

# Improved Prediction of IDA\*’s Performance via $\epsilon$ -Truncation

**Levi Lelis**

Department of Computing Science  
University of Alberta  
Edmonton, AB, Canada T6G 2E8  
santanad@cs.ualberta.ca

**Sandra Zilles**

Department of Computer Science  
University of Regina  
Regina, SK, Canada S4S 0A2  
zilles@cs.uregina.ca

**Robert C. Holte**

Department of Computing Science  
University of Alberta  
Edmonton, AB, Canada T6G 2E8  
holte@cs.ualberta.ca

## Abstract

Korf, Reid, and Edelkamp launched a line of research aimed at predicting how many nodes IDA\* will expand with a given cost bound. This paper advances this line of research in three ways. First, we identify a source of prediction error that has hitherto been overlooked. We call it the “discretization effect”. Second, we disprove the intuitively appealing idea that a “more informed” prediction system cannot make worse predictions than a “less informed” one. More informed systems are more susceptible to the discretization effect, and in several of our experiments the more informed system makes poorer predictions. Our third contribution is a method, called “ $\epsilon$ -truncation”, which makes a prediction system less informed, in a carefully chosen way, so as to improve its predictions by reducing the discretization effect. In our experiments  $\epsilon$ -truncation rarely degraded predictions; in the vast majority of cases it improved predictions, often substantially.

## Introduction

Tree search is a popular technique for solving combinatorial problems (Pearl 1984). A frequent impediment of the application of tree searching algorithms is the inability to quickly predict the running time of an algorithm on a particular problem instance. While one instance of a problem might be solved in a blink of an eye, another instance of the same problem might take centuries. Korf, Reid, and Edelkamp (2001) launched a line of research aimed at creating a method to predict exactly how many nodes IDA\* would expand on an iteration with cost bound  $d$  given a particular heuristic function. This was in contrast with the traditional approach to search complexity analysis, which focused on “big-O” complexity typically parameterized by the accuracy of the heuristic (Dinh, Russell, and Su 2007; Gaschnig 1979; Pearl 1984; Pohl 1977).

Korf, Reid, and Edelkamp developed a prediction formula, KRE, for the special case of consistent heuristics, proved that it was exact asymptotically (in the limit of large  $d$ ), and experimentally showed that it was extremely accurate even at depths of practical interest. Zahavi, Felner, Burch, and Holte (2010) generalized KRE to work with inconsistent heuristics and to account for the heuristic values of the start states. Their formula, CDP, is described below.

Copyright © 2011, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

$d$	IDA*	Coarse	Refined	$\epsilon$ -Coarse	$\epsilon$ -Refined
18	14.5	0.72	1.16	0.79	<b>1.02</b>
19	22.2	0.73	1.19	0.74	<b>1.01</b>
20	27.4	0.74	1.24	0.82	<b>0.99</b>
21	43.3	0.74	1.28	0.77	<b>0.99</b>
22	58.5	0.75	1.34	0.84	<b>0.96</b>
23	95.4	0.76	1.39	0.81	<b>0.98</b>
24	135.7	0.76	1.45	0.85	<b>0.93</b>
25	226.7	0.77	1.52	0.84	<b>0.98</b>
26	327.8	0.77	1.58	0.86	<b>0.92</b>
27	562.0	0.77	1.64	0.86	<b>0.98</b>

Table 1: 8-puzzle, Inconsistent Heuristic. Signed Error.

The present paper advances this research in three ways. First, we identify a source of prediction error that has hitherto been overlooked. We call it the “discretization effect”. Second, we disprove the intuitively appealing idea, specifically asserted to be true by Zahavi *et al.*, that a “more informed” system cannot make worse predictions than a “less informed” system.<sup>1</sup> The possibility of this statement being false follows directly from the discretization effect, because a more informed system is more susceptible to the discretization effect than a less informed one. We will show several cases of this statement being false and use the phrase “informativeness pathology” to refer to this situation. Our final contribution is a method for counteracting the discretization effect, which we call “ $\epsilon$ -truncation”. One way to view  $\epsilon$ -truncation is that it makes a prediction system *less informed*, in a carefully chosen way, so as to improve its predictions by reducing the discretization effect. In our experiments  $\epsilon$ -truncation rarely degraded predictions; in the vast majority of cases it improved predictions, often substantially.

These contributions are illustrated in Table 1 for the 8-puzzle using the inconsistent heuristic defined by Zahavi *et al.* (see their Table 10). Column “ $d$ ” shows the cost bound for which the prediction is being made and column “IDA\*” shows how many nodes IDA\* actually expanded, on average, on its iteration with cost bound  $d$  for the given set of start states. The remaining columns represent the quality of the predictions of four prediction systems (defined below under different labels). A perfect score is 1.0 (we explain how this score is calculated in the experimental section of

<sup>1</sup>“More informed” is defined formally in Definition 3 below.

this paper). The *Refined* system is “more informed” than the *Coarse* system;  $\epsilon$ -*Coarse* and  $\epsilon$ -*Refined* are the resulting systems when  $\epsilon$ -truncation is applied to *Coarse* and *Refined*, respectively. The number given is the ratio of the average predicted number of nodes expanded to the average number of nodes actually expanded. For example, the entry in the bottom row of the *Coarse* column being 0.77 means that for  $d = 27$ , the prediction by the *Coarse* system was 77% of the actual IDA\* value (*Coarse* predicted 432.7 nodes would be expanded, on average, whereas in reality 562.0 were). *Coarse* is the  $\text{CDP}_2$  system of Zahavi *et al.* as it was applied to the 8-puzzle and the *Coarse* column here exactly reproduces the  $\text{CDP}_2$  column in their Table 10. These are among the least accurate predictions Zahavi *et al.* report for  $\text{CDP}$ , and they offer no explanation for these poor results. Even though *Refined* is “more informed” than *Coarse*, Table 1 shows that its predictions are worse than *Coarse*’s for  $d \geq 21$ . The last two columns show the effect of applying  $\epsilon$ -truncation: in both cases the predictions are substantially improved, and the more informed system ends up outperforming the less informed system.

## The CDP Prediction Framework

Here we briefly sketch the  $\text{CDP}$  system using our own notation. The reader is referred to the original paper (Zahavi *et al.* 2010) for full explanations and illustrative examples. For reference, Table 2 summarizes the notation introduced in this section.

Let  $S$  be the set of states,  $E \subseteq S \times S$  the set of directed edges over  $S$  representing the parent-child relation in the underlying state space, and  $h : S \rightarrow \mathbb{N}$  the heuristic function.

**Definition 1**  $T = \{t_1, \dots, t_n\}$  is a type system for  $(S, E)$  if it is a disjoint partitioning of  $E$ . For every  $(\hat{s}, s) \in E$ ,  $T(\hat{s}, s)$  denotes the unique  $t \in T$  with  $(\hat{s}, s) \in t$ .

**Definition 2** Let  $t, t' \in T$ .  $p(t'|t)$  denotes the probability that a node  $s$  with parent  $\hat{s}$  and  $T(\hat{s}, s) = t$  generates a node  $c$  with  $T(s, c) = t'$ .  $b_t$  denotes the average number of children generated by a node  $s$  with parent  $\hat{s}$  with  $T(\hat{s}, s) = t$ .

