generation of examples that strictly conform to the required  
775 format. So, we apply a filtering process to remove any non-conforming instances. Furthermore, the  
776 generation process often produces duplicate examples, necessitating a subsequent de-duplication  
777 step.

778 You are a math question generator for teachers, and your task is to  
779 generate some questions and answers using function [FUNC] to solve and can  
780 be solved within one single step. You do not need to give specific  
781 numbers, so that the teachers can fill any numbers they want. Here are two  
782 examples that use the function [FUNC].

783

784 [EXAMPLE\_1]

785

786 [EXAMPLE\_2]

787

788 [FUNC] is a function to [DESCRIPTION]. Now, let's mimic the format of  
789 examples to generate various real world QA pairs using the function [FUNC]  
790 that can be solved within one step. The numbers should be replaced by  
791 [ARG] and [ANSWER] as the examples given above.

### 792 A.2.3 Prompt for FuncQA One-Hop

793 Prompt for Zero-Shot with ChatGPT:

794 Solve the following math problem, and then provide the final answer in the  
795 format: 'So, the answer is xxx.'

796

797 [QUESTION]

798 Prompt for Chain of Thought (CoT) and ToolkenGPT:

799 Q: If Amy's income increases by 4% annually, how many times will it  
800 multiply in 11 years?

801 A: In 11 years, Amy's income will increase by  $1.04^{11}=1.54$  times. So, the  
802 answer is 1.54.

803

804 Q: If a store sells 147 bananas today and 354 more bananas tomorrow, how  
805 many bananas does the store sell in total?

806 A: The store sells 147 bananas today and 354 more bananas tomorrow, so the  
807 total number of bananas sold is  $147+354=501$ . So, the answer is 501.

808

809 Q: A man had 789.4 dollars in his wallet. He spent 11.99 dollars on a  
810 movie ticket. How much money does he have left now?

811 A: The man had 789.4 dollars in his wallet and spent 11.99 dollars on a  
812 movie ticket, so he has  $789.4-11.99=777.41$  dollars left. So, the answer is  
813 777.41 dollars.

814

815 Q: If a cake weighs 3.77 pounds and is divided into 13 equal pieces, how  
816 much does each piece weight?

817 A: Each piece of the cake weighs  $3.77/13=0.29$  pounds. So, the answer is  
818 0.29 pounds.

819

820 Q: [QUESTION]

821 A:

822 Prompt for ReAct:

823 Answer the following question with <add>, <subtract>, <multiply>,  
824 <divide>, <power>, <sqrt>, <log>, <lcm>, <gcd>, <ln>, <choose>,  
825 <remainder>, <permutate>:

826

827 Q: If Amy's income increases by 4% annually, how many times will it  
828 multiply in 11 years?

829 A: In 11 years, Amy's income will increase by  $1.04^{11} =$   
830 <power>(1.04,11)=1.54 times. So, the answer is 1.54.

831

832 Q: If a store sells 147 bananas today and 354 more bananas tomorrow, how  
833 many bananas does the store sell in total?

834 A: The store sells 147 bananas today and 354 more bananas tomorrow, so the  
835 total number of bananas sold is  $147+354=<add>(147,354)=501$ . So, the answer  
836 is 501.

837

838 Q: A man had 789.4 dollars in his wallet. He spent 11.99 dollars on a  
839 movie ticket. How much money does he have left now?

840 A: The man had 789.4 dollars in his wallet and spent 11.99 dollars on a  
841 movie ticket, so he has  $789.4-11.99=<subtract>(789.4,11.99)=777.41$  dollars  
842 left. So, the answer is 777.41.

843

844 Q: If a cake weighs 3.77 pounds and is divided into 13 equal pieces, how  
845 much does each piece weight?

846 A: Each piece of the cake weighs  $3.77/13=<divide>(3.77,13)=0.29$  pounds.  
847 So, the answer is 0.29.

848

849 Q: [QUESTION]

850 A:

#### 851 A.2.4 Prompt for FuncQA Multi-Hop

852 Prompt for Zero-Shot with ChatGPT:

853 Solve the following math problem step by step, and then provide the final  
854 answer in the format: 'So, the answer is xxx.'

