SUPPLEMENTARY MATERIAL FOR ICLR SUBMISSION NUMBER: 5400

CONTENTS

1 DATASET

In this section, we provide examples from the Planning dataset for each of the considered domain - **bw**, **hn**, **gr**, and **dl**. Figure 1 captures the different problem instances.

2 PLANNING VS PLANSFORMER INPUT

We have talked about how a Plansformer brings about reduced knowledge engineering effort. In Figure 2a and Figure 2b, we show the input requirement for an Automated Planner for a **driverlog** problem configuration and Figure 3 shows corresponding input required by Plansformer for the same problem. An automated planner requires two files - (a) **domain.pddl**, and (b) **problem.pddl**. We reduce the knowledge engineering efforts in Plansformer by not requiring:

- explicit mention of *predicates* which are present in **domain.pddl** file.
- explicit mention of *objects* which are present in **problem.pddl** file.

We also have a conversion mechanism for Plansformer and Planning inputs, i.e., given a **domain.pddl** and **problem.pddl** files, we can convert them automatically to input required by Plansformer and vice versa.

3 TRAINING PHASE

In this section, we describe the hardware used for computation, training parameters and time taken by different models for training.

3.1 HARDWARE

We have used 9 (Dual P-100) 44 (Dual V100) GPU nodes for running our experiments. For training all models, we have made use of 24 cores of CPU run on 1 GPU node. Compute and GPU nodes

| Task | Problem | Plan |
|------|---------|------|
| blocksworld | **\<GOAL\>** on b1 b2, ontable b2, on b3 b1, on b4 b5, clear b4, on b5 b3<br>**\<INIT\>** handempty, ontable b1, clear b1, ontable b2, clear b2, on b3 b5, clear b3, ontable b4, on b5 b4<br>**\<ACTION\>** pick-up<br> **\<PRE\>** clear x, ontable x, handempty<br> **\<EFFECT\>** not ontable x, not clear x, not handempty, holding x<br>**\<ACTION\>** put-down<br> **\<PRE\>** holding x<br> **\<EFFECT\>** not holding x, clear x, handempty, ontable x<br>**\<ACTION\>** stack<br> **\<PRE\>** holding x, clear y<br> **\<EFFECT\>** not holding x, not clear y, clear x, handempty, on x y<br>**\<ACTION\>** unstack<br> **\<PRE\>** on x y, clear x, handempty<br> **\<EFFECT\>** holding x, clear y, not clear x, not handempty, not on x y | unstack b4 b2, put-down b4, pick-up b1, stack b1 b2, pick-up b4, stack b4 b1 |

(a) Example from **blocksworld** dataset

| Task | Problem | Plan |
|------|---------|------|
| hanoi | **\<GOAL\>** on d1 d2, clear d1, on d2 d4, on d3 peg2 , clear d3, on d4 peg1 , on d5 peg3 , clear d5<br>**\<INIT\>** smaller peg1 d1, smaller peg1 d2, smaller peg1 d3, smaller peg1 d4, smaller peg1 d5, smaller peg2 d1, smaller peg2 d2, smaller peg2 d3, smaller peg2 d4, smaller peg2 d5, smaller peg3 d1, smaller peg3 d2, smaller peg3 d3, smaller peg3 d4, smaller peg3 d5, smaller d2 d1, smaller d3 d1, smaller d4 d1, smaller d5 d1, smaller d3 d2, smaller d4 d2, smaller d5 d2, smaller d4 d3, smaller d5 d3, smaller d5 d4, on d1 d2, clear d1, on d2 d5, on d3 peg1 , clear d3, on d4 peg2 , clear d4, on d5 peg3<br>**\<ACTION\>** move<br> **\<PRE\>** smaller to disc, on disc from, clear disc, clear to<br> **\<EFFECT\>** clear from, on disc to, not on disc from, not clear to | move d1 d2 d3, move d2 d5 d4, move d1 d3 d2, move d3 peg1 d5, move d1 d2 peg1, move d2 d4 d3, move d1 peg1 d2, move d4 peg2 peg1, move d1 d2 peg2, move d2 d3 d4, move d1 peg2 d2, move d3 d5 peg2 |

(b) Example from **hanoi** dataset

| Task | Problem | Plan |
|------|---------|------|
| grippers | **\<GOAL\>** at ball1 room3, at ball2 room2, at ball3 room3, at ball4 room2, at ball5 room3<br>**\<INIT\>** at-robby robot1 room2, free robot1 lgripper1, free robot1 rgripper1, at-robby robot2 room1, free robot2 lgripper2, free robot2 rgripper2, at ball1 room3, at ball2 room1, at ball3 room1, at ball4 room1, at ball5 room3<br>**\<ACTION\>** move<br> **\<PRE\>** at-robby r from<br> **\<EFFECT\>** at-robby r to, not at-robby r from<br>**\<ACTION\>** pick<br> **\<PRE\>** at obj room, at-robby r room, free r g<br> **\<EFFECT\>** carry r obj g, not at obj room, not free r g<br>**\<ACTION\>** drop<br> **\<PRE\>** carry r obj g, at-robby r room<br> **\<EFFECT\>** at obj room, free r g, not carry r obj g | pick robot2 ball2 room1 lgripper2, move robot1 room2 room1, pick robot1 ball3 room1 lgripper1, move robot1 room1 room3, drop robot1 ball3 room3 lgripper1, pick robot2 ball4 room1 rgripper2, move robot2 room1 room2, drop robot2 ball2 room2 lgripper2, drop robot2 ball4 room2 rgripper2 |

(c) Example from **grippers** dataset

| Task | Problem | Plan |
|------|---------|------|
| driverlog | **\<GOAL\>** at package1 s2, at package2 s2, at package3 s3, at package4 s1, at package5 s3<br>**\<INIT\>** at driver1 s2, at driver2 s4, at truck1 s4, empty truck1, at truck2 s3, empty truck2, at truck3 s3, empty truck3, link s1 s2, link s2 s1, link s1 s3, link s3 s1, link s2 s4, link s4 s2, link s3 s4, link s4 s1, link s1 s4, at package1 s2, at package2 s3, at package3 s1, at package4 s2, at package5 s4<br>**\<ACTION\>** load-truck<br>**\<PRE\>** at truck loc, at obj loc<br>**\<EFFECT\>** not at obj loc, in obj truck<br>**\<ACTION\>** unload-truck<br>**\<PRE\>** at truck loc, in obj truck<br>**\<EFFECT\>** not in obj truck, at obj loc<br>**\<ACTION\>** board-truck<br>**\<PRE\>** at truck loc, at driver loc, empty truck<br>**\<EFFECT\>** not at driver loc, driving driver truck, not empty truck<br>**\<ACTION\>** disembark-truck<br>**\<PRE\>** at truck loc, driving driver truck<br>**\<EFFECT\>** not driving driver truck, at driver loc, empty truck<br>**\<ACTION\>** drive-truck<br>**\<PRE\>** at truck loc-from, driving driver truck, link loc-from loc-to<br>**\<EFFECT\>** not at truck loc-from, at truck loc-to<br>**\<ACTION\>** walk<br>**\<PRE\>** at driver loc-from, path loc-from loc-to<br>**\<EFFECT\>** not at driver loc-from, at driver loc-to | board-truck driver2 truck1 s4, load-truck package5 truck1 s4, drive-truck truck1 s4 s1 driver2, load-truck package3 truck1 s1, drive-truck truck1 s1 s3 driver2, unload-truck package5 truck1 s3, unload-truck package3 truck1 s3, load-truck package2 truck1 s3, drive-truck truck1 s3 s1 driver2, drive-truck truck1 s1 s2 driver2, load-truck package4 truck1 s2, unload-truck package2 truck1 s2, drive-truck truck1 s2 s1 driver2, unload-truck package4 truck1 s1 |

(d) Example from **driverlog** dataset

Figure 1: Problem instances from four different planning domains

```pddl
1  (define (domain driverlog)
2    (:requirements :typing)
3    (:types          location locatable - object
4       driver truck obj - locatable
5
6    )
7    (:predicates
8       (at ?obj - locatable ?loc - location)
9       (in ?obj1 - obj ?obj - truck)
10      (driving ?d - driver ?v - truck)
11      (link ?x ?y - location) (path ?x ?y
    - location)
12      (empty ?v - truck)
13  )
14
15
16  (:action LOAD-TRUCK
17    :parameters
18     (?obj - obj
19      ?truck - truck
20      ?loc - location)
21    :precondition
22     (and (at ?truck ?loc) (at ?obj ?loc))
23    :effect
24     (and (not (at ?obj ?loc)) (in ?obj ?truck
    )))
25
26  (:action UNLOAD-TRUCK
27    :parameters
28     (?obj - obj
29      ?truck - truck
30      ?loc - location)
31    :precondition
32     (and (at ?truck ?loc) (in ?obj ?truck))
33    :effect
34     (and (not (in ?obj ?truck)) (at ?obj ?loc
    )))
35
36  (:action BOARD-TRUCK
37    :parameters
38     (?driver - driver
39      ?truck - truck
40      ?loc - location)
41    :precondition
42     (and (at ?truck ?loc) (at ?driver ?loc
    ) (empty ?truck))
43    :effect
44     (and (not (at ?driver ?loc)) (driving
    ?driver ?truck) (not (empty ?truck))))
45
46  (:action DISEMBARK-TRUCK
47    :parameters
48     (?driver - driver
49      ?truck - truck
50      ?loc - location)
51    :precondition
52     (and (at ?truck ?loc) (driving ?driver
    ?truck))
53    :effect
54     (and (not (driving ?driver ?truck)) (at
    ?driver ?loc) (empty ?truck)))
55
56  (:action DRIVE-TRUCK
57    :parameters
58     (?truck - truck
59      ?loc-from - location
60      ?loc-to - location
61      ?driver - driver)
62    :precondition
63     (and (at ?truck ?loc-from)
64      (driving ?driver ?truck) (link ?loc-from
    ?loc-to))
65    :effect
66     (and (not (at ?truck ?loc-from)) (at ?truck
    ?loc-to)))
67
68  (:action WALK
69    :parameters
70     (?driver - driver
71      ?loc-from - location
72      ?loc-to - location)
73    :precondition
74     (and (at ?driver ?loc-from) (path ?loc
    -from ?loc-to))
75    :effect
76     (and (not (at ?driver ?loc-from)) (at
    ?driver ?loc-to)))
77
78  )
```

