

Document downloaded from:

<http://hdl.handle.net/10251/64413>

This paper must be cited as:

Hernández Orallo, J. (2015). C-tests revisited: back and forth with complexity. En Artificial General Intelligence 8th International Conference, AGI 2015, AGI 2015, Berlin, Germany, July 22-25, 2015, Proceedings. Springer International Publishing. 272-282. doi:10.1007/978-3-319-21365-1\_28.



The final publication is available at

[http://link.springer.com/chapter/10.1007/978-3-319-21365-1\\_28](http://link.springer.com/chapter/10.1007/978-3-319-21365-1_28)

Copyright Springer International Publishing

Additional Information

The final publication is available at Springer via [http://dx.doi.org/10.1007/978-3-319-21365-1\\_28](http://dx.doi.org/10.1007/978-3-319-21365-1_28)

# C-tests revisited: back and forth with complexity

José Hernández-Orallo

DSIC, Universitat Politècnica de València, Spain  
jorallo@dsic.upv.es

**Abstract.** We explore the aggregation of tasks by weighting them using a difficulty function that depends on the complexity of the (acceptable) policy for the task (instead of a universal distribution over tasks or an adaptive test). The resulting aggregations and decompositions are (now retrospectively) seen as the natural (and trivial) interactive generalisation of the  $C$ -tests.

**Keywords:** Intelligence evaluation, artificial intelligence,  $C$ -tests, algorithmic information theory, universal psychometrics, agent response curve.

## 1 Introduction

A first test using AIT was the  $C$ -test [32,16], where the goal was to find a continuation of a sequence of letters, as in some IQ tasks, and in the spirit of Solomonoff’s inductive inference problems: “given an initial segment of a sequence, predict its continuation” (as quoted in [47, p.332]). Levin’s  $Kt$  complexity was used to calculate the difficulty of a sequence of letters. The performance was measured as an aggregated value over a range of difficulties:

$$I(\pi) \triangleq \sum_{h=1}^H h^e \sum_{i=1}^N \frac{1}{N} \text{Hit}(\pi, x_{i,h}) \quad (1)$$

where  $\pi$  is the subject, the difficulties range from  $h = 1$  to  $H$  and there are  $N$  sequences  $x_{i,k}$  per difficulty  $h$ . The function  $\text{hit}$  returns 1 if  $\pi$  is right with the continuation and 0 otherwise. If  $e = 0$  we have that all difficulties have the same weight. The  $N$  sequences per difficulty were chosen (uniformly) randomly.

This contrasts with a more common evaluation in artificial intelligence based on average-case performance according to a probability of problems or tasks:

$$\Psi(\pi) \triangleq \sum_{\mu \in M} p(\mu) \cdot \mathbb{E}[R(\pi, \mu)] \quad (2)$$

where  $p$  is a probability distribution on the set of tasks  $M$ , and  $R$  is a result function of agent  $\pi$  on task  $\mu$ . Actually, eq. 2 can also be combined with AIT, in a different way, by using a universal distribution [54,47], i.e.,  $p(\mu) = 2^{-K(\mu)}$ , where  $K(\mu)$  is the Kolmogorov complexity of  $\mu$ , as first chosen by [44].

The work in [44] has been considered a generalisation of [32,16], from static sequences (predicting a continuation of a sequence correctly) to dynamic environments (maximising rewards in an MDP). In this paper we challenge this

interpretation and look for a proper generalisation of [32,16] using the notion of difficulty in the outer sum, as originally conceived and seen in eq. 1. The key idea is the realisation that for the  $C$ -test the task and the solution were the same thing. This meant that the difficulty was calculated as the size of the simplest programs that generates the sequence, which is both the task and the solution. Even if the complexity of the task and the solution coincide here, it is *the complexity of the solution what determines the difficulty of the problem*.

However, when we move from sequences to environments or other kind of interactive tasks, the complexity of the policy that solves the task and the complexity of the environment are no longer the same. In fact, this is discussed in [24,22]: the complexity of the environment is roughly an upper bound of the complexity of the acceptable policies (any agent that reach an acceptable performance value), but very complex environments can have very simple acceptable policies. In fact, the choice of  $p(\mu) = 2^{-K(\mu)}$  has been criticised for giving too much weight to a few environments. Also, it is important to note that the invariance theorem is more meaningful for Kolmogorov Complexity than for Algorithmic Probability, as for the former it gives some stability for values of  $K$  that are not very small, but for a probability it is precisely the small cases that determine most of the distribution mass. In fact, any computable distribution can be approximated (to whatever required precision) by a universal distribution. This means that the choice of  $p(\mu) = 2^{-K(\mu)}$  for Eq. 2 is actually a metadefinition, which leads to virtually any performance measure, depending on the UTM that is chosen as reference.

By decoupling the complexity of task and policy we can go back to eq. 1 and work out a notion of difficulty of environments that depends on the complexity of the policy. While this may look retrospectively trivial, and the natural extension in hindsight, we need to solve and clarify some issues, and properly analyse the relation of the two different philosophies given by eq. 2 and eq. 1.

Section 2 discusses some previous work, introduces some notation and recovers the difficulty-based decomposition of aggregated performance. Section 3 introduces several properties about difficulty functions and the view of difficulty as policy complexity. Section 4 discusses the choices for the difficulty-dependent probability. Section 5 briefly deals with the role of computational steps for difficulty. Section 6 closes the paper with a discussion.

## 2 Background

AI evaluation has been performed in many different ways (for a recent account of AI evaluation, see [21]), but a common approach is based on averaging performance on a range of tasks, as in eq. 2.

