# A  Appendix

## A.1  Example from PRONTOQA-OOD

Q: Everything that is a lorpus, a brimpus, or a jompus is a shumpus. Every wumpus is a vumpus and a sterpus and a brimpus. Everything that is a vumpus, a grimpus, or a brimpus is a lempus. Everything that is a lempus or a jompus or a lorpus is a dumpus. Vumpuses are rompuses. Every sterpus is a gorpus. Everything that is a vumpus, a grimpus, or a brimpus is a dumpus. Wumpuses are shumpuses. Polly is a rompus. Polly is a wumpus. Prove: Polly is a lempus or an impus or a yumpus.

A: Polly is a wumpus. Every wumpus is a vumpus and a sterpus and a brimpus. Polly is a vumpus and a sterpus and a brimpus. Polly is a brimpus. Polly is a vumpus, a grimpus, or a brimpus. Everything that is a vumpus, a grimpus, or a brimpus is a lempus. Polly is a lempus. Polly is a lempus, an impus, or a yumpus.

FIGURE 11: An example from PRONTOQA-OOD. This is a compositional example with a min depth of 4 and 3 rule types. The given answer is the expected chain-of-thought. The question is shown in blue, the query in red, and the chain-of-thought/answer in green.

## A.2  Generative process details

In this section, we describe the process to generate examples of each deduction rule.

***Implication elimination*** (i.e. *modus ponens*)    Given $f(c)$ and $\forall x(f(x) \rightarrow g(x))$, prove $g(c)$. These are the examples in the original PRONTOQA. We follow the same process here:

1. Generate an ontology. For simplicity, we generate linear ontologies, consisting of a collection of concepts, as well as subtype-supertype relations between those concepts (i.e. concept $f$ is a subtype of the supertype $g$ if every instance of $f$ is an instance of $g$). For simplicity, we limit each type to have at most one supertype.
2. Perform a random walk of length $k$ from a randomly selected start vertex, where $k$ is the desired proof depth.
3. Traverse the edges of the ontology and convert each into a sentence of the question.
4. Convert each step of the random walk into a sentence of the gold chain-of-thought.

Note that this process allows us to generate proofs of any depth, but the width is fixed to 1.

***Conjunction introduction***    Given $A$ and $B$, prove $A \wedge B$. The generative process is a modified version of that for implication elimination. Instead of generating rules of the form $\forall x(f(x) \rightarrow g(x))$, we generate rules of the form $\forall x(f_1(x) \wedge \ldots \wedge f_n(x) \rightarrow g(x))$, where $n$ is the proof width. Given, $f_1(c)$, ..., and $f_n(c)$, the model must first prove $f_1(c) \wedge \ldots \wedge f_n(c)$ before applying implication elimination to prove $g(c)$. To increase the depth of the proof, $g(c)$ itself can be part of a conjunct in the antecedent of another rule.

***Conjunction elimination***    Given $A \wedge B$, prove $A$. These examples are identical to those in conjunction introduction, except the conjunction appears in the consequent of each rule, rather than in the antecedent: $\forall x(f(x) \rightarrow g_1(x) \wedge \ldots \wedge g_n(x))$ where $n$ is the proof width.

***Disjunction introduction***    Given $A$, prove $A \vee B$. These examples are identical to those in conjunction introduction, except the conjunction is replaced with disjunction: $\forall x(f_1(x) \vee \ldots \vee f_n(x) \rightarrow g(x))$ where $n$ is the proof width. But note that grounded axioms are not necessary for every disjunct: To apply the rule $\forall x(f_1(x) \vee \ldots \vee f_n(x) \rightarrow g(x))$, knowing $f_n(c)$ is sufficient, and we do not need to generate grounded axioms for the other disjuncts $f_i(c)$ for $i < n$.

***Disjunction elimination*** (i.e. *proof by cases*)    Given $A_1 \vee \ldots \vee A_n$, and $A_i \vdash C$ for all $i$, prove $C$. Here, $n$ is the proof width. While it is possible to construct proofs containing multiple nested

14

applications of disjunction elimination, such proofs are quite complex, even for humans to understand, and so we fix the depth of these examples to 1. To generate an example, we first generate the disjunction: $f_1(c) \vee \ldots \vee f_n(c)$. Next, generate the rules for each case: $\forall x(f_i(x) \rightarrow g(x))$ for all $i$. The goal is to prove $g(c)$.