Parent-pruning is an easy-to-implement enhancement of IDA\* that avoids re-expanding the parent of a node. IDA\* with parent-pruning will not generate a node  $\hat{s}$  from  $s$  if  $\hat{s}$  is the parent of  $s$ . In making the prediction of the number of nodes expanded on an iteration of IDA\* we are interested in estimating the number of nodes in the subtree below a given node. Because of parent-pruning the subtree below a node differs depending on the node from which it was generated. Like Zahavi *et al.*, to account for this effect of parent-pruning we define types over node pairs instead of just nodes.<sup>2</sup>

All type systems considered in this paper have the property that  $h(s) = h(s')$  if  $T(\hat{s}, s) = T(\hat{s}', s')$ . We assume this property in the formulae below, and denote by  $h(t)$  the value  $h(s)$  for any  $s, \hat{s}$  such that  $T(\hat{s}, s) = t$ .

<sup>2</sup>Zahavi *et al.* call the type definition over pairs a “two-step” model.

$\text{CDP}$  samples the state space in order to estimate  $p(t'|t)$  and  $b_t$  for all  $t, t' \in T$ . We denote by  $\pi(t'|t)$  and  $\beta(t)$  the respective estimates thus obtained. The predicted number of nodes expanded by IDA\* (with parent pruning) for start state  $\hat{s}^*$ , cost bound  $d$ , heuristic  $h$ , and type system  $T$  is

$$\text{CDP}(\hat{s}^*, d, h, T) = \sum_{(\hat{s}^*, s^*) \in E} \sum_{i=1}^d \sum_{t \in T} N(i, t, (\hat{s}^*, s^*), d).$$

Here  $N(i, t, (\hat{s}^*, s^*), d)$  is the number of pairs  $(\hat{s}, s)$  with  $T(\hat{s}, s) = t$  and  $s$  at level  $i$  of the search tree rooted at  $s^*$ . It is computed recursively as follows.

$$N(1, t, (\hat{s}^*, s^*), d) = \begin{cases} 0 & \text{if } T(\hat{s}^*, s^*) \neq t, \\ 1 & \text{if } T(\hat{s}^*, s^*) = t, \end{cases}$$

and, for  $i > 1$ , the value  $N(i, t, (\hat{s}^*, s^*), d)$  is given by

$$\sum_{u \in T} N(i-1, u, (\hat{s}^*, s^*), d) \pi(t|u) \beta_t P(t, i, d) \quad (1)$$

where  $P(t, i, d) = 1$  if  $h(t) + i \leq d$ , and is 0 otherwise.

According to the formulae above, in order to predict the number of nodes IDA\* expands with a cost bound  $d$ , for every level  $i \leq d$ ,  $\text{CDP}$  predicts how many instances of each type will be generated; *i.e.*, it predicts a vector  $(N[1], \dots, N[|T|])$  of numbers of instances of each type on a level.<sup>3</sup> We will call such a vector a *type allocation vector*. The type allocation vector for the first level of prediction is computed by verifying the type of the children of the start state (the  $i = 1$  base case of the recursive calculation shown above). Once the allocation vector is calculated for the first level, the vector for the next level is estimated according to Equation 1. At level  $i$ , for each type  $t$  such that  $h(t) + i$  exceeds the cost bound  $d$ , the corresponding entry in the type allocation vector,  $N[t]$ , is set to zero to indicate that IDA\* will prune nodes of this type from its search.<sup>4</sup> The prediction continues to deeper and deeper levels as long as there is an entry in the type allocation vector greater than zero.

As our basic type system,  $T_h$ , we use Zahavi *et al.*’s basic “two-step” model, defined (in our notation) as  $T_h(\hat{s}, s) = (h(\hat{s}), h(s))$ . Two new domain-independent type systems we will also use, which are “more informed” than  $T_h$ , are:

$T_c(\hat{s}, s) = (T_h(\hat{s}, s), c((\hat{s}, s), 0), \dots, c((\hat{s}, s), H))$ , where  $c((\hat{s}, s), k)$  is the number of children of  $s$ , considering parent pruning, whose  $h$ -value is  $k$ , and  $H$  is the maximum  $h$ -value observed in the sampling process;

$T_{gc}(\hat{s}, s) = (T_c(\hat{s}, s), gc((\hat{s}, s), 0), \dots, gc((\hat{s}, s), H))$ , where  $gc((\hat{s}, s), k)$  is the number of grandchildren of  $s$ , considering parent pruning, whose  $h$ -value is  $k$ .

The intuitive concept of one type system being “more informed” than another is captured formally as follows.

<sup>3</sup>We use  $N[t]$  to denote  $N(i, t, (\hat{s}^*, s^*), d)$  when  $i, (\hat{s}^*, s^*)$  and  $d$  are clear from context.

<sup>4</sup>That is why we ensure all nodes mapped to a type have the same heuristic value, as mentioned above.

Notation	Meaning
$T$	type system
$T(\hat{s}, s)$	type of a pair $(\hat{s}, s) \in E$
$p(t' t)$	probability of type $t$ generating type $t'$
$b_t$	average branching factor of node pairs of type $t$
$\pi(t' t)$	approximation of $p(t' t)$
$\beta_t$	approximation of $b_t$
$P(t, i, d)$	pruning function
$N(i, t, (\hat{s}^*, s^*), d)$	number of node pairs of type $t$ at level $i$
$T_h(\hat{s}, s)$	type defined as $(h(\hat{s}), h(s))$
$T_c(\hat{s}, s)$	type defined by $T_h$ and the $h$ -values of $s$ 's children
$T_{gc}(\hat{s}, s)$	type defined by $T_c$ and the $h$ -values of $s$ 's grandchildren
$T_1 \prec T_2$	$T_1$ is a refinement of $T_2$

Table 2: Notation used in the CDP prediction framework.

**Definition 3** Let  $T_1, T_2$  be type systems.  $T_1$  is a refinement of  $T_2$ , denoted  $T_1 \prec T_2$ , if  $|T_1| > |T_2|$  and for all  $t_1 \in T_1$  there is a  $t_2 \in T_2$  with  $\{(\hat{s}, s) | T_1(\hat{s}, s) = t_1\} \subseteq \{(\hat{s}, s) | T_2(\hat{s}, s) = t_2\}$ . If  $t_1 \in T_1$  and  $t_2 \in T_2$  are related in this way, we write  $T_2(t_1) = t_2$ .

Note that  $T_{gc} \prec T_c \prec T_h$ , and so, by transitivity,  $T_{gc} \prec T_h$ .

Intuitively, if  $T_1 \prec T_2$  one would expect predictions using  $T_1$  to be at least as accurate as the predictions using  $T_2$ , since all the information that is being used by  $T_2$  to condition its predictions is also being used by  $T_1$  ((Zahavi et al. 2010), p. 59). However, our experiments show that this is not always true (e.g. Table 1, where *Refined*  $\prec$  *Coarse*). The underlying cause of poorer predictions by  $T_1$  when  $T_1 \prec T_2$  is the discretization effect, which we will now describe.

### The Discretization Effect

Consider for example the situation depicted in Figure 1. Here type  $t_1$  generates type  $t_2$  with probability  $p(t_2|t_1) = p$  and type  $t_3$  with probability  $p(t_3|t_1) = 1 - p$ . Let  $b_{t_1} = 1$ , i.e., one instance of type  $t_1$  generates 1 new instance on average. Assuming, for simplicity, that the estimates  $\pi(t'|t)$  and  $\beta(t)$  are accurate, CDP predicts that  $p$  instances of type  $t_2$  and  $(1 - p)$  instances of type  $t_3$  will be generated by one instance of  $t_1$ . If level  $i - 1$  consists only of one instance, which is of type  $t_1$ , CDP's prediction  $(N[1], N[2], N[3]) = (0, p, (1 - p))$  for level  $i$  minimizes the mean squared error (MSE) compared to the actual expected type allocation vector  $(a_1, a_2, a_3)$  on level  $i$ . Note that the ultimate task of CDP is to minimize the error in predicting the number of nodes expanded by IDA\*. There is no guarantee that minimizing MSE for type allocation vectors at each level will actually minimize the overall MSE of the predicted number of nodes expanded, as will be illustrated below.