(a) Capturing **driverlog's** environment in `domain.pddl`



```pddl
1  (define (problem problem_3_2_4_34291)
2  (:domain driverlog)
3  (:objects
4   driver1 driver2 driver3 - driver
5   truck1 truck2 - truck
6   package1 package2 package3 package4 - obj
7   s1 s2 s3 - location
8   )
9
10 (:init
11  (at driver1 s3)
12  (at driver2 s3)
13  (at driver3 s3)
14  (at truck1 s3)
15  (empty truck1)
16  (at truck2 s3)
17  (empty truck2)
18  (link s1 s2)
19  (link s2 s1)
20  (link s2 s3)
21  (link s3 s2)
22  (link s3 s1)
23  (link s1 s3)
24  (at package1 s3)
25  (at package2 s3)
26  (at package3 s2)
27  (at package4 s1)
28  )
29
30 (:goal (and
31  (at package1 s1)
32  (at package2 s3)
33  (at package3 s3)
34  (at package4 s3)
35  )
36  )
37  )
```

(b) Capturing the current state and desired goal of an object in **driverlog's** environment in `problem.pddl`

Figure 2: Files required to model a problem from **driverlog** in PDDL for execution by an Automated Planner

```
<GOAL> at package1 s1, at package2 s3, at package3 s3, at package4 s3
<INIT> at driver1 s3, at driver2 s3, at driver3 s3, at truck1 s3, empty truck1, at truck2 s3, empty truck2,
link s1 s2, link s2 s1, link s2 s3, link s3 s2, link s3 s1, link s1 s3, at package1 s3, at package2 s3, at
package3 s2, at package4 s1
<ACTION> load-truck
    <PRE> at truck loc, at obj loc
    <EFFECT> not at obj loc, in obj truck
<ACTION> unload-truck
    <PRE> at truck loc, in obj truck
    <EFFECT> not in obj truck, at obj loc
<ACTION> board-truck
    <PRE> at truck loc, at driver loc, empty truck
    <EFFECT> not at driver loc, driving driver truck, not empty truck
<ACTION> disembark-truck
    <PRE> at truck loc, driving driver truck
    <EFFECT> not driving driver truck, at driver loc, empty truck
<ACTION> drive-truck
    <PRE> at truck loc-from, driving driver truck, link loc-from loc-to
    <EFFECT> not at truck loc-from, at truck loc-to
<ACTION> walk
    <PRE> at driver loc-from, path loc-from loc-to
    <EFFECT> not at driver loc-from, at driver loc-to
```

Figure 3: Plansformer's input for the same problem defined in Figure 2a

| Hyperparameter | Value |
|---|---|
| Train Batch Size | 8 |
| Validation Batch Size | 8 |
| Train Epochs | 3 |
| Validation Epochs | 1 |
| Learning Rate | 1e-4 |
| Max Source Text Length | 512 |
| Max Target Text Length | 150 |

Table 1: Hyperparameters used for Training

have 128 GB of RAM and Big Data nodes have 1.5 TB RAM. All nodes have EDR infiniband (100 Gb/s) interconnects, and access to 1.4 PB of GPFS storage. The processor speed is 2.8 GHz.

## 3.2 TRAINING HYPERPARAMETERS

Table 1 captures the hyperparameters set for training our models. For plan generation by all models apart from Codex, we have used beam search with number of beams set to 2, repetition penalty of 2.5, and length penalty set to 1.0. Codex doesn't have the functionality to change the parameters, thus, we have used it in the default setting. On the parameters constituting the considered models, Codex has 12 billion parameters, GPT-2 has 1.2 billion parameters, T5-base has 220 million parameters, and CodeT5-base has 8.35 million parameters.

## 3.3 TRAINING TIME

Figure 4 presents the training time taken by different Plansformer variations. Base models are Plansformer variants directly trained on each of the planning domains with CodeT5 as base. Derived models use a Plansformer base model as a starting point, and further pretrain on other domains. We can see a considerable drop in training time taken by derived models. It is also to be noted that these derived models outperform the base models when entire training data points are used.

4

| | Models | Training Time |
|---|---|---|
| | plansformer-bw | 1hr 1min |
| | plansformer-hn | 1hr 10mins |
| | plansformer-gr | 1hr 12mins |
| *Base Models* | plansformer-dl | 1hr 15mins |
| | plansformer-bw-dl[500] | 3 mins |
| | plansformer-bw-dl[1000] | 5 mins |
| | plansformer-bw-dl[1500] | 7 mins |
| | plansformer-bw-dl[2000] | 13 mins |
| | plansformer-bw-dl[5000] | 21 mins |
| | plansformer-bw-dl[10000] | 38 mins |
| | plansformer-bw-dl[14400] | 46 mins |
| | plansformer-bw-gr[500] | 3 mins |
| | plansformer-bw-gr[1000] | 5 mins |
| | plansformer-bw-gr[1500] | 7 mins |
| | plansformer-bw-gr[2000] | 13 mins |
| | plansformer-bw-gr[5000] | 19 mins |
| | plansformer-bw-gr[10000] | 37 mins |
| | plansformer-bw-gr[14400] | 44 mins |
| | plansformer-bw-hn[500] | 2 mins |
| | plansformer-bw-hn[1000] | 3 mins |
| | plansformer-bw-hn[1500] | 5 mins |
| | plansformer-bw-hn[2000] | 9 mins |
| | plansformer-bw-hn[5000] | 18 mins |
| | plansformer-bw-hn[10000] | 23 mins |
| | plansformer-bw-hn[14400] | 32 mins |
| | plansformer-hn-bw[500] | 2 mins |
| | plansformer-hn-bw[1000] | 2 mins |
| | plansformer-hn-bw[1500] | 3 mins |
| | plansformer-hn-bw[2000] | 5 mins |
| | plansformer-hn-bw[5000] | 8 mins |
| | plansformer-hn-bw[10000] | 12 mins |
| | plansformer-hn-bw[14400] | 17 mins |
| | plansformer-hn-gr[500] | 2 mins |
| | plansformer-hn-gr[1000] | 3 mins |
| | plansformer-hn-gr[1500] | 5 mins |
| | plansformer-hn-gr[2000] | 9 mins |
| | plansformer-hn-gr[5000] | 18 mins |
| | plansformer-hn-gr[10000] | 23 mins |
| | plansformer-hn-gr[14400] | 29 mins |
| | plansformer-hn-dl[500] | 1 min |
| | plansformer-hn-dl[1000] | 5 mins |
| | plansformer-hn-dl[1500] | 5 mins |
| | plansformer-hn-dl[2000] | 12 mins |
| | plansformer-hn-dl[5000] | 19 mins |
| | plansformer-hn-dl[10000] | 33 mins |
| | plansformer-hn-dl[14400] | 39 mins |
| | plansformer-gr-bw[500] | 2 mins |
| | plansformer-gr-bw[1000] | 2 mins |
| | plansformer-gr-bw[1500] | 3 mins |
| | plansformer-gr-bw[2000] | 5 mins |
| | plansformer-gr-bw[5000] | 8 mins |
| | plansformer-gr-bw[10000] | 14 mins |
| | plansformer-gr-bw[14400] | 18 mins |
| | plansformer-gr-hn[500] | 1 min |
| | plansformer-gr-hn[1000] | 2 mins |
| | plansformer-gr-hn[1500] | 2 mins |
| | plansformer-gr-hn[2000] | 5 mins |
| | plansformer-gr-hn[5000] | 10 mins |
| | plansformer-gr-hn[10000] | 13 mins |
| | plansformer-gr-hn[14400] | 16 mins |
| | plansformer-gr-dl[500] | 3 mins |
| | plansformer-gr-dl[1000] | 5 mins |
| | plansformer-gr-dl[1500] | 7 mins |
| | plansformer-gr-dl[2000] | 13 mins |
| | plansformer-gr-dl[5000] | 21 mins |
| | plansformer-gr-dl[10000] | 38 mins |
| | plansformer-gr-dl[14400] | 46 mins |
| | plansformer-dl-bw[500] | 1 min |
| | plansformer-dl-bw[1000] | 1 min |
| | plansformer-dl-bw[1500] | 2 mins |
| | plansformer-dl-bw[2000] | 2 mins |
| | plansformer-dl-bw[5000] | 4 mins |
| | plansformer-dl-bw[10000] | 7 mins |
| | plansformer-dl-bw[14400] | 13 mins |
| | plansformer-dl-hn[500] | 1 min |
| | plansformer-dl-hn[1000] | 2 mins |
| | plansformer-dl-hn[1500] | 2 mins |
| | plansformer-dl-hn[2000] | 5 mins |
| | plansformer-dl-hn[5000] | 8 mins |
| | plansformer-dl-hn[10000] | 12 mins |
| | plansformer-dl-hn[14400] | 16 mins |
| | plansformer-dl-gr[500] | 2 mins |
| | plansformer-dl-gr[1000] | 5 mins |
| | plansformer-dl-gr[1500] | 7 mins |
| | plansformer-dl-gr[2000] | 12 mins |
| | plansformer-dl-gr[5000] | 21 mins |
| | plansformer-dl-gr[10000] | 29 mins |
| *Derived Models* | plansformer-dl-gr[14400] | 37 mins |

Figure 4: Training time of different Plansformer variants

# 4 EXTENDED RESULTS

This section adds additional qualitative and quantitative results that cover all the domains and model configurations tried and tested during our experimentation with Plansformer. In our initial testing phase, we have fine-tuned both T5 and CodeT5 with the same blocksworld dataset and hyperparameters and found that fine-tuned T5 gave 32% valid plans, whereas fine-tuned CodeT5 generated 90% valid plans. This is because CodeT5 has syntactically meaningful sequences for code-like structured inputs well defined as opposed to T5, which only deals with natural language. With this

intuition that models pre-trained on code have an advantage for plan generation, we proceeded with the choice of using CodeT5 for all our experiments.