***Proof by contradiction***   Given $A \vdash B$ and $\neg B$, prove $\neg A$. Note that this is a rule composed of two natural deduction rules: negation elimination and introduction. But since those individual rules do not lend themselves to a natural text representation, we choose to study their composition. Similar to disjunction elimination, it is possible to construct proofs containing multiple nested applications of proof by contradiction, but such proofs are unnaturally complex. So we fix the depth to 1. To generate an example, we first generate an axiom $\neg g(c)$. Next, for each subproof, we generate a rule $\forall x(f_1(x) \vee \ldots \vee f_n(x) \rightarrow g(x))$, where $n$ is the proof width. The goal is to prove $\neg f_1(c) \wedge \ldots \wedge \neg f_n(c)$. Note that in addition to proof by contradiction, this proof requires disjunction introduction, implication elimination, and conjunction introduction.

Note that the above list constitutes a complete set of deduction rules from propositional natural deduction, save for one rule: implication introduction. However, it is unclear how to construct an example with this deduction rule where its difficulty can be controlled by increasing the width or depth of the proof (e.g. how can a statement of the form $A_1 \rightarrow A_2 \rightarrow \ldots \rightarrow A_n$ be expressed in natural language?).

## A.3   Generating compositional proofs

We use a simple recursive procedure to generate compositional proofs: (1) select a deduction rule uniformly at random, (2) select the premises for the selected rule, (3) recursively generate a subproof for each premise. A consistency checking step is required to make sure we avoid generating contradictory axioms.[6] In addition, we avoid generating an elimination rule directly following an introduction rule (or vice versa).[7] See Algorithm 1 in the appendix for pseudocode of this procedure. To test compositional proofs of various sizes, we implement a parameter that controls the minimum depth of the proof tree, and another parameter that controls the number of distinct rule types in each proof.

## A.4   Further details on evaluation of CoT

We aim to test whether LLMs are able to use deduction rules OOD, where the rules do not appear in the in-context examples, and we take care not to be overly strict. For example, we wish to avoid penalizing the model for formatting differences, so long as the reasoning is correct. To this end, in determining whether a logical form follows from previous logical forms, we consider any deduction rule listed in Table 2. We also allow for two additional rules: (1) given $\forall x(f(x) \rightarrow g(x))$ and $\forall x(g(x) \rightarrow h(x))$ conclude $\forall x(f(x) \rightarrow h(x))$, and (2) given $\forall x(f(x) \rightarrow g(x))$ and $\neg g(c)$ conclude $\neg f(c)$ (i.e. modus tollens).[8] Additionally, we are flexible with respect to the ordering of conjuncts and disjuncts. For example, given the previous steps $f(a) \wedge g(a)$ and $\forall x(g(x) \wedge f(x) \rightarrow u(x) \vee v(x))$, we consider $v(a) \vee u(a)$ to be valid.

## A.5   Generating distractors

***Implication elimination***   For any rule $\forall x(f(x) \rightarrow g(x))$ in the gold proof, we generate a distractor rule $\forall x(f(x) \rightarrow h(x))$ where the concept $h$ is a distractor and is not helpful in completing the proof. In addition, for any ground logical form in the gold proof $f(c)$, we generate a distractor logical form $h(c)$ as well as a rule $\forall x(h(x) \rightarrow h'(x))$. Note that the original PRONTOQA only adds a single distractor, whereas we add multiple, one for each hop in the proof.

***Conjunction introduction***   Similar to those in implication elimination. For any rule $\forall x(f_1(x) \wedge \ldots \wedge f_n(x) \rightarrow g(x))$, we generate a rule of the form $\forall x(h_1(x) \wedge \ldots \wedge h_{n-1}(x) \wedge f_n(x) \rightarrow g(x))$ where

---

[6]An example is: Suppose we select conjunction introduction as the first rule; next, we recursively generate the proof of each conjunct; suppose for each of these, we choose to generate the axioms `cat(alex)` and `¬cat(alex)`.

[7]In the following example, a conjunction introduction step immediately follows a conjunction elimination step: "Jay is a cat and orange. Jay is a cat. Jay is orange. Jay is a cat and orange."

[8]Analogous to the *broadly-valid* steps in PRONTOQA.