In the example above, if  $p = 0.05$ , CDP predicts that 0.05 instances of type  $t_2$  and 0.95 instances of type  $t_3$  will be generated, and thus minimizes the expected MSE of the type allocation vector. In an actual search tree, only integer values are possible, hence it is most likely that 1 instances of type  $t_3$  and no instances of type  $t_2$  are generated. Intuitively, the case that an instance of type  $t_2$  is generated happens so rarely in practice that it might be better to predict that it never happens. That does not mean that fractional predictions in the type allocation vectors should be avoided

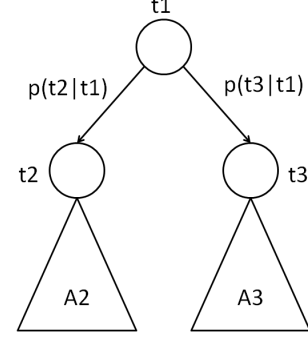


Figure 1: A hypothetical situation in CDP prediction.

in general, but in practice it might be advantageous to avoid very small fractions.

To illustrate that it is at least possible that minimizing MSE for type allocation vectors does not minimize MSE of the number of nodes predicted to be expanded, consider again Figure 1. Let us denote by  $A_2$  and  $A_3$  the actual average number of nodes expanded below an instance of type  $t_2$  and type  $t_3$ , respectively, and by  $E_2$  and  $E_3$  CDP's predictions for the number of nodes in the subtrees below an instance of type  $t_2$  and an instance of type  $t_3$ , respectively. The expected MSE for CDP's prediction would be

$$p(pE_2 + (1-p)E_3 - A_2)^2 + (1-p)(pE_2 + (1-p)E_3 - A_3)^2. \quad (2)$$

The MSE for predicting only 1 instance of type  $t_3$  (and no instances of type  $t_2$ ) of a system estimating  $A_3$  with  $E'_3$  would be

$$p(E'_3 - A_2)^2 + (1-p)(E'_3 - A_3)^2. \quad (3)$$

It is mathematically possible for the value of (2) to be larger than that of (3), even if  $E_3 = E'_3$ , as can be verified by a simple algebraic calculation, which we skip because of space constraints. This by itself does not justify changing the way CDP makes predictions, but it shows at least that minimizing MSE for type allocation vectors does not guarantee minimizing MSE in the overall prediction.

Hence we would like to examine the possibility that ignoring very rare types, i.e., in the example above, predicting that type  $t_2$  will not be generated at all if  $p$  is very small, might in practice lead to better predictions—we call this phenomenon the discretization effect. Empirical evidence for the discretization effect would also explain what we observed in the introduction, namely that more refined type systems may lead to worse CDP predictions. The reason is that more refined type systems sometimes contain many very small type generation probabilities  $p(t'|t)$ . For example, for the 8-puzzle using the inconsistent heuristic defined by Zahavi et al., 7.3% of the types generated with probability greater than zero are generated with probability lower

Level ( $i$ )	10	11-12	13	14	15	16-17	18-19	20-23	24
$\epsilon_i$	0.05	0.07	0.08	0.07	0.05	0.04	0.03	0.01	0.00

Table 3:  $\epsilon_i$  values for the 8-puzzle with Manhattan Distance.

than 1% for the more refined type system. This percentage is only 0.3% for the less refined type system.

The discretization effect hence would (i) explain why more refined type systems may behave poorly, and (ii) suggest that ignoring rare types may improve predictions overall. The goal of this paper is to show that our intuition is correct and to provide a method for systematically ignoring rare types in CDP, thus improving CDP predictions substantially in practice.

## The $\epsilon$ -Truncation Prediction Method

Our approach to avoiding the discretization effect, which we call  $\epsilon$ -truncation, can be summarized as follows.

1. As before, sample the state space to obtain  $\pi(t|u)$ .
2. Compute a cutoff value  $\epsilon_i$  for each  $i$  between 1 and  $d$ .
3. Use  $\epsilon_i$  to define  $\pi^i(t|u)$ , a version of  $\pi(t|u)$  that is specific to level  $i$ . In particular, if  $\pi(t|u) < \epsilon_i$  then  $\pi^i(t|u) = 0$ ; the other  $\pi^i(t|u)$  are set by scaling up the corresponding  $\pi(t|u)$  values so that the  $\pi^i(t|u)$  sum to 1.
4. In computing CDP use  $\pi^i(t|u)$  at level  $i$  instead of  $\pi(t|u)$ .

The key step in this process is the calculation of the  $\epsilon_i$  values. A full description of this calculation is given in the Appendix. It requires computing CDP predictions for a set of start states and, for each level  $i$  in each of these prediction calculations, solving a set of small linear programs (one for each “supertype” at that level, as defined in the Appendix). Table 3 shows the  $\epsilon_i$  values calculated using 10,000 randomly generated start states for the 8-puzzle with Manhattan Distance. In practice, useful  $\epsilon_i$  values can be computed using a much smaller number of start states. Also, as observed in Table 3 and in all the experiments in this paper, the  $\epsilon_i$  values converge to zero as  $i$  gets larger.

## Experimental Results

This section presents the results of experiments showing that: (a) refining a type system often reduces prediction accuracy; and (b)  $\epsilon$ -truncation often substantially improves predictions. Each experiment will use two type systems, a basic one and a refinement of the basic one, and will compare the predictions made by CDP with each type system and with  $\epsilon$ -truncation applied to the refined type system. We omit the results of applying  $\epsilon$ -truncation to the basic type system for lack of space; as in Table 1 it never produced large improvements and never produced predictions that were more accurate than those of the refined system with  $\epsilon$ -truncation. Our experiments are run on two domains: the sliding-tile puzzle, which has a small branching factor and deep solutions, and the pancake puzzle, which has a large branching factor and shallow solutions. We use at least two sizes for each domain: one that is small enough that the

entire reachable portion of the state space can be enumerated and used in lieu of “sampling”, and one that is large enough to be of practical interest. The small domains are an important element of the experiments because phenomena witnessed in them cannot be attributed to sampling effects. For each domain we use at least one consistent heuristic and one inconsistent one.

The choice of the set of start states will be described in the specific sections below, but we always applied the same principle as Zahavi *et al.* (2010): start state  $s$  is included in the experiment with cost bound  $d$  only if IDA\* would actually have used  $d$  as a cost bound in its search with  $s$  as the start state. Unlike an actual IDA\* run, we count the number of nodes expanded in the entire iteration for a start state even if the goal is encountered during the iteration.

As in Table 1 for each prediction system we will report the ratio of the predicted number of nodes expanded, averaged over all the start states, to the actual number of nodes expanded, on average, by IDA\*. This ratio will be rounded to two decimal places. Thus a ratio of 1.00 does not necessarily mean the prediction is perfect. This ratio we call the (average) signed error. It is the same as the “Ratio” reported by Zahavi *et al.* (2010) and is appropriate when one is interested in predicting the total number of nodes that will be expanded in solving a set of start states. It is not appropriate for measuring the accuracy of the predictions on individual start states because errors with a positive sign cancel errors with a negative sign. If these exactly balance out, a system will appear to have no error (a ratio of 1.00) even though there might be substantial error in every single prediction. To evaluate the accuracy of individual predictions, an appropriate measure is absolute error. For each instance one computes the absolute value of the difference between the predicted and the actual number of nodes expanded, divides this difference by the actual number of nodes expanded, adds these up over all start states, and divides by the total number of start states. A perfect score according to this measure is 0.0. Appendix B shows the corresponding results when root mean squared error (RMSE) is chosen as a measure (recall that CDP minimizing MSE and thus RMSE for type allocation vectors does *not* imply that CDP minimizes RMSE for the overall prediction; the numbers in the appendix clearly demonstrate that CDP is outperformed by  $\epsilon$ -truncation in terms of RMSE).