## 4.1 QUALITATIVE ANALYSIS

Figure 5 shows the output generations obtained by different models under study for a problem instance from each of the planning domains. When reporting Plansformer results, we take the *base models* corresponding to the domain being tested.

## 4.2 QUANTITATIVE ANALYSIS

Figure 6 captures the performance of all models in their ability to generate plans for multiple domains.

Figures 7 to 15 represent the performance of different base models fine-tuned and tested on other domains in graphical manner. Additionally, we also wanted to check Plansformer's capability in plan generation when the plan length of test set is considerably larger than that of the train set. For this purpose, we have created a train set for blocksworld consisting of 2,3 block configurations and a test set with 4,5 block configurations. The test set consists of 100 total instances, 50 from 4 block configuration and 50 from 5 block configuration. The average plan length for the train set and test set is 4 and 10 respectively. Plansformer trained on 2,3 block configurations was able to generate 64% valid plans on problem instances from the test set, showing that our approach can generate valid plans even if the plan length for test set $>>>$ train set. We were able to achieve 64% valid plans in this experimentation as the train set consists of only 162 data points (all possible 2,3 block configurations) as compared to 87% obtained by Plansformer-bw (trained on 14,400 data points) on the same test set.

### 4.2.1 TRANSFER LEARNING

We have reported the advantages of transfer learning when using Plansformer-bw as the base model and further fine-tuning it on **hn** for all the varying data points. Figures 16 and 17 show a similar trend for the domains **dl** and **gr** on further fine-tuning of the Plansformer-bw base model.

### 4.2.2 MULTI-DOMAIN PLAN GENERATION USING A SINGLE MODEL

We have additionally trained a single Plansformer model on all the training data points belonging to all the four domains - **bw, hn, dl,** and **gr**. Each domain consists of 14,400 training data points, thus, the single Plansformer model is trained on a total of 57,600 data points. The obtained single Plansformer model is tested on the validation datasets corresponding to all the four domains. Each dataset consists of 3000 problem instances. Table 2 reports the performance obtained by the single Plansformer model on all four domains in terms of plan validation. We observe that the single model is able to perform relatively comparable to the *base models*. The single Plansformer trained on all models outperforms Plansformer-bw by 5.79% (in terms of valid plans) and has around a 4.5% to 13.5% decrease in valid plans for the other three domains.

| Test Domain | Valid Plans | Cost-Optimal Plans | Invalid Plans |
|---|---|---|---|
| **bw** | 95.83% | 93.75% | Failed = 4.17%, Incomplete = 0% |
| **hn** | 79.25% | 76.72% | Failed = 2.34%, Incomplete = 18.41% |
| **gr** | 78.44% | 50.61% | Failed = 15.03%, Incomplete = 6.53% |
| **dl** | 63.03% | 55.25% | Failed = 2.81%, Incomplete = 34.16% |

Table 2: Results of TREAT System

## 5 PLANSFORMER ARCHITECTURE

Figure 18 shows the layer-wise architecture that makes up Plansformer. All these layers are updated during the fine-tuning process. We inherit this architecture from CodeT5, and do not freeze any layers during our fine-tuning process involved for constructing Plansformer.

**Legend:** text represents optimal and valid plans, text represents invalid plans, text represents valid but not optimal plans

| Problem Instance | Actual Plan (FastDownward) | GPT-2 | T5 | CodeT5 | Codex | Plansformer |
|---|---|---|---|---|---|---|
| *blocksworld* | | | | | | |
| <GOAL> ontable b1, clear b1, ontable b2, clear b2, ontable b3, ontable b3, clear b3, ontable b4 <INIT> handempty, ontable b1, on b2, b3, clear b3, ontable b3, on b4 b1, clear b4 <ACTION> pick-up <PRE> clear x, ontable x, handempty <EFFECT> not ontable x, not clear x, not handempty, holding x <ACTION> put-down <PRE> holding x <EFFECT> not holding x, clear x, handempty, ontable x <ACTION> stack <PRE> holding x, clear y <EFFECT> not holding x, not clear y, clear x, handempty, on x y <ACTION> unstack <PRE> on x y, clear x, handempty <EFFECT> holding x, clear y, not clear x, not handempty, not on x y | unstack b2 b3, put-down b2, unstack b4 b1, put-down b4 | move - left, right, up, down, right, down, left, right, down | clear b3, clear b4, ontable b5 | ontableontableontableontableontableontable leontableontableontableontableontableontable ableontableontableontableontableontableontableo ntableontableontableontableontableontable eontableontable x, ontable x, ontable x, ontable x, ontable x | unstack b2 b3, put-down b2, unstack b4 b1, put-down b4, unstack b1 b3, stack b1 b2, unstack b3 b2, put-down b3, unstack b1 b2, stack b1 b3, pick-up b4, stack b4 b1, put-down b2, stack b2 b4, pick-up b3, stack b3 b2 | unstack b2 b3, put-down b2, unstack b4 b1, put-down b4 |
| *hanoi* | | | | | | |
| <GOAL> on d1 peg2, clear d1, on d2 peg1, clear d2, clear peg3 <INIT> smaller peg1 d1, smaller peg1 d2, smaller peg2 d1, smaller peg2 d2, smaller peg3 d1, smaller peg3 d2, smaller d2 d1, on d1 d2, clear d1, on d2 peg3, clear peg1, clear peg2 <ACTION> move <PRE> smaller to disc, on disc from, clear disc to <EFFECT> clear from, on disc to, not on disc from, not clear to | move smaller peg1 d1 d2 d3 , move d2 d3 peg1 d4 , move d1 peg1 d2 | 2, move d1 peg2 peg2, move d2 d3 | | | move d1 d2 peg3, move d2 peg1 peg2, move d1 d2 peg3, move peg1 peg2 peg3, move d1 d2 peg2, move peg1 peg3 peg3, move d1 d2 peg2, move peg1 peg3 peg2, ... move d1 peg | move d1 d2 peg2, move d2 peg3 peg1 |
| *grippers* | | | | | | |
| <GOAL> at ball1 room2, at ball2 room1, at ball3 room4 <INIT> at-robby robot1 room2, free robot1 lgripper1, free robot1 rgripper1, at-robby robot2 room1, free robot2 lgripper2, free robot2 rgripper2, at-robby robot3 room2, free robot3 lgripper3, free robot3 rgripper3, at ball1 room3, at ball2 room4, at ball3 room1 <ACTION> move <PRE> at-robby r from <EFFECT> at-robby r to, not at-robby r from <ACTION> pick <PRE> at obj room, at-robby r room, free r g <EFFECT> carry r obj g, not at obj room, not free r g <ACTION> drop <PRE> carry r obj g, at-robby r room <EFFECT> at obj room, free r g, not carry r obj g | pick robot2 ball3 room1 lgripper2, move robot2 room1 room4, pick robot2 ball2 room4 rgripper2, move robot1 room2 room2, move robot1 room1 drop robot2 ball1 room3, drop robot2 ball3 room3 lgripper2, pick robot2 ball1 room1, drop robot1 ball1 room2 lgripper1, move robot1 room3 room3, drop robot1 ball1 room2 lgripper1 | at ball 1 from - rob by r gpo -rob by r gpo gpo at ball 2 from - rob by r gpo goo at ball lgripper2, | 1 lgripper2, drop robot2 ball4 room1 lgripper2, | point:point;point;point:point: point: point: point:point: point: point: point: point: point: point: point: point: point:point: point: point: point: point: point:point: point: point: point: point: point: | pick robot2 ball2 room1 lgripper2, move robot1 room2 room1, pick robot1 ball3 room1 lgripper1, move robot1 room1 room3, drop robot1 ball3 room3 lgripper1, pick robot2 ball4 room1 rgripper2, move robot2 room2, drop robot2 ball4 room2 rgripper2 | pick robot3 ball1 room2 lgripper3, move robot3 room2 room1, pick robot3 ball2 room1 rgripper3, drop robot3 ball1 room1 lgripper3, move robot3 room1 room4, pick robot3 ball3 room4 lgripper3, drop robot3 ball3 room4 lgripper3, drop robot3 ball2 room4 rgripper3, move robot3 room4 rgripper3, move robot3 room4 room3, drop robot3 ball3 room3 lgripper3 |
| *driving* | | | | | | |
| <GOAL> at package1 s1, at package2 s2, at package3 s2, at package4 s1 <INIT> at driver1 s3, at driver2 s1, at truck1 s4, empty truck1, at truck2 s3, empty truck2, at truck3 s1, empty truck3, link s1 s2, link s2 s1, link s1 s3, link s3 s1, link s4 s1, link s1 s4, link s3 s2, link s2 s3, at package1 s1, at package2 s2, at package3 s4, at package4 s1 <ACTION> load-truck <PRE> at truck loc, at obj loc <EFFECT> not at obj loc, in obj truck <ACTION> board-truck <PRE> at truck loc, at driver loc, empty truck <EFFECT> not at driver loc, driving driver truck, not empty truck <ACTION> disembark-truck <PRE> at truck loc, driving driver truck <EFFECT> at driver loc, empty truck <ACTION> walk <PRE> at driver loc-from, path loc-from loc-to <EFFECT> not at driver loc-from, at driver loc-to | board-truck driver1 truck2 s3, drive-truck truck2 s3 s4 driver1, load-truck package3 truck2 s4, drive-truck truck2 s4 s2 driver1, unload-truck package3 truck2 s2 | if load -1 ruck -truck (empty ) i=> not empty truck, at obj loc- from obj loc- to => path- from- to => | at package1 s1, at package2 s2, at package3 s2, at | atatatatatatatatatatatatatatatatatatatatatatata tatatatatatatatatatatatatatatatatatatatatatat atatatat | board-truck driver2 truck1 s1, load-truck package2 truck1 s1, drive-truck truck1 s1 s2 driver2, load-truck package3 truck1 s2, drive-truck truck1 s2 s4 driver2, unload-truck package3 truck1 s1, unload-truck package1 truck1 s1 | board-truck driver2 truck3 s1, drive-truck truck3 s1 s4 driver2, load-truck package3 truck3 s4, drive-truck truck3 s4 s2 driver2, unload-truck package3 truck3 s2 |