---

**Algorithm 1:** Pseudocode for generating examples of compositional proofs in PRONTOQA-OOD. In this algorithm, $\Omega$ denotes the set of all logical forms. The function `generate_compositional_proof` is initially called with parameters $\Omega, \varnothing, d, e$, and `false`, where $d$ is the requested depth and $e$ is a randomly selected entity name (e.g. `alex`, `fae`, etc). `sample` is a helper function that, given an input set of logical forms $S$ and an entity $e$, returns `sample_uniform`({set of logical forms in $S$ with minimal depth where all atoms are of the form $t(e)$ where $t$ is a predicate}).

---

1 | **function** `generate_compositional_proof` (*set of possible conclusions (logical forms) $C$,*
                                           *disallowed deduction rules $R$,*
                                           *requested depth d,*
                                           *ground entity e,*
                                           *is proof hypothetical h*)
2 | initialize $A$ as the set of all deduction rules excluding those in $R$
    /* filter deduction rules such that: (1) an element of $C$ can be a conclusion of the rule, (2) for which we have sufficient depth, and (3) we don't create overly complex logical forms */
3 | **if** *$C$ does not contain a conjunction*
4 |    set $A = A \setminus \{$`conjunction_introduction`$\}$
5 | **if** *$C$ does not contain a disjunction*
6 |    set $A = A \setminus \{$`disjunction_introduction`$\}$
7 | **if** *$h = $ **true** or $d = 1$ or $C$ does not contain a negation*
8 |    set $A = A \setminus \{$`proof_by_contradiction`$\}$
9 | **if** *$h = $ **true** or $d = 1$ or $C$ contains only conjunctions or only disjunctions*
10 |    set $A = A \setminus \{$`disjunction_elimination`$\}$
11 | **if** *$C$ contains only conjunctions or only disjunctions*
12 |    set $A = A \setminus \{$`conjunction_elimination`$\}$
13 | **if** *$C$ contains only conjunctions or only disjunctions and any operand is negated*
14 |    set $A = A \setminus \{$`implication_elimination`$\}$
15 | **if** *$d = 0$ or $C$ contains a singleton logical form or $A = \varnothing$*
16 |    **return** *axiom step with conclusion given by* $sample(C, e)$
17 | $r = $ `sample_uniform`$(A)$
18 | **if** *$r = $ implication_elimination*
19 |    **do**
20 |       for any $c \in C$, $a$ and $c$ share any operands or negations of operands
21 |    **while** $a = $ `generate_compositional_proof`$(\Omega, \varnothing, d-1, e, h)$
22 |    **do**
23 |       $a$ and $s$ do not share any operands or negations of operands
24 |    **while** $s = sample(C, e)$
25 |    **return** *implication_elimination with premises $a$ and* $\forall x (a[e \to x] \to s[e \to x])$
26 | **else if** *$r = $ conjunction_introduction*
27 |    initialize $P$ as an empty list, and $i = 0$
28 |    $L = |C|$ if $C$ contains only conjunctions, else $L = 3$
29 |    **do**
30 |       let $C_i = i^{th}$ operand of $C$ if $C$ contains only conjunctions, else $C_i = \Omega$
31 |       $a = $ `generate_compositional_proof`$(C_i, \{$`conjunction_elimination`$\}, d-1, e, h)$
32 |       **if** *$a$ is atomic and $a$ is not any other operand of $C$*
33 |          append $a$ to $P$
34 |          $i = i + 1$
35 |    **while** $i < L$
36 |    **return** *conjunction_introduction with premises $P$*
37 | **else if** *$r = $ conjunction_elimination*
38 |    let $C'$ be the set of conjunctions of length 3, $i = $ `sample_uniform`$(\{1, 2, 3\})$
39 |    $C' = \{c \in C' : $ the $i^{th}$ operand of $c'$ is in $C\}$
40 |    **do**
41 |       $a = $ `generate_compositional_proof`$(C', \{$`conjunction_introduction`$\}, d-1, e, h)$
42 |    **while** *$a$ has no duplicate operands, and each operand of $a$ is not itself a conjunction or disjunction*
43 |    **return** *conjunction_elimination with premise $a$ and conclusion given by the $i^{th}$ operand of $a$*

---

$h_i$ are distractor concepts. Grounded distractor conjuncts are also generated as axioms $h_i(c)$, so that, given $f_n(c)$, both the gold rule and distractor rule are valid proof steps.