Zahavi *et al.* (2010) introduced a method for improving predictions for single start states. Instead of directly predicting how many nodes will be expanded for cost bound  $d$  and start state  $s$ , all states,  $S_r$ , at depth  $r < d$  are enumerated and one then predicts how many nodes will be expanded for depth bound  $d - r$  when  $S_r$  is the set of start states. We applied this technique in all our experiments. The value of  $r$  for each experiment is specified below.

The number of start states used to determine the  $\epsilon_i$  is closely related to the  $r$ -value that will be used in the experiment. For example, the number of states (and thus the number of linear programs) at level 10 of the 8-puzzle is expected to be much lower than the number of states at level 25 of the 15-puzzle. Therefore, in order to find a suitable  $\epsilon_{10}$ -value for the 8-puzzle we have to use more start states than

is required to determine the  $\epsilon_{25}$ -value for the 15-puzzle. The number of states used to determine the  $\epsilon_i$  is stated below for each experiment.

### Small Domains

The 10-pancake puzzle has  $10!$  states and a maximum optimal solution depth of 11. We used 100 random start states to determine the  $\epsilon_i$  and 5,000 to measure prediction accuracy. The heuristic used was a pattern database (PDB (Culberson and Schaeffer 1996)) based on the smallest four pancakes. Heuristics defined by PDBs are admissible and consistent. We used  $T_h$  and  $T_{gc}$  as the type systems. The results of this experiment are shown in the upper part of Table 4. The lower part of Table 4 shows the same experiment when the heuristic is multiplied by 1.5, which makes it inconsistent and inadmissible (e.g., there is a  $d=12$  row even though the maximum optimal solution depth is 11). The bold entries in this and all other tables of results indicate the best predictions. In both these experiments the refinement  $T_{gc}$  substantially improves on  $T_h$  and  $\epsilon$ -truncation automatically detects that  $T_{gc}$  contains almost no rare types and therefore makes little change to its predictions. For all the type systems the absolute error is substantially higher than the signed error suggests it should be, indicating that a combination of over- and under-estimation is occurring.

d	IDA*	Signed Error			Absolute Error		
		$T_h$	$T_{gc}$	$\epsilon-T_{gc}$	$T_h$	$T_{gc}$	$\epsilon-T_{gc}$
Admissible and Consistent Heuristic							
7	1,710.8	1.01	<b>1.00</b>	<b>1.00</b>	0.33	<b>0.02</b>	0.03
8	12,082.4	1.03	<b>1.00</b>	<b>1.00</b>	0.27	<b>0.03</b>	<b>0.03</b>
9	79,650.5	1.08	1.01	<b>1.00</b>	0.25	<b>0.03</b>	0.04
10	507,640.4	1.14	1.01	<b>1.00</b>	0.24	<b>0.03</b>	0.04
11	3,449,158.0	1.20	1.01	<b>1.00</b>	0.22	<b>0.02</b>	0.03
The heuristic above multiplied by 1.5							
7	528.5	1.03	<b>1.01</b>	<b>1.01</b>	0.68	<b>0.13</b>	<b>0.13</b>
8	3,740.3	1.06	<b>1.01</b>	<b>1.01</b>	0.51	<b>0.11</b>	<b>0.11</b>
9	13,924.7	1.13	<b>1.02</b>	<b>1.02</b>	0.57	<b>0.12</b>	<b>0.12</b>
10	64,611.8	1.22	1.04	<b>1.03</b>	0.50	<b>0.11</b>	<b>0.11</b>
11	356,366.1	1.36	1.06	<b>1.05</b>	0.46	<b>0.10</b>	<b>0.10</b>
12	2,397,748.6	1.36	1.09	<b>1.05</b>	0.39	<b>0.09</b>	<b>0.09</b>

Table 4: 10-pancake puzzle.  $r=1$ .

For the sliding-tile puzzle we ran experiments on two small sizes, the (3x3)-puzzle (the 8-puzzle) and the (3x4)-puzzle. We used the same type system as Zahavi *et al.* (2010), which is a refinement of  $T_h$  we call  $T_{h,b}$ .  $T_{h,b}$  is defined by  $T_{h,b}(s, s') = (T_h, \text{blank}(s), \text{blank}(s'))$  where  $\text{blank}(s)$  returns the kind of location (corner, edge, or middle) the blank occupies in state  $s$ .<sup>5</sup>  $T_{gc,b}$  is defined analogously. For square versions of the puzzle  $T_{gc}$  is exactly the same as  $T_{gc,b}$  and therefore  $T_{gc} \prec T_{h,b}$ . However, for the (3x4)-puzzle,  $T_{gc}$  and  $T_{gc,b}$  are not the same. For the 8-puzzle we used 1,000 random start states to determine the  $\epsilon_i$  and every solvable state in the space to measure prediction

<sup>5</sup>For the (3x4)-puzzle there are two kinds of edge locations that  $\text{blank}(s)$  needs to distinguish—edge locations on the short side (length 3) and edge locations on the long side (length 4).

accuracy. Table 5 shows the results for Manhattan Distance, which is admissible and consistent, with  $r=10$ . Here we see the informativeness pathology:  $T_{gc}$ 's predictions are worse than  $T_{h,b}$ 's, despite its being a refinement of  $T_{h,b}$ . Applying  $\epsilon$ -truncation substantially reduces  $T_{gc}$ 's prediction error.

d	IDA*	Signed Error			Absolute Error		
		$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
18	134.4	1.01	1.01	<b>1.00</b>	0.03	<b>0.01</b>	<b>0.01</b>
19	238.4	1.01	1.01	<b>1.00</b>	0.04	<b>0.02</b>	<b>0.02</b>
20	360.1	1.01	1.02	<b>0.99</b>	0.06	<b>0.03</b>	<b>0.03</b>
21	630.7	1.01	1.02	<b>0.99</b>	0.06	<b>0.03</b>	<b>0.03</b>
22	950.6	1.02	1.04	<b>0.99</b>	0.08	0.05	<b>0.04</b>
23	1,649.5	<b>1.01</b>	1.04	<b>0.99</b>	0.08	0.06	<b>0.04</b>
24	2,457.5	<b>1.01</b>	1.07	0.98	0.09	0.08	<b>0.06</b>
25	4,245.5	<b>1.00</b>	1.08	0.98	0.09	0.08	<b>0.05</b>
26	6,294.4	<b>1.00</b>	1.10	0.97	0.10	0.11	<b>0.07</b>
27	10,994.9	<b>0.99</b>	1.11	0.97	0.11	0.11	<b>0.06</b>

Table 5: 8-puzzle, Manhattan Distance.  $r=10$ .