Figure 5: Output generations from different models for planning problem instances

**Test Domain** - *(tested using 3600 problems for each domain)*

Trained using 14,400 problem instances for each domain

| Variant | blocksworld (bw) | | | hanoi (hn) | | | gripper (gr) | | | driverlog (dl) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Valid Plans | Invalid Plans | Cost-Optimal Plans | Valid Plans | Invalid Plans | Cost-Optimal Plans | Valid Plans | Invalid Plans | Cost-Optimal Plans | Valid Plans | Invalid Plans | Cost-Optimal Plans |
| plansformer-bw | 90.04% | Failed = 9.94%, Incomplete = 0.02% | 88.44% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-hn | 0.00% | Failed = 2.58%, Incomplete = 97.42% | 0.00% | 84.97% | Failed = 14.72%, Incomplete = 0.31% | 82.58% | 0.00% | Failed = 1.14%, Incomplete = 98.86% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-gr | 0.00% | Failed = 1.75%, Incomplete = 98.25% | 0.00% | 0.00% | Failed = 8.83%, Incomplete = 91.16% | 0.00% | 82.97% | Failed = 16.61%, Incomplete = 0.42% | 69.47% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-dl | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 76.56% | Failed = 23.44%, Incomplete = 0% | 52.61% |
| plansformer-bw-dl[500] | 68.28% | Failed = 31.72%, Incomplete = 0% | 66.56% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 23.17%, Incomplete = 76.83% | 0.00% | 0.00% | Failed = 99.81%, Incomplete = 0.19% | 0.00% |
| plansformer-bw-dl[1000] | 59.22% | Failed = 40.78%, Incomplete = 0% | 58.44% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 1.92% | Failed = 97.97%, Incomplete = 0.11% | 1.58% |
| plansformer-bw-dl[1500] | 61.67% | Failed = 38.33%, Incomplete = 0% | 60.86% | 0.00% | Failed = 0.17%, Incomplete = 99.83% | 0.00% | 0.00% | Failed = 83.94%, Incomplete = 16.06% | 0.00% | 17.57% | Failed = 82.26%, Incomplete = 0.17% | 15.99% |
| plansformer-bw-dl[2000] | 49.57% | Failed = 50.43%, Incomplete = 0% | 48.66% | 0.00% | Failed = 0.17%, Incomplete = 99.83% | 0.00% | 0.00% | Failed = 83.94%, Incomplete = 16.06% | 0.00% | 37.81% | Failed = 61.94%, Incomplete = 0.25% | 27.39% |
| plansformer-bw-dl[5000] | 33.44% | Failed = 56.69%, Incomplete = 9.86% | 32.50% | 0.00% | Failed = 0.14%, Incomplete = 99.86% | 0.00% | 0.00% | Failed = 94.69%, Incomplete = 5.31% | 0.00% | 82.72% | Failed = 17.06%, Incomplete = 0.22% | 64.81% |
| plansformer-bw-dl[10000] | 3.25% | Failed = 10.64%, Incomplete = 86.11% | 3.25% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 78.56% | Failed = 10.22%, Incomplete = 11.22% | 62.33% |
| plansformer-bw-dl[14400] | 0.00% | Failed = 1%, Incomplete = 99% | 0.00% | 0.00% | Failed = 2.56%, Incomplete = 97.44% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 93.75% | Failed = 6%, Incomplete = 0.25% | 74.56% |
| plansformer-bw-gr[500] | 32.33% | Failed = 40.03%, Incomplete = 27.64% | 32.14% | 0.00% | Failed = 4.89%, Incomplete = 95.11% | 0.00% | 2.86% | Failed = 96.64%, Incomplete = 0.5% | 0.97% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-gr[1000] | 17.03% | Failed = 13.81%, Incomplete = 69.17% | 17.00% | 0.00% | Failed = 0.67%, Incomplete = 99.33% | 0.00% | 9.42% | Failed = 89.47%, Incomplete = 1.11% | 7.61% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-gr[1500] | 6.39% | Failed = 10.78%, Incomplete = 82.83% | 6.28% | 0.00% | Failed = 9.72%, Incomplete = 90.28% | 0.00% | 20.53% | Failed = 79.42%, Incomplete = 0.05 | 9.03% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |

Figure 6: Plan Validation metrics for all Plansformer variants

| Model | Col1 | Col2 | Col3 | Col4 | Col5 | Col6 | Col7 | Col8 | Col9 | Col10 | Col11 | Col12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| plansformer-bw-gr[2000] | 9.31% | Failed = 18.42%, Incomplete = 72.28% | 9.31% | 0.00% | Failed = 0.92%, Incomplete = 99.08% | 0.00% | 26.56% | Failed = 73.11%, Incomplete = 0.33% | 11.31% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-gr[5000] | 1.28% | Failed = 3.64%, Incomplete = 95.08% | 1.28% | 0.00% | Failed = 39.03%, Incomplete = 60.97% | 0.00% | 64.22% | Failed = 35.78%, Incomplete = 0% | 38.22% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-gr[10000] | 0.00% | Failed = 1.56%, Incomplete = 98.44% | 0.00% | 0.00% | Failed = 6.06%, Incomplete = 93.94% | 0.00% | 77.22% | Failed = 22.75%, Incomplete = 0.03% | 47.83% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-gr[14400] | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 2.03%, Incomplete = 97.97% | 0.00% | 87.17% | Failed = 12.80%, Incomplete = 0.03% | 60.86% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[500] | 40.86% | Failed = 58.33%, Incomplete = 0.81% | 37.42% | 3.11% | Failed = 96.89%, Incomplete = 0% | 2.92% | 0.00% | Failed = 4.92%, Incomplete = 95.08% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[1000] | 37.81% | Failed = 57.81%, Incomplete = 4.39% | 35.61% | 14.72% | Failed = 85.28%, Incomplete = 0% | 13.18% | 0.00% | Failed = 4.92%, Incomplete = 95.08% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[1500] | 31.64% | Failed = 60.53%, Incomplete = 7.83% | 27.70% | 35.28% | Failed = 64.44%, Incomplete = 0.28% | 32.58% | 0.00% | Failed = 14.58%, Incomplete = 85.42% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[2000] | 17.97% | Failed = 56.22%, Incomplete = 25.81% | 16.56% | 48.19% | Failed = 51.78%, Incomplete = 0.03% | 44.50% | 0.00% | Failed = 9.06%, Incomplete = 90.94% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[5000] | 1.18% | Failed = 2.74%, Incomplete = 96.08% | 1.01% | 60.25% | Failed = 39.52%, Incomplete = 0.22% | 57.27% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[10000] | 0.00% | Failed = 1.22%, Incomplete = 98.78% | 0.00% | 88.28% | Failed = 11.66%, Incomplete = 0.06% | 86.64% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-bw-hn[14400] | 0.00% | Failed = 1.55%, Incomplete = 98.44% | 0.00% | 97.05% | Failed = 2.94%, Incomplete = 0% | 95.22% | 0.00% | Failed = 1.44%, Incomplete = 98.56% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% |
| plansformer-hn-bw[500] | 11.75% | Failed = 88.14%, Incomplete = 0.11% | 11.06% | 0.00% | Failed = 2.34%, Incomplete = 97.61% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-hn-bw[1000] | 20.61% | Failed = 78.97%, Incomplete = 0.42% | 18.28% | 0.00% | Failed = 0.28%, Incomplete = 99.72% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-hn-bw[1500] | 27.58% | Failed = 70.36%, Incomplete = 2.06% | 22.00% | 0.00% | Failed = 0.14%, Incomplete = 99.86% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |
| plansformer-hn-bw[2000] | 46.53% | Failed = 53.31%, Incomplete = 0.17% | 37.89% | 0.00% | Failed = 0.56%, Incomplete = 99.44% | 0.00% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% |