***Conjunction elimination*** Distractors are generated similarly to the conjunction introduction case.

16

**Algorithm 1:** (continued from previous page)

```
44    else if r = disjunction_introduction
45        if C = Ω  let C be the set of disjunctions of length 3
46        i = sample_uniform( number of disjuncts in C)
47        do
48          │ let C_i = i^{th} operand of C
49          │ a = generate_compositional_proof(C_i, {disjunction_elimination}, d − 1, e, h)
50        while a is atomic and a is not any other operand of C
51        replace i^{th} operand of C with a
52        do
53          │ x = sample(C, e)
54        while i^{th} disjunct of x is distinct from all other disjuncts
55        return disjunction_introduction with premise given by the i^{th} operand of x and conclusion x
56    else if r = disjunction_elimination
57        initialize P as an empty list
58        while |P| < 2 do
59          │ p = generate_compositional_proof(C, {disjunction_introduction}, d − 1, e, true)
60          │ if p is not a conjunction or disjunction and p has an axiom that is not an axiom of any q ∈ P
61          │   └ append p to P
62        let A_i be the set of axioms of P_i that are not axioms of P_j for i ≠ j
63        let a_i = sample_uniform(A_i) for all i
64        let a′ be a disjunction with disjuncts a_i
65        a = generate_compositional_proof({a′}, {disjunction_introduction}, d − 1, e, h)
66        └ return disjunction_introduction with premises a and P_i
67    else if r = proof_by_contradiction
68        let N be the set of all negated logical forms
69        a = generate_compositional_proof(N, {proof_by_contradiction}, d − 1, h)
70        do
71          │ let a = ¬s
72          │ b = generate_compositional_proof({s}, {proof_by_contradiction}, d − 1, e, true)
73        while b has an atomic non-negated axiom that is not an axiom of a
74        s′ = sample_uniform({atomic non-negated axioms of b that are not axioms of a})
75        └ return proof_by_contradiction with premises a and b and conclusion ¬s′
```

**Disjunction introduction** Distractors are generated similarly to the conjunction introduction case.

**Disjunction elimination** Since this deduction step has many premises, multiple distractors are necessary to ensure the model doesn't resort to heuristics. For every rule of the form $\forall x(f_i(x) \rightarrow g(x))$, two distractor rules are generated: $\forall x(f_i(x) \rightarrow h'(x))$ and $\forall x(h_i(x) \rightarrow g(x))$. A distractor disjunction is also generated: $h''(c) \vee h_1(c) \vee \ldots \vee h_{n-1}(c)$.

**Proof by contradiction** As with disjunction elimination, multiple distractors are necessary here. We generate two distractor rules $\forall x(f_1(x) \vee \ldots \vee f_n(x) \rightarrow h(x))$ and $\forall x(h_1(x) \vee \ldots \vee h_n(x) \rightarrow g(x))$. We also generate the distractor axiom $\neg h'(c)$ so that the model is forced to choose between two axioms for the first step of the proof.

To avoid creating inconsistencies when generating a distractor rule, we avoid using existing predicates in the consequent of each rule.