The inconsistent heuristic we used for the 8-puzzle is the one defined by Zahavi *et al.* (2010). Two PDBs were built, one based on tiles 1-4, and one based on tiles 5-8. The locations in the puzzle are numbered left-to-right and top-to-bottom and the first PDB is consulted for states having the blank in an even location and the second PDB is consulted otherwise. Since the blank's location changes parity every time it moves, we are guaranteed that the heuristic value of a child node will be drawn from a different PDB than its parent. Again we used 1,000 random start states to determine the  $\epsilon_i$  and every solvable state in the space to measure prediction accuracy. The results of this experiment, with  $r=1$ , are shown in Table 6. The Signed Error columns of Table 6 reproduce the results shown in Table 1. As discussed in the Introduction, they exhibit the informativeness pathology and demonstrate that  $\epsilon$ -truncation can substantially reduce prediction error. The pathology and the  $\epsilon$ -truncation improvements are also observed when absolute error is the measure.

d	IDA*	Signed Error			Absolute Error		
		$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
18	14.5	0.72	1.16	<b>1.02</b>	0.38	0.44	<b>0.36</b>
19	22.2	0.73	1.19	<b>1.01</b>	0.45	0.49	<b>0.39</b>
20	27.4	0.74	1.24	<b>0.99</b>	0.45	0.55	<b>0.41</b>
21	43.3	0.74	1.28	<b>0.99</b>	0.49	0.57	<b>0.40</b>
22	58.5	0.75	1.34	<b>0.96</b>	0.47	0.62	<b>0.41</b>
23	95.4	0.76	1.39	<b>0.98</b>	0.48	0.62	<b>0.38</b>
24	135.7	0.76	1.45	<b>0.93</b>	0.46	0.66	<b>0.37</b>
25	226.7	0.77	1.52	<b>0.98</b>	0.44	0.65	<b>0.34</b>
26	327.8	0.77	1.58	<b>0.92</b>	0.41	0.68	<b>0.31</b>
27	562.0	0.77	1.64	<b>0.98</b>	0.40	0.67	<b>0.27</b>

Table 6: 8-puzzle, Inconsistent Heuristic.  $r=1$ .

For the (3x4)-puzzle we used 10 random start states to determine the  $\epsilon_i$  and 10,000 to measure prediction accuracy. The upper part of Table 7 shows the results for Manhattan Distance.  $T_{gc,b}$ 's absolute error is very close to  $T_{h,b}$ 's, so being more informed provides no advantage.  $\epsilon$ -truncation very substantially improves  $T_{gc,b}$ 's absolute error. The lower

part of the table is for Manhattan Distance multiplied by 1.5, which is inadmissible and inconsistent. Here  $T_{gc,b}$ 's predictions are considerably more accurate than  $T_{h,b}$ 's and are substantially improved by  $\epsilon$ -truncation.

		Signed Error			Absolute Error		
d	IDA*	$T_{h,b}$	$T_{gc,b}$	$\epsilon-T_{gc,b}$	$T_{h,b}$	$T_{gc,b}$	$\epsilon-T_{gc,b}$
Manhattan Distance							
33	30,461.9	1.02	1.02	<b>1.01</b>	0.02	<b>0.01</b>	<b>0.01</b>
34	49,576.8	1.03	1.03	<b>1.01</b>	0.03	0.02	<b>0.01</b>
35	80,688.2	1.04	1.04	<b>1.02</b>	0.04	0.03	<b>0.02</b>
36	127,733.4	1.05	1.06	<b>1.02</b>	0.05	0.05	<b>0.03</b>
37	201,822.7	1.08	1.09	<b>1.04</b>	0.07	0.07	<b>0.04</b>
38	327,835.3	1.10	1.11	<b>1.04</b>	0.09	0.09	<b>0.05</b>
39	478,092.5	1.13	1.16	<b>1.06</b>	0.11	0.13	<b>0.06</b>
40	822,055.4	1.17	1.20	<b>1.07</b>	0.15	0.17	<b>0.08</b>
41	1,163,312.1	1.21	1.27	<b>1.10</b>	0.18	0.23	<b>0.10</b>
42	1,843,732.2	1.27	1.34	<b>1.13</b>	0.23	0.30	<b>0.13</b>
Weighted Manhattan Distance ( $w=1.5$ )							
33	926.2	1.05	1.01	<b>1.00</b>	0.02	<b>0.01</b>	<b>0.01</b>
34	1,286.9	1.07	1.02	<b>1.01</b>	0.03	<b>0.01</b>	<b>0.01</b>
35	2,225.6	1.09	1.03	<b>1.00</b>	0.05	<b>0.02</b>	<b>0.02</b>
36	2,670.7	1.12	1.05	<b>1.00</b>	0.06	0.03	<b>0.02</b>
37	3,519.5	1.14	1.07	<b>0.99</b>	0.08	0.04	<b>0.02</b>
38	5,570.8	1.19	1.10	<b>0.99</b>	0.13	0.07	<b>0.03</b>
39	6,983.8	1.23	1.13	<b>0.98</b>	0.15	0.09	<b>0.04</b>
40	9,103.3	1.29	1.18	<b>0.97</b>	0.20	0.12	<b>0.05</b>
41	13,635.3	1.36	1.24	<b>0.96</b>	0.28	0.18	<b>0.06</b>
42	16,634.2	1.43	1.31	<b>0.95</b>	0.33	0.23	<b>0.07</b>

Table 7: (3x4)-puzzle.  $r=20$ .

## Large Domains

For the 15-pancake puzzle, we used 10 random start states to determine the  $\epsilon_i$  and 1,000 to measure prediction accuracy while using the consistent heuristic and 1,000 while using the inadmissible and inconsistent heuristic. We used  $T_h$  and  $T_c$  as the type systems. To define  $\pi(t|u)$  and  $\beta_t$ , 100 million random states were sampled and, in addition, we used the process described by Zahavi *et al.* (2010) to non-randomly extend the sampling: we sampled the child of a sampled state if the type of that child had not yet been sampled. The results with  $r=4$  and a PDB heuristic based on the smallest eight pancakes are shown in the upper part of Table 8. In both cases  $T_c$  outperforms  $T_h$  but is also substantially improved by  $\epsilon$ -truncation.

For the 15-puzzle, we used 5 random start states to determine the  $\epsilon_i$  and 1,000 to measure prediction accuracy. To define  $\pi(t|u)$  and  $\beta_t$ , one billion random states were sampled and, in addition, we used the extended sampling process described for the 15-pancake puzzle. Table 9 gives the results when Manhattan Distance is the heuristic and  $T_{h,b}$  and  $T_{gc}$  are the type systems. Here again we see the informativeness pathology ( $T_{h,b}$ 's predictions are better than  $T_{gc}$ 's) which is eliminated by  $\epsilon$ -truncation.

Like for the 8-puzzle, an inconsistent heuristic for the 15-puzzle was created with one PDB based on tiles 1-7, and the other based on tiles 9-15, exactly as Zahavi *et al.* did (see their Table 11). The results with  $T_{h,b}$  and  $T_c$  as the type systems are shown in Table 10. Here we see that even though

		Signed Error			Absolute Error		
d	IDA*	$T_h$	$T_c$	$\epsilon-T_c$	$T_h$	$T_c$	$\epsilon-T_c$
Admissible and Consistent Heuristic (PDB)							
9	1,035.2	1.11	1.06	<b>1.04</b>	0.23	0.10	<b>0.09</b>
10	5,547.9	1.11	1.06	<b>1.02</b>	0.36	0.15	<b>0.11</b>
11	46,009.3	1.11	1.05	<b>0.99</b>	0.51	0.19	<b>0.13</b>
12	322,426.7	1.15	1.07	<b>0.98</b>	0.58	0.23	<b>0.14</b>
13	2,480,436.7	1.27	1.14	<b>1.02</b>	0.61	0.25	<b>0.15</b>
14	19,583,169.2	1.36	1.17	<b>1.03</b>	0.66	0.27	<b>0.17</b>
15	133,596,114.2	1.61	1.31	<b>1.16</b>	0.74	0.32	<b>0.20</b>
Same Heuristic, Weighted ( $w=1.5$ )							
12	188,177.1	2.00	1.25	<b>1.14</b>	1.51	0.51	<b>0.38</b>
13	398,418.8	2.13	1.32	<b>1.12</b>	1.62	0.53	<b>0.39</b>
14	3,390,387.6	2.31	1.38	<b>1.11</b>	1.62	0.51	<b>0.33</b>
15	6,477,150.7	2.24	1.28	<b>0.98</b>	1.74	0.54	<b>0.37</b>
16	16,848,215.1	2.79	1.49	<b>1.12</b>	1.98	0.55	<b>0.37</b>

Table 8: 15-pancake puzzle.  $r=4$ .