| Model | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| plansformer-hn-bw[5000] | 69.64% | Failed = 29.56%, Incomplete = 0.81% | 64.72% | Failed = 0.22%, Incomplete = 99.78% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% |
| plansformer-hn-bw[10000] | 90.36% | Failed = 9.58%, Incomplete = 0.06% | 87.97% | Failed = 4.06%, Incomplete = 95.94% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% |
| plansformer-hn-bw[14400] | 95.44% | Failed = 4.56%, Incomplete = 0% | 93.03% | Failed = 2.28%, Incomplete = 97.72% | 0.00% | Failed = 100%, Incomplete = 0% | 0.00% | 0.00% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% |
| plansformer-hn-gr[500] | 0.00% | Failed = 1.56%, Incomplete = 98.44% | 0.00% | Failed = 75.81%, Incomplete = 17.28% | 6.92% | Failed = 95.58%, Incomplete = 1.25% | 3.17% | 1.78% | Failed = 1.56%, Incomplete = 98.44% | 0.00% | 0.00% |
| plansformer-hn-gr[1000] | 0.00% | Failed = 1.06%, Incomplete = 98.94% | 0.00% | Failed = 31.42%, Incomplete = 55.69% | 12.89% | Failed = 89.22%, Incomplete = 0.14% | 10.64% | 5.92% | Failed = 0%, Incomplete = 100% | 0.00% | 0.00% |
| plansformer-hn-gr[1500] | 0.00% | Failed = 0.5%, Incomplete = 99.5% | 0.00% | Failed = 10.31%, Incomplete = 81.61% | 8.08% | Failed = 84.36%, Incomplete = 0.03% | 15.61% | 10.08% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% |
| plansformer-hn-gr[2000] | 0.00% | Failed = 1.97%, Incomplete = 98.03% | 0.00% | Failed = 1.5%, Incomplete = 91.78% | 6.72% | Failed = 77.36%, Incomplete = 0.28% | 22.36% | 16.06% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% |
| plansformer-hn-gr[5000] | 0.00% | Failed = 27.69%, Incomplete = 72.31% | 0.00% | Failed = 1.08%, Incomplete = 96.28% | 2.64% | Failed = 50.33%, Incomplete = 0.36% | 49.31% | 33.72% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% |
| plansformer-hn-gr[10000] | 0.00% | Failed = 41.14%, Incomplete = 58.86% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | Failed = 24.25%, Incomplete = 0.17% | 75.58% | 48.56% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% |
| plansformer-hn-gr[14400] | 0.00% | Failed = 5.28%, Incomplete = 94.72% | 0.00% | Failed = 0.11%, Incomplete = 99.89% | 0.00% | Failed = 34.31%, Incomplete = 0.19% | 65.50% | 39.44% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% |
| plansformer-hn-dl[500] | 0.00% | Failed = 0.97%, Incomplete = 99.03% | 0.00% | Failed = 87.44%, Incomplete = 3.36% | 9.19% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 99.44%, Incomplete = 0.44% | 0.11% | 0.08% |
| plansformer-hn-dl[1000] | 0.00% | Failed = 1.11%, Incomplete = 98.89% | 0.00% | Failed = 84.86%, Incomplete = 3.89% | 11.25% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 92.47%, Incomplete = 0.56% | 6.97% | 4.64% |
| plansformer-hn-dl[1500] | 0.00% | Failed = 1.67%, Incomplete = 98.33% | 0.00% | Failed = 56.89%, Incomplete = 32.92% | 10.19% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 72.72%, Incomplete = 0.17% | 27.11% | 18.19% |
| plansformer-hn-dl[2000] | 0.00% | Failed = 1.92%, Incomplete = 98.08% | 0.00% | Failed = 50.0%, Incomplete = 38.44% | 11.56% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 59.14%, Incomplete = 0.17% | 40.69% | 30.06% |
| plansformer-hn-dl[5000] | 0.00% | Failed = 37.17%, Incomplete = 62.83% | 0.00% | Failed = 0.14%, Incomplete = 99.11% | 0.75% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 26.72%, Incomplete = 0.53% | 72.75% | 55.83% |

| Model | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| plansformer-hn-dll[10000] | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 0.03% | Failed = 0.11%, Incomplete = 99.86% | 0.03% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 86.00% | Failed = 13.92%, Incomplete = 0.08% | 67.89% |
| plansformer-hn-dll[14400] | 0.00% | Failed = 0.31%, Incomplete = 99.69% | 0.00% | 0.00% | Failed = 0.92%, Incomplete = 99.08% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 90.06% | Failed = 9.64%, Incomplete = 0.31% | 72.61% |
| plansformer-gr-bw[500] | 22.33% | Failed = 77.61%, Incomplete = 0.06% | 17.69% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 1.92% | Failed = 97.86%, Incomplete = 0.22% | 0.75% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[1000] | 48.11% | Failed = 51.89%, Incomplete = 0.0% | 35.75% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 1.50% | Failed = 98.06%, Incomplete = 0.44% | 0.75% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[1500] | 75.19% | Failed = 23.06%, Incomplete = 1.75% | 60.56% | 0.00% | Failed = 0.08%, Incomplete = 99.92% | 0.00% | 3.39% | Failed = 96.25%, Incomplete = 0.36% | 2.56% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[2000] | 58.33% | Failed = 40.53%, Incomplete = 1.14% | 51.89% | 0.00% | Failed = 0.14%, Incomplete = 99.86% | 0.00% | 0.81% | Failed = 98.36%, Incomplete = 0.83% | 0.42% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[5000] | 79.00% | Failed = 20.89%, Incomplete = 0.11% | 69.00% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[10000] | 93.89% | Failed = 6.11%, Incomplete = 0.0% | 91.39% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-bw[14400] | 95.94% | Failed = 4.06%, Incomplete = 0.0% | 93.72% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[500] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 3.61% | Failed = 96.33%, Incomplete = 0.06% | 3.39% | 3.50% | Failed = 96.03%, Incomplete = 0.47% | 1.39% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[1000] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 9.33% | Failed = 90.64%, Incomplete = 0.03% | 8.81% | 3.36% | Failed = 96.06%, Incomplete = 0.58% | 2.06% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[1500] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 24.78% | Failed = 75.19%, Incomplete = 0.03% | 23.94% | 2.92% | Failed = 96.22%, Incomplete = 0.86% | 1.03% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[2000] | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 34.64% | Failed = 65.36%, Incomplete = 0.0% | 34.08% | 1.58% | Failed = 97.58%, Incomplete = 0.83% | 1.06% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[5000] | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 71.36% | Failed = 28.56%, Incomplete = 0.08% | 68.83% | 0.00% | Failed = 97.89%, Incomplete = 2.11% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-hn[10000] | 0.00% | Failed = 0.75%, Incomplete = 99.25% | 0.00% | 88.47% | Failed = 11.53%, Incomplete = 0.0% | 86.56% | 0.00% | Failed = 63.72%, Incomplete = 36.28% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |

| Model | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| plansformer-gr-hn[14400] | 0.00% | Failed = 0.08%, Incomplete = 99.92% | 0.00% | Failed = 8.94%, Incomplete = 0.0% | 91.06% | 89.86% | Failed = 6.0%, Incomplete = 94.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-gr-dl[500] | 0.00% | Failed = 0.06%, Incomplete = 99.94% | 0.00% | Failed = 0.11%, Incomplete = 99.89% | 0.00% | 0.00% | Failed = 25.28%, Incomplete = 0.0% | 43.36% | 74.72% | 31.08% | Failed = 64.28%, Incomplete = 4.64% | 18.44% |
| plansformer-gr-dl[1000] | 0.00% | Failed = 0.39%, Incomplete = 99.61% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 25.22%, Incomplete = 0.03% | 40.14% | 74.75% | 50.81% | Failed = 48.92%, Incomplete = 0.28% | 37.72% |
| plansformer-gr-dl[1500] | 0.00% | Failed = 0.06%, Incomplete = 99.94% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 0.00% | Failed = 33.67%, Incomplete = 0.03% | 31.39% | 66.31% | 61.42% | Failed = 38.47%, Incomplete = 0.11% | 40.19% |
| plansformer-gr-dl[2000] | 0.00% | Failed = 0.06%, Incomplete = 99.94% | 0.00% | Failed = 0.06%, Incomplete = 99.97% | 0.00% | 0.00% | Failed = 43.14%, Incomplete = 0.03% | 24.25% | 53.69% | 65.03% | Failed = 34.64%, Incomplete = 0.33% | 45.86% |
| plansformer-gr-dl[5000] | 0.00% | Failed = 0.08%, Incomplete = 99.92% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 0.00% | Failed = 96.0%, Incomplete = 1.42% | 1.14% | 2.58% | 76.36% | Failed = 11.08%, Incomplete = 12.56% | 60.72% |
| plansformer-gr-dl[10000] | 0.00% | Failed = 1.03%, Incomplete = 98.97% | 0.00% | Failed = 1.03%, Incomplete = 98.97% | 0.00% | 0.00% | Failed = 99.94%, Incomplete = 0.06% | 0.00% | 0.00% | 90.44% | Failed = 9.31%, Incomplete = 0.25% | 74.06% |
| plansformer-gr-dl[14400] | 0.00% | Failed = 0.86%, Incomplete = 99.14% | 0.00% | Failed = 3.22%, Incomplete = 96.78% | 0.00% | 0.00% | Failed = 98.5%, Incomplete = 1.5% | 0.00% | 0.00% | 91.97% | Failed = 7.94%, Incomplete = 0.08% | 76.14% |
| plansformer-dl-bw[500] | 20.83% | Failed = 79.17%, Incomplete = 0.0% | 15.89% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 67.89%, Incomplete = 32.11% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-bw[1000] | 40.36% | Failed = 59.36%, Incomplete = 0.28% | 32.50% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 98.36%, Incomplete = 1.64% | 0.00% | 0.00% | 0.03% | Failed = 99.53%, Incomplete = 0.44% | 0.00% |
| plansformer-dl-bw[1500] | 39.86% | Failed = 60.08%, Incomplete = 0.06% | 32.47% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-bw[2000] | 68.00% | Failed = 31.22%, Incomplete = 0.78% | 62.42% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-bw[5000] | 90.97% | Failed = 9.03%, Incomplete = 0.0% | 82.69% | Failed = 0.47%, Incomplete = 99.53% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-bw[10000] | 90.17% | Failed = 9.81%, Incomplete = 0.03% | 88.22% | Failed = 0.06%, Incomplete = 99.94% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-bw[14400] | 93.19% | Failed = 6.81%, Incomplete = 0.0% | 91.28% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |

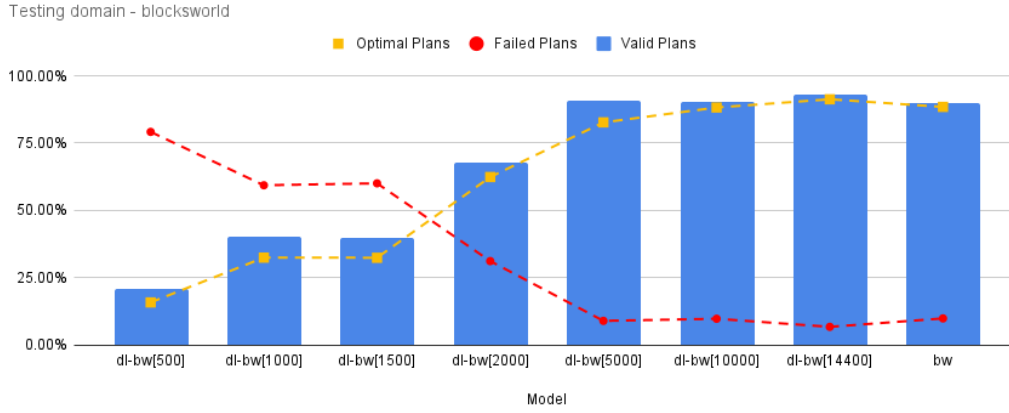| Model | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| plansformer-dl-hn[500] | 0.00% | Failed = 1.31%, Incomplete = 98.69% | 0.00% | 1.00% | Failed = 98.94%, Incomplete = 0.06% | 0.61% | 0.00% | Failed = 99.89%, Incomplete = 0.11% | 0.00% | 33.25% | Failed = 48.72%, Incomplete = 18.03% | 28.97% |
| plansformer-dl-hn[1000] | 0.00% | Failed = 1.69%, Incomplete = 98.31% | 0.00% | 3.31% | Failed = 96.69%, Incomplete = 0.0% | 3.28% | 0.00% | Failed = 3.83%, Incomplete = 96.17% | 0.00% | 10.28% | Failed = 49.78%, Incomplete = 39.94% | 9.39% |
| plansformer-dl-hn[1500] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 16.28% | Failed = 83.61%, Incomplete = 0.11% | 15.69% | 0.00% | Failed = 0.03%, Incomplete = 99.97% | 0.00% | 11.72% | Failed = 41.44%, Incomplete = 46.83% | 11.03% |
| plansformer-dl-hn[2000] | 0.00% | Failed = 0.17%, Incomplete = 99.83% | 0.00% | 13.03% | Failed = 86.69%, Incomplete = 0.28% | 12.83% | 0.00% | Failed = 4.97%, Incomplete = 95.03% | 0.00% | 5.50% | Failed = 54.81%, Incomplete = 39.69% | 5.17% |
| plansformer-dl-hn[5000] | 0.00% | Failed = 2.19%, Incomplete = 97.81% | 0.00% | 42.97% | Failed = 56.97%, Incomplete = 0.06% | 38.75% | 0.00% | Failed = 2.69%, Incomplete = 97.31% | 0.00% | 0.00% | Failed = 99.22%, Incomplete = 0.78% | 0.00% |
| plansformer-dl-hn[10000] | 0.00% | Failed = 1.28%, Incomplete = 98.72% | 0.00% | 89.67% | Failed = 10.31%, Incomplete = 0.03% | 87.72% | 0.00% | Failed = 0.22%, Incomplete = 99.78% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-hn[14400] | 0.00% | Failed = 0.53%, Incomplete = 99.47% | 0.00% | 81.86% | Failed = 18.03%, Incomplete = 0.11% | 79.81% | 0.00% | Failed = 1.06%, Incomplete = 98.94% | 0.00% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| plansformer-dl-gr[500] | 0.00% | Failed = 3.5%, Incomplete = 96.5% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 18.39% | Failed = 81.36%, Incomplete = 0.25% | 14.25% | 29.11% | Failed = 39.14%, Incomplete = 31.75% | 24.89% |
| plansformer-dl-gr[1000] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 48.72% | Failed = 51.22%, Incomplete = 0.06% | 27.56% | 45.92% | Failed = 45.03%, Incomplete = 9.06% | 38.14% |
| plansformer-dl-gr[1500] | 0.00% | Failed = 0.22%, Incomplete = 99.78% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 61.47% | Failed = 37.56%, Incomplete = 0.97% | 34.42% | 20.22% | Failed = 68.03%, Incomplete = 11.75% | 15.97% |
| plansformer-dl-gr[2000] | 0.00% | Failed = 1.72%, Incomplete = 98.28% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 62.08% | Failed = 37.83%, Incomplete = 0.08% | 33.69% | 4.50% | Failed = 91.81%, Incomplete = 3.69% | 3.67% |
| plansformer-dl-gr[5000] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 77.72% | Failed = 22.25%, Incomplete = 0.03% | 42.44% | 0.00% | Failed = 99.97%, Incomplete = 0.03% | 0.00% |
| plansformer-dl-gr[10000] | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 85.00% | Failed = 15.0%, Incomplete = 0.0% | 51.42% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |
| Further fine-tuned models — plansformer-dl-gr[14400] | 0.00% | Failed = 0.06%, Incomplete = 99.94% | 0.00% | 0.00% | Failed = 0.0%, Incomplete = 100.0% | 0.00% | 85.47% | Failed = 14.5%, Incomplete = 0.03% | 51.47% | 0.00% | Failed = 100.0%, Incomplete = 0.0% | 0.00% |

Figure 7: Plansformer-dl variants performance on blocksworld at various stages of fine-tuning
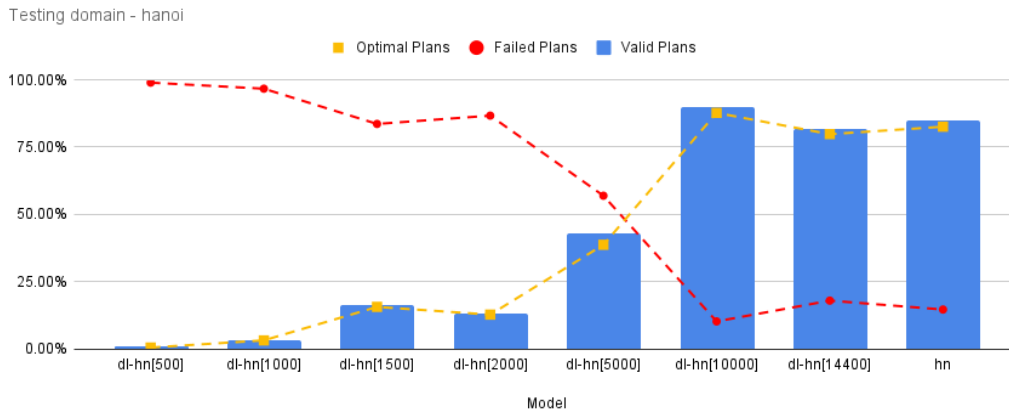


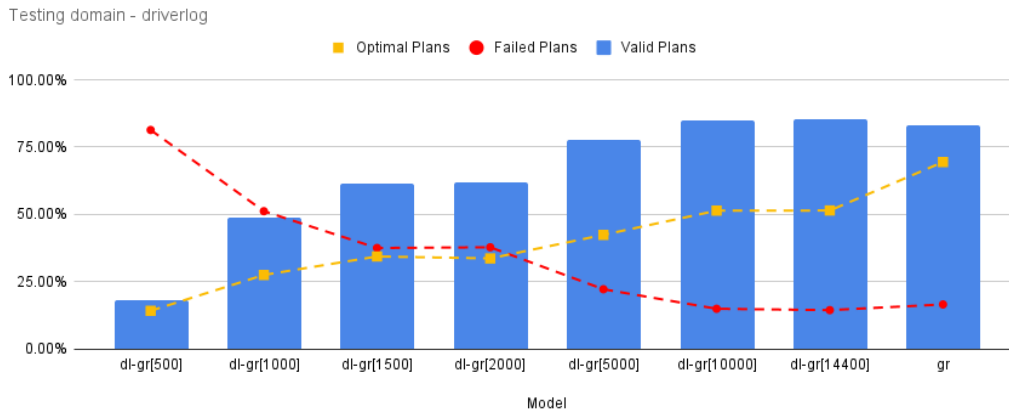Figure 8: Plansformer-dl variants performance on hanoi at various stages of fine-tuning



Figure 9: Plansformer-dl variants performance on grippers at various stages of fine-tuning
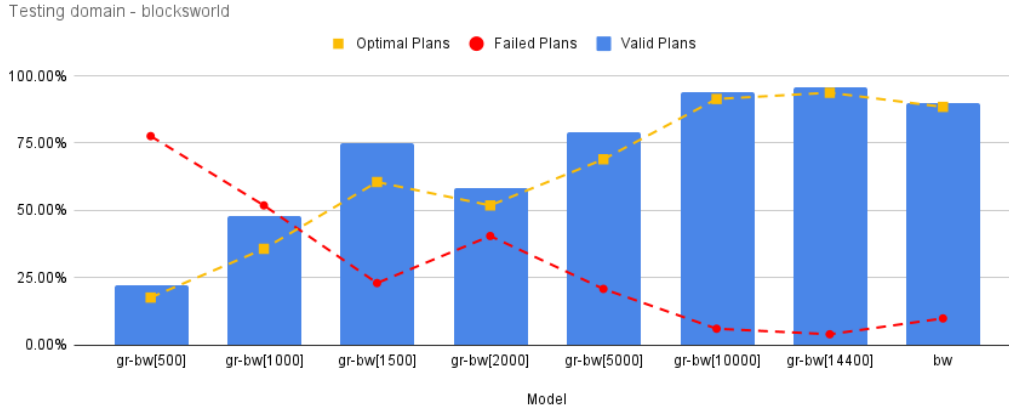
14