855

856 [QUESTION]

857 Prompt for Chain of Thought (CoT) and ToolkenGPT:

858 Answer the following questions step by step:

859

860 Question: A coin is tossed 8 times, what is the probability of getting  
861 exactly 7 heads ?

862 Answer: The total number of possible outcomes to toss a coin 8 times is  
863  $2^8=256$ . The number of ways of getting exactly 7 heads is  $8C7=8$ . The  
864 probability of getting exactly 7 heads is  $8/256=0.03125$ . ##### 0.03125

865

866 Question: If paint costs \$3.2 per quart, and a quart covers 12 square  
867 feet, how much will it cost to paint the outside of a cube 10 feet on each  
868 edge?

869 Answer: The total surface area of the 10 ft cube is  $6*10^2=6*100=600$   
870 square feet. The number of quarts needed is  $600/12=50$ . The cost is  
871  $50*3.2=160$ . ##### 160

872

873 Question:  $\log(x)=2$ ,  $\log(y)=0.1$ , what is the value of  $\log(x-y)$  ?

874 Answer:  $\log(x)=2$ , so  $x=10^2=100$ ;  $\log(y)=0.1$ , so  $y=10^{0.1}=1.26$ ;  
875  $x-y=100-1.26=98.74$ , so  $\log(x-y)=\log(98.74)=1.99$ . ##### 1.99

876

877 Question: How many degrees does the hour hand travel when the clock goes  
878 246 minutes?

879 Answer: The hour hand travels 360 degrees in 12 hours, so every hour it  
880 travels  $360/12=30$  degrees. 246 minutes is  $246/60=4.1$  hours. The hour hand  
881 travels  $4.1*30=123$  degrees. ##### 123

882

883 Question: [QUESTION]

884 Answer:

885 Prompt for ReAct:

886 Answer the following questions with <add>, <subtract>, <multiply>,  
887 <divide>, <power>, <sqrt>, <log>, <lcm>, <gcd>, <ln>, <choose>,  
888 <remainder>, and <permutate>:

889

890 Question: A coin is tossed 8 times, what is the probability of getting  
891 exactly 7 heads?

892 Answer: The total number of possible outcomes to toss a coin 8 times is  
893  $2^8=<power>(2,8)=256$ . The number of ways of getting exactly 7 heads is  
894  $8C7=<choose>(8,7)=8$ . The probability of getting exactly 7 heads is  
895  $8/256=<divide>(8,256)=0.03125$ . ##### 0.03125

896

897 Question: If paint costs \$3.2 per quart, and a quart covers 12 square  
898 feet, how much will it cost to paint the outside of a cube 10 feet on each  
899 edge?

900 Answer: The total surface area of the 10 ft cube is  
901  $6*10^2=6*<power>(10,2)=100=<multiply>(6,100)=600$  square feet. The number  
902 of quarts needed is  $600/12=<divide>(600,12)=50$ . The cost is  
903  $50*3.2=<multiply>(50,3.2)=160$ . ##### 160

904  
905 Question:  $\log(x)=2$ ,  $\log(y)=0.1$ , what is the value of  $\log(x-y)$  ?  
906 Answer:  $\log(x)=2$ , so  $x=10^2=\text{power}(10,2)=100$ ;  $\log(y)=0.1$ , so  
907  $y=10^{0.1}=\text{power}(10,0.1)=1.26$ ;  $x-y=100-1.26=\text{subtract}(100,1.26)=98.74$ , so  
908  $\log(x-y)=\log(98.74)=\text{log}(98.74)=1.99$ . #### 1.99  
909  
910 Question: How many degrees does the hour hand travel when the clock goes  
911 246 minutes?  
912 Answer: The hour hand travels 360 degrees in 12 hours, so every hour it  
913 travels  $360/12=\text{divide}(360,12)=30$  degrees. 246 minutes is  
914  $246/60=\text{divide}(246,60)=4.1$  hours. The hour hand travels  
915  $4.1*30=\text{multiply}(4.1,30)=123$  degrees. #### 123  
916  
917 Question: [QUESTION]  
918 Answer:

## 919 B Details of Knowledge-based QA

920 For KAMEL, we first transform each wikidata relation identifier (e.g. P1346) into a natural language  
921 description. Note that the natural language descriptions are not necessary for ToolkenGPT which  
922 is directly trained with demonstrations, but they are crucial for the in-context learning baselines,  
923 especially for ICL (desc), which can only understand tools through descriptions. We then describe  
924 the process of synthesizing training data.

### 925 B.1 Getting Text Description

926 KAMEL provides a question template for each relations. We randomly sample 3 facts from the  
927 dataset and instantiate them into question-answer pair, and use the following prompt to generate a  
928 description for them with ChatGPT:

929 Given a question template and some example answer, you need to define an  
930 API that can help you answer the question.  
931 Q 1.1: What is the original language of The Wonderful Galaxy of Oz  
932 A 1.1: Japanese  
933 Q 1.2: What is the original language of Wild Field?  
934 A 1.2: Russian  
935 Q 1.3: What is the original language of Nadigan?  
936 A 1.3: Tamil  
937 API 1: original\_language(title): gets the original language of an art work  
938 Q 2.1: What languages does Judah Maccabee speak?  
939 A 2.1: Hebrew  
940 Q 2.2: What languages does Ronelda Kamfer speak?  
941 A 2.2: Afrikaans  
942 Q 2.3: What languages does Leibush Lehrer speak?  
943 A 2.3: Yiddish  
944 API 2: spoken\_languages(name): gets the spoken languages of a person  
945 Q 3.1: [Q1]  
946 A 3.1: [A1]  
947 Q 3.2: [Q2]  
948 A 3.2: [A2]  
949 Q 3.3: [Q3]  
950 A 3.3: [A3]  
951 API 3:

### 952 B.2 Synthetic Data

953 We use two prompts to synthesize diverse training data, and aggregate the samples from each prompt.

954 Here are some examples of using functions for text generation (after the  
955 function call, the sentence should continue with the returned value of the  
956 function call):  
957 1. `star_rating(product)`: gets the star rating of the product on a scale  
958 from 0 to 5.  
959 Example 1.1: The new iPhone 12 Pro Max is already generating a lot of  
960 buzz, thanks to its `<f>star_rating("iPhone 12 Pro Max")="4.7"</f>`4.7 star  
961 rating.  
962 2. `literary_genre(book)`: gets the literary genre of a book  
963 Example 2.1: Literature is often categorized by genre, such as drama,  
964 romance, or science fiction. The Harry Potter series is a popular example  
965 of the `<f>literary_genre("Harry Potter")="fantasy"</f>`fantasy genre, which  
966 features imaginary worlds and magical elements.  
967 3. `current_location(user_id)`: gets the current location of a user.  
968 Example 3.1: If you're trying to coordinate plans with a friend, it's  
969 helpful to know their current location. You can ask the question "Where  
970 are you right now?" and use the function `<f>current_location("1234")="New  
971 York"</f>`New York as an example response.  
972 4. `number_of_movies(director)`: gets the number of movies directed by a  
973 specific director.  
974 Example 4.1: Martin Scorsese is one of the most celebrated movie directors  
975 of all time. He has directed a total of `<f>number_of_movies("Martin  
976 Scorsese")="78"</f>`78 movies throughout his career.  
977 5. `word_definition(word)`: gets the definition of a particular word  
978 Example 5.1: Writers and English language learners can enhance their  
979 vocabulary by knowing the definition of unfamiliar words. The definition  
980 of the word "eccentric" is `<f>word_definition("eccentric")="unconventional  
981 and slightly strange"</f>`unconventional and slightly strange.  
982 6. `number_of_spotify_followers(artist)`: gets the number of Spotify  
983 followers for the artist.  
984 Example 6.1: Taylor Swift's latest album is a hit and her fan base is  
985 growing rapidly. In fact, her number of Spotify followers as of today is  
986 `<f>number_of_spotify_followers("Taylor Swift")="49,879,220"</f>`49,879,220.  
987 6. [DESCRIPTION]  
988 Please continue to generate 10 examples using the function [NAME],  
989 starting with 6.1 to 6.10.