		Signed Error			Absolute Error		
d	IDA*	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
48	2,958,898.5	1.10	1.11	<b>1.05</b>	0.05	0.05	<b>0.03</b>
49	5,894,396.1	1.13	1.15	<b>1.06</b>	0.07	0.07	<b>0.04</b>
50	8,909,564.5	1.16	1.18	<b>1.08</b>	0.09	0.10	<b>0.05</b>
51	15,427,786.9	1.15	1.19	<b>1.07</b>	0.11	0.12	<b>0.07</b>
52	28,308,808.8	1.25	1.28	<b>1.14</b>	0.14	0.17	<b>0.09</b>
53	45,086,452.6	1.23	1.29	<b>1.13</b>	0.16	0.20	<b>0.11</b>
54	85,024,463.5	1.36	1.41	<b>1.22</b>	0.21	0.27	<b>0.15</b>
55	123,478,361.5	1.36	1.45	<b>1.24</b>	0.24	0.31	<b>0.17</b>
56	261,945,964.0	1.44	1.54	<b>1.30</b>	0.28	0.39	<b>0.21</b>
57	218,593,372.3	1.43	1.57	<b>1.32</b>	0.33	0.45	<b>0.26</b>
58	531,577,032.2	1.47	1.64	<b>1.37</b>	0.35	0.51	<b>0.30</b>

Table 9: 15-puzzle. Manhattan Distance.  $r=25$ .

$T_{h,b}$  presents a reasonable signed error, it has in fact a very large absolute error, and once again  $\epsilon$ -truncation produced substantial improvement in prediction accuracy. These prediction results could be improved by increasing the  $r$ -value used. However, we wanted our results to be comparable those in Zahavi *et al.*'s Table 11.

## Conclusion

In this paper we have identified a source of prediction error that has previously been overlooked, namely, that low probability events can degrade predictions in certain circumstances. We call this the discretization effect. This insight led directly to the  $\epsilon$ -truncation method for altering the probability distribution used for making predictions at level  $i$  of the search tree by setting to zero all probabilities smaller than  $\epsilon_i$ , an automatically derived threshold for level  $i$ . Our experimental results showed that more informed type systems for prediction often suffer more from the discretization effect than less informed ones, sometimes leading to the pathological situation that predictions based on the more informed system are actually worse than those based on the less informed system. In our experiments  $\epsilon$ -truncation rarely degraded predictions; in the vast majority of cases it improved predictions, often substantially.

d	IDA*	Signed Error			Absolute Error		
		$T_{h,b}$	$T_c$	$\epsilon-T_c$	$T_{h,b}$	$T_c$	$\epsilon-T_c$
48	193,396.1	0.36	1.50	<b>1.00</b>	232.23	1.13	<b>1.08</b>
49	433,915.3	0.35	<b>1.13</b>	0.82	210.85	1.21	<b>1.11</b>
50	562,708.5	0.55	1.77	<b>1.20</b>	537.97	1.29	<b>1.17</b>
51	965,792.6	0.70	1.39	<b>1.04</b>	812.37	1.32	<b>1.12</b>
52	1,438,694.0	0.96	1.68	<b>1.23</b>	513.99	1.52	<b>1.35</b>
53	2,368,940.3	1.29	1.75	<b>1.32</b>	694.34	1.56	<b>1.26</b>
54	3,749,519.9	1.64	2.03	<b>1.54</b>	647.24	1.77	<b>1.53</b>
55	7,360,297.6	1.90	2.07	<b>1.59</b>	650.59	1.72	<b>1.35</b>
56	12,267,171.0	2.30	2.19	<b>1.61</b>	927.71	2.16	<b>1.86</b>
57	23,517,650.8	2.69	2.29	<b>1.78</b>	600.13	2.02	<b>1.55</b>
58	24,607,970.9	5.26	2.54	<b>1.90</b>	700.29	2.73	<b>2.38</b>

Table 10: 15-puzzle. Inconsistent Heuristic.  $r=1$ .

## Appendix A: Deriving $\epsilon_i$ for $1 \leq i \leq d$

Our aim is to compute a value  $\epsilon_i$ , for  $1 \leq i \leq d$ , below which values of  $\pi(t|u)$  for  $t, u \in T$  are ignored at level  $i$ . To do so, we propose an objective function whose minimizer defines a set of  $\pi$ -values that yield predictions with a minimal absolute expected error for the type allocation vectors.

Often a type at level  $i+1$  can be generated by several different types at level  $i$ . To account for this we define a system of “super types”. A *super type system* groups types at level  $i$  so that two different super types (at level  $i$ ) will never generate children of the same type. In Figure 2, the pairs  $(tree_x, tree_y)$  for  $(x, y) \in \{(1, 3), (1, 4), (2, 3), (2, 4)\}$  would be in one super type.

**Definition 4** Let  $T, T'$  be type systems,  $T \prec T'$ , and  $t' \in T'$ . For all  $i$  and all  $(\hat{s}^*, s^*) \in E$ , the super type  $st(t', i, \hat{s}^*, s^*)$  over  $T$  contains exactly the pairs  $(t_1, t_2) \in T \times T$  for which  $T'(t_2) = t'$  and  $t_1$  occurs at level  $i$  starting the prediction from  $(\hat{s}^*, s^*)$ . The super type system  $ST(i, \hat{s}^*, s^*)$  over  $T$  with respect to  $T'$  is defined by  $ST(i, \hat{s}^*, s^*) = (st(t'_1, i, \hat{s}^*, s^*), \dots, st(t'_z, i, \hat{s}^*, s^*))$  where  $T' = \{t'_1, \dots, t'_z\}$ . We write  $st$  instead of  $st(t', i, \hat{s}^*, s^*)$  whenever  $t', i, \hat{s}^*, s^*$  are clear from context.

In our experiments, for the pancake domain, when  $T = T_{gc}$ , we used  $T' = T_c$ ; for the sliding tile puzzle,  $T = T_{gc,b}$  was used with  $T'$  given by  $T_c$  augmented with the kind of blank location of the parent of the node.

We adapt CDP to super types by estimating, for state pair  $(\hat{s}^*, s^*)$ , level  $i$ , depth limit  $d$ , type  $t$ , and super type  $st \in ST(i, \hat{s}^*, s^*)$ , the probability of generating a node of type  $t$  after seeing super type  $st$  at level  $i$ . We denote the estimate by  $\pi_{\hat{s}^*, s^*}^{i,d}(t|st)$ , defined as

$$\frac{\sum_{\{t_p | (t_p, t) \in st\}} \pi(t|t_p) \beta_{t_p} N(i, t_p, (\hat{s}^*, s^*), d)}{\sum_{\{t_p | (t_p, t) \in st\}} \beta_{t_p} N(i, t_p, (\hat{s}^*, s^*), d)}.$$

The number of nodes  $N(i, st, (\hat{s}^*, s^*), d)$  of a super type  $st \in ST$  at a level  $i$  of prediction is given by

$$N(i, st, (\hat{s}^*, s^*), d) = \sum_{(t_p, t_c) \in st} N(i, t_p, (\hat{s}^*, s^*), d) \beta_{t_p}.$$

We then reformulate the CDP formula equivalently by computing  $N(i, t, (\hat{s}^*, s^*), d)$  as

$$\sum_{st \in ST} N(i-1, st, (\hat{s}^*, s^*), d) \pi_{\hat{s}^*, s^*}^{i-1,d}(t|st) P(t, i, d).$$

Next, we compute a redistribution of  $\pi(t, st)$ -values, for each  $st$  and each level  $i$ , from a solution to the following optimization problem.