Figure 10: Plansformer-gr variants performance on blocksworld at various stages of fine-tuning
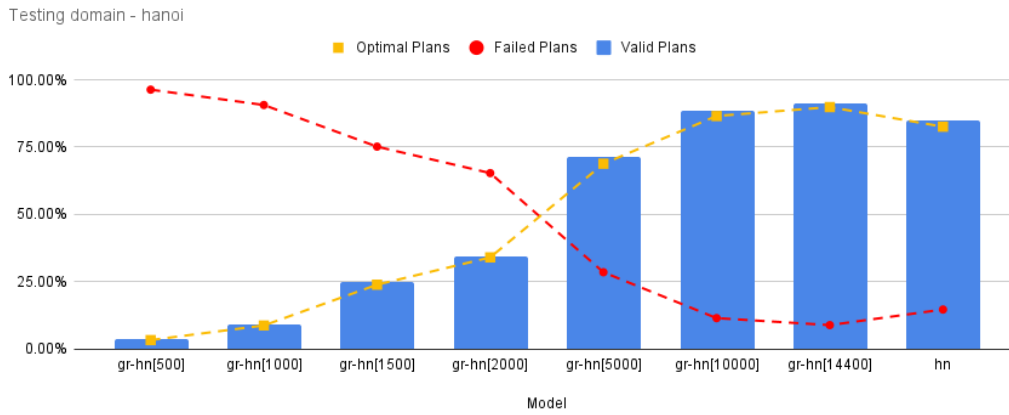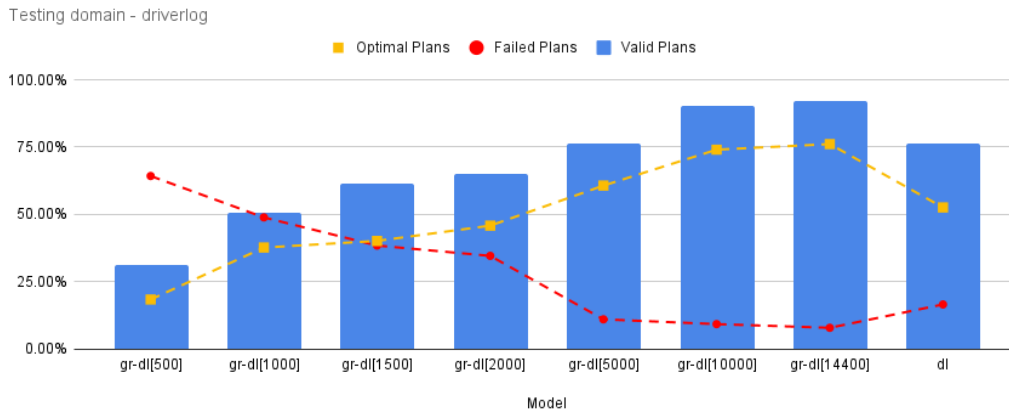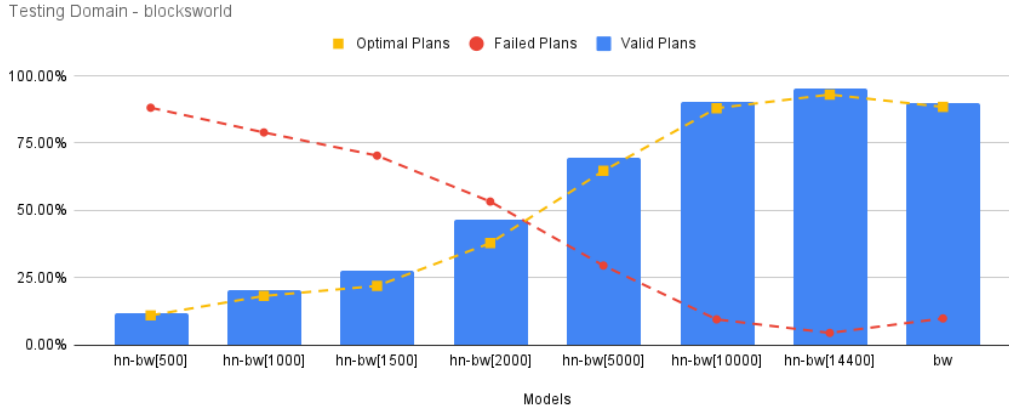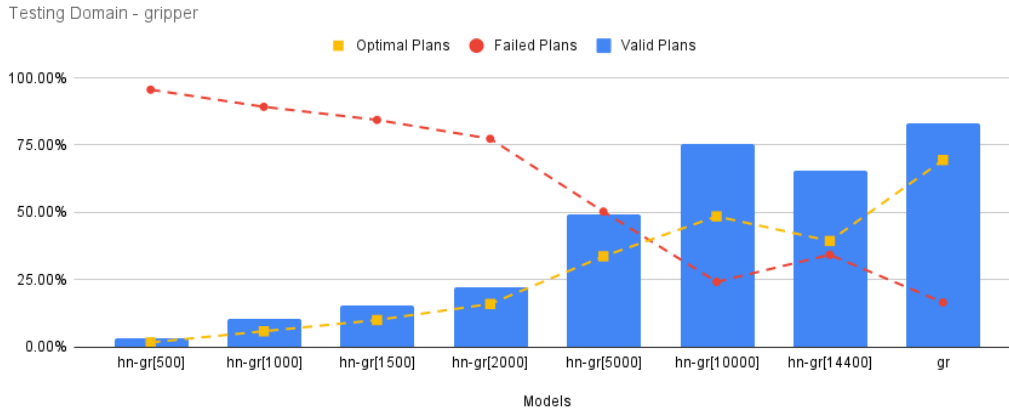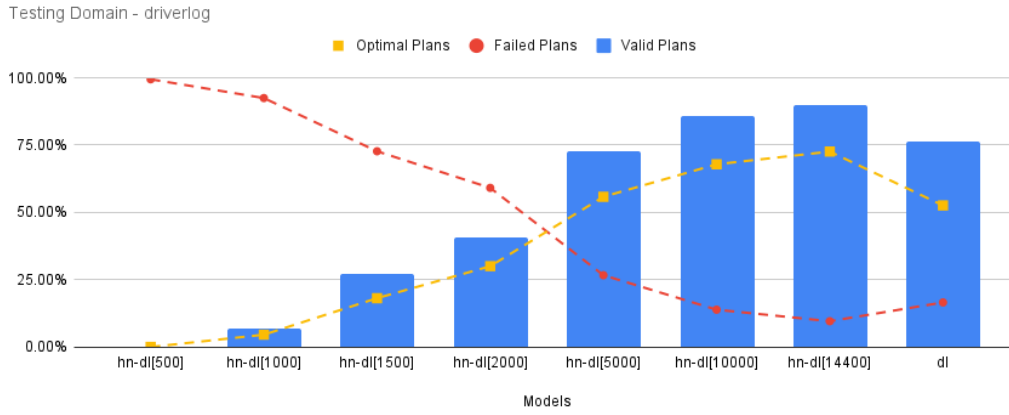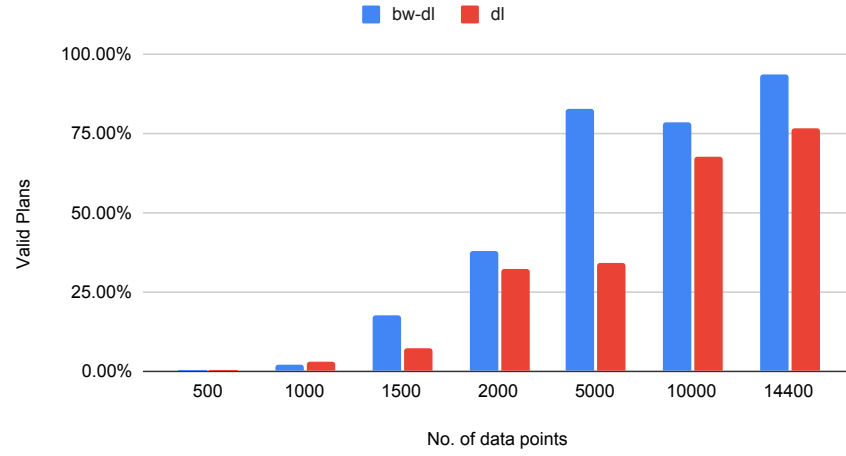


Figure 11: Plansformer-gr variants performance on hanoi at various stages of fine-tuning



Figure 12: Plansformer-gr variants performance on driverlog at various stages of fine-tuning

Figure 13: Plansformer-hn variants performance on blocksworld at various stages of fine-tuning



Figure 14: Plansformer-hn variants performance on gripper at various stages of fine-tuning



Figure 15: Plansformer-dl variants performance on driverlog at various stages of fine-tuning

Figure 16: Comparison of valid plans generated by Plansformer-bw-dl derived models with Plansformer-dl trained using similar data points.