**Algorithm 2:** Pseudocode for evaluating the output chain-of-thought. Here, the comparison operations between logical forms ignore the order of conjuncts if both operands are conjunctions; and similarly for disjunctions. In addition, when iterating over previous steps in the proof, we consider them in reverse order, so that more recent steps are prioritized. The helper function `negate` is defined, in order of precedence: $\texttt{negate}(\neg A) = A$, $\texttt{negate}(A \lor B) = \texttt{negate}(A) \land \texttt{negate}(B)$, $\texttt{negate}(A \land B) = \texttt{negate}(A) \lor \texttt{negate}(B)$, or $\texttt{negate}(A) = \neg A$.

```
1  function evaluate_cot(context sentences Q₁,...,Qₘ,
                         predicted chain-of-thought sentences C₁,...,Cₙ,
                         goal sentence g)
2      Lᵍ = semantic_parse(g)                               /* parse the goal */
3      for i ∈ 1,...,m do                                   /* parse the context */
4          Lᵢᵠ = semantic_parse(Qᵢ)
5      for i ∈ 1,...,m do                           /* parse the predicted chain-of-thought */
6          Lᵢᶜ = semantic_parse(Cᵢ)
7      initialize S as an empty set, and H as an empty map
8      for i ∈ 1,...,n do                 /* reconstruct the proof from the chain-of-thought */
9          if Lᵢᶜ indicates 'this is a contradiction'
10             if negate(Lᵢ₊₁ᶜ) ∈ H(Lᵢ₋₁ᶜ)
11                 (P, D, k) = ({Lᵢ₋₁ᶜ, negate(Lᵢ₊₁ᶜ)}, {negate(Lᵢ₊₁ᶜ)}, 1)
12             else continue
13          else
14             (P, D, k) = is_provable(Lᵢᶜ, {L₁ᵠ,...,Lₘᵠ}, S, H)
15          set H(Lᵢᶜ) = ⋃ₚ∈ₚ H(p) \ D
16          if k ≥ 0
17             add Lᵢᶜ to S
18      return Lᵍ ∈ S            /* the proof is correct if the final conclusion is provable */

19 function is_provable(logical form φ,
                        set of axioms A,
                        previous conclusions S,
                        hypothesis map H)
20     if φ ∈ A
21         return ({φ}, 1)                                /* this is an axiom */
22     else if φ is a conjunction or disjunction
23         initialize P' as an empty list, and k' = 0
24         for φᵢ operand in φ do
25             (P, k) = is_provable(φᵢ, A, S, H)
26             if φ is a conjunction
27                 if k ≥ 0 and the step immediately preceding φ in the proof is in P
28                     append P to P'
29                     set k' = k' + k
30                 else break
31             else if k > 0 and φ is a disjunction
32                 return (P, ∅, k + 1)          /* provable by disjunction introduction */

33         if P' has the same size as φ has operands
34             return (⋃ P', ∅, k')             /* provable by conjunction introduction */
35     for a ∈ S ∪ A do
36         if a is a conjunction and φ = aᵢ for some i
37             return ({a}, ∅, 1 + 𝟙{a ∈ A})        /* provable by conjunction elimination */
38         else if a has form ∀x(ψ → γ) where γ[x ↦ c] = φ
39             (P, k) = is_provable(ψ[x ↦ c], A, S, H)
40             if k ≥ 0 and the step immediately preceding φ in the proof is in P ∪ {a}
41                 return (P ∪ {a}, ∅, k + 𝟙{a ∈ A})   /* provable by conjunction elimination */

42     for s ∈ S where s is a disjunction do
43         if for all disjuncts sᵢ, there is a sⱼ ∈ S such that sⱼ = φ and sᵢ ∈ H(sⱼ)
44             return ({sⱼ}, {sᵢ}, 1)              /* provable by disjunction elimination */
```

**Algorithm 2:** (continued from previous page)

```
45   for a ∈ S ∪ A do
46   |   if a has form ∀x(ψ → γ) where γ[x ↦ c] = φ
47   |   |   (P, k) = is_provable(ψ[x ↦ c], A, S, H)
48   |   |   if k ≥ 0 and the step immediately preceding φ in the proof is in P ∪ {a}
49   |   |   └   return (P ∪ {a}, ∅, k + 𝟙{a ∈ A})      /* provable by implication elimination */
     |
50   |   else if a has form ∀x(ψ → γ) where negate(ψ[x ↦ c]) = φ
51   |   |   (P, k) = is_provable(negate(γ[x ↦ c]), A, S, H)
52   |   |   if k ≥ 0 and the step immediately preceding φ in the proof is in P ∪ {a}
     |   |   |   /* provable with additional deduction rules (modus tollens) */
53   |   |   └   return (P ∪ {a}, ∅, k + 𝟙{a ∈ A})
     |
54   if φ ∈ S
55   |   return ({φ}, ∅, 0)                              /* proved by previous step */
56   else if φ has form ∀x(ψ → γ)
     |   /* note:  we precompute this graph */
57   |   let G be the graph where for any axiom in A with form ∀x(α → β), α and β are vertices and there is a
     |     directed edge from α to β
58   |   if there is a path in G from ψ to γ
     |   |   /* provable with additional deduction rules */
59   |   └   return ( axioms corresponding to path edges , ∅, length of path )
     |
60   return (∅, ∅, −1)                     /* this step is not provable (i.e., invalid) */
```