Find a type allocation vector  $(a_1, \dots, a_{|T|})$  that minimizes:

$$\lceil N(i, st, (\hat{s}^*, s^*), d) \rceil \sum_{j=0} \sum_{t \in T} |pr(j, t)(a_t - j)| \quad (4)$$

$$s.t.: \sum_{t \in T} a_t = \lceil N(i, st, (\hat{s}^*, s^*), d) \rceil \text{ and } a_t \geq 0 \text{ for } t \in T.$$

Here  $pr(j, t)$  is short for the probability of generating exactly  $j$  children of type  $t$  from  $\lceil N(i, st, (\hat{s}^*, s^*), d) \rceil$  many parents of super type  $st$ . The objective corresponds to replacing, for each super type  $st$ , each  $\pi(t|st)$  by  $a_t / \sum_{t \in T} a_t$  so as to minimize the expected *absolute* error of the type vector allocation.

We deliberately chose to minimize the absolute error rather than MSE here, because we expect small fractions in the type allocation vectors to be harmful, yet favoured by MSE minimization. Consider for example a case of only two types, one being generated with probability  $p$  and the other being generated with probability  $1-p$ , i.e., the expected vector would be  $(p, 1-p)$ . Assume a system predicts a vector  $(q, 1-q)$ . Then its MSE would be  $p((1-q)^2 + (0 - (1-q))^2) + (1-p)((0-q)^2 + (1 - (1-q))^2) = 2q^2 + 2p - 4pq$ , which would be minimized by setting  $q = p$ . Thus MSE minimization would never suggest to ignore a small  $p$  value. The absolute error of the same prediction would be  $p(|1-q| + |0 - (1-q)|) + (1-p)(|0-q| + |1 - (1-q)|) = 2p(1-q) + 2(1-p)q$ . For  $p < 0.5$  this error is minimized at  $q = 0$  while for  $p > 0.5$  error is minimized at  $q = 1$  (and thus  $1-q = 0$ ). In this example we assume that  $N(i, st, (\hat{s}^*, s^*), d) = 1$ . For larger values of  $N(i, st, (\hat{s}^*, s^*), d)$  this error is minimized at  $q = 0$  for smaller  $p$ . Minimizing the absolute error hence encourages the system to ignore types less likely to be generated given an  $N(i, st, (\hat{s}^*, s^*), d)$ -value, as is our intention.

Solving Equation 4 in every step of the prediction would be too time-consuming. Instead, it plays a role in preprocessing, when we use the modified  $\pi$ -values for *super types* to derive  $\epsilon$ -cuts for  $\pi$ -values for *types* as follows. First, we solve a small number of instances of Equation (4) for each level  $i$ . For every value  $\pi(t|st)$  (which is the  $x$ -axis in the example in Figure 3) estimated before solving (4), we compute the fraction of times that this  $\pi(\cdot)$ -value was set to zero in the solutions to (4) for the given level  $i$  (this is the  $y$ -axis in Figure 3). The largest  $\pi(t|st)$  that was set to zero in 50% or more of these instances at level  $i$  is  $\hat{\epsilon}_i$ , a potential cutoff value for level  $i$  (in Figure 3 this is the  $x$ -value where

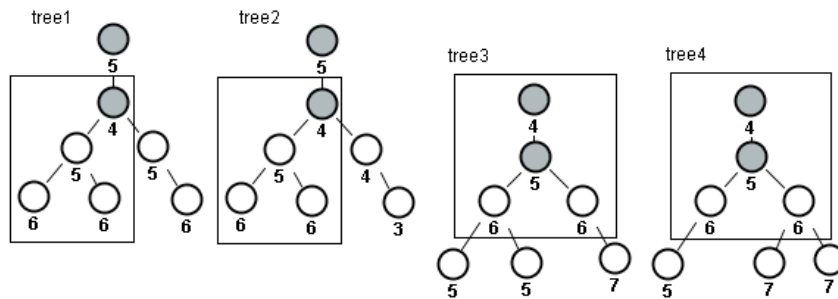


Figure 2: Each of the trees above contains the information necessary to determine the  $T_{gc}$  type of the highlighted pair. Numbers denote  $h$ -values. The left branch of *tree1* and *tree2* can potentially generate children of the same  $T_{gc}$  type, illustrated by *tree3* and *tree4*. Therefore, the pairs  $(tree_x, tree_y)$  for  $(x, y) \in \{(1, 3), (1, 4), (2, 3), (2, 4)\}$  will be in the same super type.

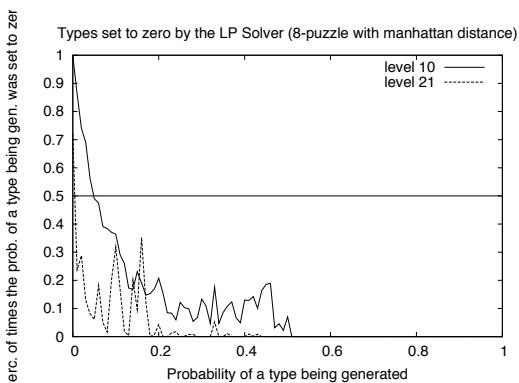


Figure 3:  $\hat{\epsilon}_i$  calculation for  $i = 10$  and  $i = 21$  (8-puzzle with Manhattan distance).

the curve intersects the horizontal  $y=0.5$  line). Note that it can happen that, for some fixed type  $u$ , all  $\pi(t|u)$ -values are below  $\hat{\epsilon}_i$ . Since we must not set  $\pi^i(t|u)$  to zero for all  $t$ , we define  $\epsilon_i$  to be the minimum of  $\hat{\epsilon}_i$  and the largest value  $\delta_i$  such that, for all  $u \in T$  there is some  $t \in T$  such that  $\pi(t|u) > \delta_i$ .

## Appendix B: Results for RMSE

Tables 11 through 17 show results on the same domains, heuristics, and type systems as used in our experimental section above but in terms of *relative root mean squared error* (RRMSE). The relative root mean squared error is calculated by dividing the sum of the RMSE by the sum of the number of nodes expanded by IDA\* for all states in the set of start states. As in the absolute error score, a perfect score for this measure is 0.0. Here we observe the same improvements of  $\epsilon$ -truncation discussed above in terms of signed and absolute errors.

## Acknowledgements

This work was supported by the Laboratory for Computational Discovery at the University of Regina. The authors gratefully acknowledge the research support provided by Alberta Innovates - Technology Futures, the Alberta Ingenuity

d	IDA*	RRMSE		
		$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
17	12.9	0.53	0.44	<b>0.41</b>
18	14.5	0.52	0.48	<b>0.43</b>
19	22.2	0.51	0.47	<b>0.40</b>
20	27.4	0.50	0.50	<b>0.40</b>
21	43.3	0.49	0.50	<b>0.38</b>
22	58.5	0.48	0.53	<b>0.37</b>
23	95.4	0.47	0.54	<b>0.34</b>
24	135.7	0.45	0.58	<b>0.34</b>
25	226.7	0.44	0.61	<b>0.30</b>
26	327.8	0.42	0.65	<b>0.29</b>
27	562.0	0.41	0.68	<b>0.25</b>
28	818.4	0.39	0.73	<b>0.24</b>
29	1,431.7	0.37	0.75	<b>0.20</b>

Table 11: 8-puzzle, Inconsistent Heuristic.  $r=1$ . RRMSE.

Centre for Machine Learning (AICML), and Canada's Natural Sciences and Engineering Research Council (NSERC).