Figure 17: Comparison of valid plans generated by Plansformer-bw-gr derived models with Plansformer-gr trained using similar data points.

```
T5ForConditionalGeneration(
  (shared): Embedding(32100, 768)
  (encoder): T5Stack(
    (embed_tokens): Embedding(32100, 768)
    (block): ModuleList(
      (0): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
              (relative_attention_bias): Embedding(32, 12)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerFF(
            (DenseReluDense): T5DenseActDense(
              (wi): Linear(in_features=768, out_features=3072, bias=False)
              (wo): Linear(in_features=3072, out_features=768, bias=False)
              (dropout): Dropout(p=0.1, inplace=False)
              (act): ReLU()
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (1): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerFF(
            (DenseReluDense): T5DenseActDense(
              (wi): Linear(in_features=768, out_features=3072, bias=False)
              (wo): Linear(in_features=3072, out_features=768, bias=False)
              (dropout): Dropout(p=0.1, inplace=False)
              (act): ReLU()
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (2): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
```

Figure 18: Plansformer Layers

```
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
    (3): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
    (4): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
```

```
      )
    )
    (5): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
    (6): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
    (7): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
```

```
    (1): T5LayerFF(
      (DenseReluDense): T5DenseActDense(
        (wi): Linear(in_features=768, out_features=3072, bias=False)
        (wo): Linear(in_features=3072, out_features=768, bias=False)
        (dropout): Dropout(p=0.1, inplace=False)
        (act): ReLU()
      )
      (layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
)
(8): T5Block(
  (layer): ModuleList(
    (0): T5LayerSelfAttention(
      (SelfAttention): T5Attention(
        (q): Linear(in_features=768, out_features=768, bias=False)
        (k): Linear(in_features=768, out_features=768, bias=False)
        (v): Linear(in_features=768, out_features=768, bias=False)
        (o): Linear(in_features=768, out_features=768, bias=False)
      )
      (layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (1): T5LayerFF(
      (DenseReluDense): T5DenseActDense(
        (wi): Linear(in_features=768, out_features=3072, bias=False)
        (wo): Linear(in_features=3072, out_features=768, bias=False)
        (dropout): Dropout(p=0.1, inplace=False)
        (act): ReLU()
      )
      (layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
)
(9): T5Block(
  (layer): ModuleList(
    (0): T5LayerSelfAttention(
      (SelfAttention): T5Attention(
        (q): Linear(in_features=768, out_features=768, bias=False)
        (k): Linear(in_features=768, out_features=768, bias=False)
        (v): Linear(in_features=768, out_features=768, bias=False)
        (o): Linear(in_features=768, out_features=768, bias=False)
      )
      (layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (1): T5LayerFF(
      (DenseReluDense): T5DenseActDense(
        (wi): Linear(in_features=768, out_features=3072, bias=False)
        (wo): Linear(in_features=3072, out_features=768, bias=False)
        (dropout): Dropout(p=0.1, inplace=False)
        (act): ReLU()
      )
      (layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
  )
)
(10): T5Block(
  (layer): ModuleList(
```

```
            (0): T5LayerSelfAttention(
              (SelfAttention): T5Attention(
                (q): Linear(in_features=768, out_features=768, bias=False)
                (k): Linear(in_features=768, out_features=768, bias=False)
                (v): Linear(in_features=768, out_features=768, bias=False)
                (o): Linear(in_features=768, out_features=768, bias=False)
              )
              (layer_norm): T5LayerNorm()
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (1): T5LayerFF(
              (DenseReluDense): T5DenseActDense(
                (wi): Linear(in_features=768, out_features=3072, bias=False)
                (wo): Linear(in_features=3072, out_features=768, bias=False)
                (dropout): Dropout(p=0.1, inplace=False)
                (act): ReLU()
              )
              (layer_norm): T5LayerNorm()
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
        )
        (11): T5Block(
          (layer): ModuleList(
            (0): T5LayerSelfAttention(
              (SelfAttention): T5Attention(
                (q): Linear(in_features=768, out_features=768, bias=False)
                (k): Linear(in_features=768, out_features=768, bias=False)
                (v): Linear(in_features=768, out_features=768, bias=False)
                (o): Linear(in_features=768, out_features=768, bias=False)
              )
              (layer_norm): T5LayerNorm()
              (dropout): Dropout(p=0.1, inplace=False)
            )
            (1): T5LayerFF(
              (DenseReluDense): T5DenseActDense(
                (wi): Linear(in_features=768, out_features=3072, bias=False)
                (wo): Linear(in_features=3072, out_features=768, bias=False)
                (dropout): Dropout(p=0.1, inplace=False)
                (act): ReLU()
              )
              (layer_norm): T5LayerNorm()
              (dropout): Dropout(p=0.1, inplace=False)
            )
          )
        )
      )
      (final_layer_norm): T5LayerNorm()
      (dropout): Dropout(p=0.1, inplace=False)
    )
    (decoder): T5Stack(
      (embed_tokens): Embedding(32100, 768)
      (block): ModuleList(
        (0): T5Block(
          (layer): ModuleList(
            (0): T5LayerSelfAttention(
              (SelfAttention): T5Attention(
                (q): Linear(in_features=768, out_features=768, bias=False)
                (k): Linear(in_features=768, out_features=768, bias=False)
                (v): Linear(in_features=768, out_features=768, bias=False)
                (o): Linear(in_features=768, out_features=768, bias=False)
                (relative_attention_bias): Embedding(32, 12)
```

```
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (1): T5LayerCrossAttention(
        (EncDecAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (1): T5Block(
    (layer): ModuleList(
      (0): T5LayerSelfAttention(
        (SelfAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (1): T5LayerCrossAttention(
        (EncDecAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (2): T5Block(
    (layer): ModuleList(
```

```
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerCrossAttention(
          (EncDecAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (2): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
    (3): T5Block(
      (layer): ModuleList(
        (0): T5LayerSelfAttention(
          (SelfAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerCrossAttention(
          (EncDecAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (2): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
```

```
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (4): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerCrossAttention(
            (EncDecAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (2): T5LayerFF(
            (DenseReluDense): T5DenseActDense(
              (wi): Linear(in_features=768, out_features=3072, bias=False)
              (wo): Linear(in_features=3072, out_features=768, bias=False)
              (dropout): Dropout(p=0.1, inplace=False)
              (act): ReLU()
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (5): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerCrossAttention(
            (EncDecAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (2): T5LayerFF(
            (DenseReluDense): T5DenseActDense(
```

```
              (wi): Linear(in_features=768, out_features=3072, bias=False)
              (wo): Linear(in_features=3072, out_features=768, bias=False)
              (dropout): Dropout(p=0.1, inplace=False)
              (act): ReLU()
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (6): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerCrossAttention(
            (EncDecAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (2): T5LayerFF(
            (DenseReluDense): T5DenseActDense(
              (wi): Linear(in_features=768, out_features=3072, bias=False)
              (wo): Linear(in_features=3072, out_features=768, bias=False)
              (dropout): Dropout(p=0.1, inplace=False)
              (act): ReLU()
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
        )
      )
      (7): T5Block(
        (layer): ModuleList(
          (0): T5LayerSelfAttention(
            (SelfAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
            )
            (layer_norm): T5LayerNorm()
            (dropout): Dropout(p=0.1, inplace=False)
          )
          (1): T5LayerCrossAttention(
            (EncDecAttention): T5Attention(
              (q): Linear(in_features=768, out_features=768, bias=False)
              (k): Linear(in_features=768, out_features=768, bias=False)
              (v): Linear(in_features=768, out_features=768, bias=False)
              (o): Linear(in_features=768, out_features=768, bias=False)
```

```
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (8): T5Block(
    (layer): ModuleList(
      (0): T5LayerSelfAttention(
        (SelfAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (1): T5LayerCrossAttention(
        (EncDecAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (9): T5Block(
    (layer): ModuleList(
      (0): T5LayerSelfAttention(
        (SelfAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
```

```
      (1): T5LayerCrossAttention(
        (EncDecAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (10): T5Block(
    (layer): ModuleList(
      (0): T5LayerSelfAttention(
        (SelfAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (1): T5LayerCrossAttention(
        (EncDecAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
          (v): Linear(in_features=768, out_features=768, bias=False)
          (o): Linear(in_features=768, out_features=768, bias=False)
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
      (2): T5LayerFF(
        (DenseReluDense): T5DenseActDense(
          (wi): Linear(in_features=768, out_features=3072, bias=False)
          (wo): Linear(in_features=3072, out_features=768, bias=False)
          (dropout): Dropout(p=0.1, inplace=False)
          (act): ReLU()
        )
        (layer_norm): T5LayerNorm()
        (dropout): Dropout(p=0.1, inplace=False)
      )
    )
  )
  (11): T5Block(
    (layer): ModuleList(
      (0): T5LayerSelfAttention(
        (SelfAttention): T5Attention(
          (q): Linear(in_features=768, out_features=768, bias=False)
          (k): Linear(in_features=768, out_features=768, bias=False)
```

```
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (1): T5LayerCrossAttention(
          (EncDecAttention): T5Attention(
            (q): Linear(in_features=768, out_features=768, bias=False)
            (k): Linear(in_features=768, out_features=768, bias=False)
            (v): Linear(in_features=768, out_features=768, bias=False)
            (o): Linear(in_features=768, out_features=768, bias=False)
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
        (2): T5LayerFF(
          (DenseReluDense): T5DenseActDense(
            (wi): Linear(in_features=768, out_features=3072, bias=False)
            (wo): Linear(in_features=3072, out_features=768, bias=False)
            (dropout): Dropout(p=0.1, inplace=False)
            (act): ReLU()
          )
          (layer_norm): T5LayerNorm()
          (dropout): Dropout(p=0.1, inplace=False)
        )
      )
    )
  )
  (final_layer_norm): T5LayerNorm()
  (dropout): Dropout(p=0.1, inplace=False)
)
(lm_head): Linear(in_features=768, out_features=32100, bias=False)
)
```