990 Here are some examples of using functions for text generation (after the  
991 function call, the sentence should continue with the returned value of the  
992 function call):  
993 1. `current_weather(city)`: gets the current weather of a city.  
994 Example 1.1: What's the weather in Beijing now? The weather is  
995 `<f>current_weather("Beijing")="sunny"</f>`sunny now. Example 1.2: Do you  
996 know what's the weather in San Diego now? The weather is  
997 `<f>current_weather("San Diego")="cloudy"</f>`cloudy now.  
998 2. `calculator(formula)`: gets the calculation result of a formula.  
999 Example 2.1: What's sum of 213 and 5032? The answer is  
1000 `<f>calculator("213+5032")="5245"</f>`5245.  
1001 Example 2.2: What's difference between 2015 and 33? The answer is  
1002 `<f>calculator("2015-33")="1982"</f>`1982.  
1003 3. [DESCRIPTION]  
1004 Please continue to generate 10 examples using the function [NAME],  
1005 starting with 3.1 to 3.10.

### 1006 B.3 Training Details

1007 Toolken embeddings are trained with a learning rate of  $1e-4$ , performing early stopping based on the  
1008 development set, and trained for a maximum of 5 epochs.

## 1009 C Details of Embodied AI

### 1010 C.1 Preprocessing

1011 We collect all scripts from ActivityPrograms [21] and filter the dataset with the following steps: (1)  
1012 filter out all the scripts that are not executable, or don't cause any state changes in VirtualHome (2)  
1013 deduplicate the scripts with the same goal and instruction. (3) discard the script that involves two  
1014 different objects of the same name (4) find the verbs and objects that occur more than 10 times in the  
1015 data, and keep the scripts composed of only these verbs and objects.

1016 Note that our preprocessing is different from Huang et al. [21], where they regard a high-level goal as  
1017 a task. We treat two scripts with the same goal but different instructions as distinct tasks because  
1018 different instructions often indicate different action sequences, which may lead to different final  
1019 state graphs, e.g., for a high-level goal "Reading", some of the instructions mention "Turn on desk  
1020 lamp" while others don't. Huang et al. [21] relies on human annotation to evaluate the correctness  
1021 of the generated script, which actually lowers the difficulties of learning the environment, because  
1022 humans may assign a correct label as long as the plan looks "reasonable". On the contrary, we can  
1023 use the Evolving Graph <sup>3</sup> to strictly match the resulting state and ground truth state. This serves as an  
1024 automatic and more objective evaluation.

### 1025 C.2 Prompts

1026 We show the prompts for LLMs to generate plans below. Note that all methods use the same prompts  
1027 in this experiment.

```
1028 I am a household robot and I can take actions from '[FIND]', '[SIT]',  
1029 '[SWITCHON]', '[TURNTO]', '[LOOKAT]', '[TYPE]', '[WALK]', '[LIE]',  
1030 '[GRAB]', '[READ]', '[WATCH]', '[POINTAT]', '[TOUCH]', '[SWITCHOFF]',  
1031 '[OPEN]', '[PUSH]', '[PUTOBJBACK]', '[CLOSE]', '[DRINK]', '[RUN]',  
1032 '[DROP]', '[PULL]'.
```

1033  
1034 Task 1:

```
1035 I am in ['bathroom']. The objects I can manipulate are ['faucet',  
1036 'keyboard', 'television', 'coffe_maker', 'chair', 'button', 'pillow',  
1037 'phone', 'cup', 'couch', 'freezer', 'desk', 'oven', 'light', 'table',  
1038 'bedroom', 'dining_room', 'cupboard', 'computer', 'sink', 'mail', 'bed',  
1039 'mouse', 'home_office'].
```