## References

- Culberson, J. C., and Schaeffer, J. 1996. Searching with pattern databases. *Advances in Artificial Intelligence (LNAI 1081)* 402–416.
- Dinh, H. T.; Russell, A.; and Su, Y. 2007. On the value of good advice: The complexity of A\* search with accurate heuristics. In *AAAI-07*, 1140–1145.
- Gaschnig, J. 1979. *Performance Measurement and Analysis of Certain Search Algorithms*. Ph.D. Dissertation, Carnegie-Mellon University.
- Korf, R. E.; Reid, M.; and Edelkamp, S. 2001. Time complexity of Iterative-Deepening-A\*. *Artif. Intell.* 129(1-2):199–218.
- Pearl, J. 1984. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Addison & Wesley.
- Pohl, I. 1977. Practical and theoretical considerations in heuristic search algorithms. *Mach. Intell.* 8:55–72.
- Zahavi, U.; Felner, A.; Burch, N.; and Holte, R. C. 2010. Predicting the performance of IDA\* using conditional distributions. *J. Artif. Intell. Res.* 37:41–83.



		RRMSE		
d	IDA*	$T_h$	$T_{gc}$	$\epsilon-T_{gc}$
Admissible and Consistent Heuristic				
7	1,710.8	0.23	<b>0.02</b>	<b>0.02</b>
8	12,082.4	0.22	<b>0.02</b>	0.03
9	79,650.5	0.21	<b>0.03</b>	<b>0.03</b>
10	507,640.4	0.22	<b>0.03</b>	<b>0.03</b>
11	3,449,158.0	0.22	<b>0.02</b>	0.03
The heuristic above multiplied by 1.5				
7	528.5	0.39	<b>0.07</b>	<b>0.07</b>
8	3,740.3	0.37	<b>0.07</b>	<b>0.07</b>
9	13,924.7	0.38	<b>0.08</b>	<b>0.08</b>
10	64,611.8	0.40	<b>0.09</b>	<b>0.09</b>
11	356,366.1	0.46	<b>0.09</b>	<b>0.09</b>
12	2,397,748.6	0.41	0.10	<b>0.09</b>

Table 12: 10-pancake puzzle.  $r=1$ . RRMSE.

		RRMSE		
d	IDA*	$T_h$	$T_c$	$\epsilon-T_c$
Admissible and Consistent Heuristic (PDB)				
9	1,035.2	0.17	0.08	<b>0.07</b>
10	5,547.9	0.17	0.10	<b>0.08</b>
11	46,009.3	0.19	0.11	<b>0.09</b>
12	322,426.7	0.23	0.14	<b>0.12</b>
13	2,480,436.7	0.31	0.19	<b>0.14</b>
14	19,583,169.2	0.39	0.23	<b>0.17</b>
15	133,596,114.2	0.62	0.34	<b>0.23</b>
Same Heuristic, Weighted ( $w=1.5$ )				
12	188,177.1	1.06	0.36	<b>0.27</b>
13	398,418.8	1.14	0.37	<b>0.22</b>
14	3,390,387.6	1.31	0.43	<b>0.27</b>
15	6,477,150.7	1.27	0.44	<b>0.32</b>
16	16,848,215.1	1.82	0.63	<b>0.41</b>

Table 15: 15-pancake puzzle.  $r=4$ . RRMSE.

		RRMSE		
d	IDA*	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
18	134.4	0.04	<b>0.02</b>	<b>0.02</b>
19	238.4	0.05	<b>0.02</b>	<b>0.02</b>
20	360.1	0.06	<b>0.03</b>	<b>0.03</b>
21	630.7	0.06	<b>0.03</b>	<b>0.03</b>
22	950.6	0.07	0.05	<b>0.04</b>
23	1,649.5	0.08	0.06	<b>0.04</b>
24	2,457.5	0.09	0.08	<b>0.05</b>
25	4,245.5	0.09	0.08	<b>0.05</b>
26	6,294.4	0.10	0.11	<b>0.06</b>
27	10,994.9	0.10	0.12	<b>0.06</b>

Table 13: 8-puzzle, Manhattan Distance.  $r=10$ . RRMSE.

		RRMSE		
d	IDA*	$T_{h,b}$	$T_{gc}$	$\epsilon-T_{gc}$
48	2,958,898.5	0.12	0.11	<b>0.06</b>
49	5,894,396.1	0.14	0.15	<b>0.07</b>
50	8,909,564.5	0.17	0.18	<b>0.09</b>
51	15,427,786.9	0.16	0.19	<b>0.09</b>
52	28,308,808.8	0.26	0.28	<b>0.15</b>
53	45,086,452.6	0.24	0.29	<b>0.14</b>
54	85,024,463.5	0.37	0.41	<b>0.23</b>
55	123,478,361.5	0.36	0.45	<b>0.25</b>
56	261,945,964.0	0.44	0.54	<b>0.30</b>
57	218,593,372.3	0.44	0.57	<b>0.32</b>
58	531,577,032.2	0.47	0.64	<b>0.37</b>

Table 16: 15-puzzle. Manhattan Distance.  $r=25$ . RRMSE.

		RRMSE		
d	IDA*	$T_{h,b}$	$T_{gc,b}$	$\epsilon-T_{gc,b}$
Manhattan Distance				
33	30,461.9	0.03	<b>0.02</b>	<b>0.02</b>
34	49,576.8	0.04	0.03	<b>0.02</b>
35	80,688.2	0.05	0.05	<b>0.03</b>
36	127,733.4	0.07	0.06	<b>0.03</b>
37	201,822.7	0.09	0.09	<b>0.05</b>
38	327,835.3	0.11	0.12	<b>0.05</b>
39	478,092.5	0.14	0.16	<b>0.07</b>
40	822,055.4	0.18	0.21	<b>0.09</b>
41	1,163,312.1	0.22	0.27	<b>0.11</b>
42	1,843,732.2	0.28	0.34	<b>0.14</b>
Weighted Manhattan Distance ( $w=1.5$ )				
33	926.2	0.06	0.02	<b>0.01</b>
34	1,286.9	0.07	0.03	<b>0.02</b>
35	2,225.6	0.10	0.04	<b>0.02</b>
36	2,670.7	0.12	0.05	<b>0.03</b>
37	3,519.5	0.15	0.07	<b>0.04</b>
38	5,570.8	0.20	0.10	<b>0.05</b>
39	6,983.8	0.24	0.14	<b>0.06</b>
40	9,103.3	0.30	0.19	<b>0.07</b>
41	13,635.3	0.37	0.25	<b>0.08</b>
42	16,634.2	0.44	0.31	<b>0.09</b>

Table 14: (3x4)-puzzle.  $r=20$ . RRMSE.

		RRMSE		
d	IDA*	$T_{h,b}$	$T_c$	$\epsilon-T_c$
48	193,396.1	1.01	0.69	<b>0.44</b>
49	433,915.3	1.01	<b>0.46</b>	0.48
50	562,708.5	1.04	0.89	<b>0.49</b>
51	965,792.6	1.09	0.62	<b>0.50</b>
52	1,438,694.0	1.20	0.81	<b>0.55</b>
53	2,368,940.3	1.33	0.86	<b>0.57</b>
54	3,749,519.9	1.63	1.08	<b>0.71</b>
55	7,360,297.6	1.68	1.14	<b>0.74</b>
56	12,267,171.0	2.17	1.23	<b>0.72</b>
57	23,517,650.8	2.44	1.31	<b>0.83</b>
58	24,607,970.9	4.73	1.57	<b>0.98</b>

Table 17: 15-puzzle. Inconsistent Heuristic.  $r=1$ . RRMSE.