1040 Goal:

1041 Write an email

1042 Hint:

1043 i went near the computer and turned it on. then sent the mail

1044 Plan:

```
1045 [WALK] <home_office>  
1046 [WALK] <table>  
1047 [FIND] <table>  
1048 [WALK] <table>  
1049 [FIND] <computer>  
1050 [TURNTO] <computer>  
1051 [LOOKAT] <computer>  
1052 [TURNTO] <computer>  
1053 [SWITCHON] <computer>  
1054 [FIND] <mail>  
1055 [TURNTO] <mail>
```

1056  
1057 Task 2:

```
1058 I am in ['home_office']. The objects I can manipulate are ['faucet',  
1059 'novel', 'keyboard', 'television', 'newspaper', 'chair', 'coffe_maker',  
1060 'pillow', 'phone', 'check', 'couch', 'freezer', 'desk', 'toothbrush',
```

---

<sup>3</sup><https://github.com/xavierpuigf/virtualhome>

```

1061 'oven', 'light', 'food_food', 'table', 'bookmark', 'bedroom',
1062 'dining_room', 'computer', 'sink', 'mail', 'bed', 'cat', 'mouse',
1063 'home_office', 'pot'].
1064 Goal:
1065 Work
1066 Hint:
1067 Find the computer. Turn it on by pressing the on button. Wait for it to
1068 load. Use the mouse and keyboard to perform your tasks on screen.
1069 Plan:
1070 [FIND] <computer>
1071 [SWITCHON] <computer>
1072 [FIND] <mouse>
1073 [TOUCH] <mouse>
1074 [FIND] <keyboard>
1075 [TOUCH] <keyboard>
1076
1077 Task 3:
1078 I am in ['bathroom']. The objects I can manipulate are ['dishwasher',
1079 'faucet', 'keyboard', 'television', 'newspaper', 'chair', 'coffe_maker',
1080 'pillow', 'phone', 'cup', 'check', 'couch', 'freezer', 'desk', 'oven',
1081 'light', 'food_food', 'plate', 'table', 'bookmark', 'bedroom',
1082 'dining_room', 'cupboard', 'computer', 'sink', 'bed', 'cat', 'mouse',
1083 'home_office', 'pot'].
1084 Goal:
1085 Pick up phone
1086 Hint:
1087 first when i hear the ringing sound i will run to my living room and picks
1088 up and i will say hello
1089 Plan:
1090 [RUN] <home_office>
1091 [WALK] <chair>
1092 [FIND] <chair>
1093 [SIT] <chair>
1094 [FIND] <phone>
1095 [GRAB] <phone>
1096
1097 Task 4:
1098 [QUESTION]

```

### 1099 C.3 Training Details

1100 Toolken embeddings are trained with a learning rate of 1e-4, performing early stopping based on the  
1101 development set, with a maximum of 10 epochs.

## 1102 D Computational Resources

1103 In terms of computational resources, we train and test ToolkenGPT based on LLaMA-13B and  
1104 LLaMA-30B using 2 and 4 Nvidia A3090 GPUs, respectively.

## 1105 E The License of the Assets

1106 In this study, we would like to emphasize that all data used are exclusively from official public  
1107 sources. We ensure strict compliance with ethical guidelines, data access rights, and intellectual  
1108 property regulations. We adhere to the license agreements and terms of use associated with each  
1109 dataset. In addition, we acknowledge and thank the original creators and providers of the datasets for  
1110 their valuable contributions to the research community.

1111 **F Safeguard Statement**

1112 In this paper, we primarily focus on the applications of mathematical, knowledge-based, and embodied  
1113 planning problems, posing no significant ethical or harmful concerns. We recognize that future  
1114 research on border applications of tool learning may pose a risk of misuse, and we recommend careful  
1115 considerations of all aspects of safety before it's applied in the real world.