In what follows, we will focus on the approaches that are based on AIT. As mentioned above, the first intelligence test using AIT was the so-called  $C$ -test [32,16]. Figure 1 shows examples of sequences that appear in this test. The difficulty of each sequence was calculated as Levin's  $Kt$ , a time-weighted version of Kolmogorov complexity  $K$ . Some preliminary experimental results

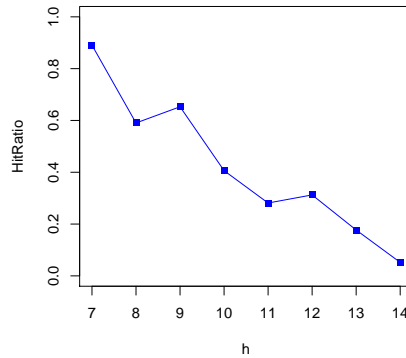
$h = 9$  : a, d, g, j, ...                      Answer: m  
 $h = 12$  : a, a, z, c, y, e, x, ...              Answer: g  
 $h = 14$  : c, a, b, d, b, c, c, e, c, d, ...      Answer: d

**Fig. 1.** Several series of different difficulties 9, 12, and 14 used in the  $C$ -test [16].

showed that human performance correlated with the absolute difficulty ( $h$ ) of each exercise and also with IQ test results for the same subjects. Figure 2 shows the results (taken from [32,16]). HitRatio is defined as the inner sum of eq. 1:

$$\text{HitRatio}(\pi, h) \triangleq \sum_{i=1}^N \frac{1}{N} \text{Hit}(\pi, x_{i,h}) \quad (3)$$

An interesting observation is that by arranging problems by difficulty we see that HitRatio seems to be very small from difficulty 8 on. This makes the estimation of the measure much easier, as we only need to focus on (the area of) a small interval of difficulties. In fact, this use of difficulty is common in psychometrics.



**Fig. 2.** Results obtained by humans on task of different difficulty in the  $C$ -test [16].

Several generalisations of the  $C$ -test were suggested (for “cognitive agents [...] with input/output devices for a complex environment” [32] where “rewards and penalties could be used instead” [18]) or extending them for other cognitive abilities [17], but not fully developed.

AIT and reinforcement learning were finally combined in [44], where all possible environments were considered in eq. 2, instantiated with a universal distribution for  $p$ , i.e.,  $p(\mu) = 2^{-K(\mu)}$ , with  $K(\mu)$  being the Kolmogorov complexity of each environment  $\mu$ . Some problems (computability, discriminating power, overweight for small environments, time, ...) were discussed with the aim of making a more applicable version of this approach by [33] and [24, secs. 3.3 and 4].

While the aim of all these proposals was to measure intelligence, many interesting things can happen if AIT is applied to cognitive abilities other than intelligence, as suggested in [17] for the passive case and hinted in [24, secs. 6.5 and 7.2] for the dynamic cases, which proposes the use of different kinds of

videogames as environments (two of the most recently introduced benchmarks and competitions in AI are in this direction [4,51]).

We consider tests that are composed of tasks (also called environments or items) and are performed by agents (also called policies or subjects). The set of tasks is denoted by  $M$ . Its elements are usually denoted by  $\mu$ . The set of agents is denoted by  $\Pi$ . Its elements are usually denoted by  $\pi$ . The length in bits of a string is denoted by  $L(x)$ . We can think of a proper encoding of tasks and agents as strings. Given a UTM  $U$  we define Kolmogorov complexity as  $K_U(y) = \min_x : U(x)=y L(x)$ . We will usually drop the subindex  $U$ . Finally, we have the expected value of the response, score or result of  $\pi$  in  $\mu$  for a time limit  $\tau$  as  $\mathbb{E}[R(\pi, \mu, \tau)]$ . The value of  $\tau$  will be usually omitted. The  $R$  function always gives values between 0 and 1 and we assume it is always defined (a value of  $R = 0$  is assumed for non-halting situations) We also define  $\mathbb{E}[LS(\pi, \mu, \tau)] \triangleq L(\pi) + \log \mathbb{E}[S(\pi, \mu, \tau)]$ . Logarithms are always binary.

It is actually in [28], where we can find a first connection between the schemas of eq. 1 and eq. 2. We adapt definition 14 in [28], which is a generalisation of eq. 2, by making the set  $M$  and the task probability  $p$  explicit as a parameters.

**Definition 1.** *The expected average result for a task class  $M$ , a distribution  $p$  and an agent  $\pi$  is:*

$$\Psi(\pi, M, p) \triangleq \sum_{\mu \in M} p(\mu) \cdot \mathbb{E}[R(\pi, \mu)] \quad (4)$$

And now we use proposition 4 in [28] that decomposes it. First, we define partial results for a given difficulty  $h$  as follows:

$$\Psi_h(\pi, M, p) \triangleq \sum_{\mu \in M, \tilde{h}(\mu)=h} p(\mu|h) \cdot \mathbb{E}[R(\pi, \mu)] \quad (5)$$

Where  $\tilde{h}$  is a difficulty function  $\tilde{h} : M \rightarrow \mathbb{R}^+ \cup 0$ . Note that this parametrises the result of eq. 4 for different difficulties. For instance, for two agents  $\pi_A$  and  $\pi_B$  we might have that  $\Psi_3(\pi_A) < \Psi_3(\pi_B)$  but  $\Psi_7(\pi_A) > \Psi_7(\pi_B)$ . If we represent  $\Psi_h(\pi, M, p)$  on the  $y$ -axis versus  $h$  on the  $x$ -axis we have a so-called agent response curve, much like Fig. 2.

If we want to get a single number from an agent response curve we can aggregate performance for a range of difficulties, e.g., as follows:

**Proposition 1.** ([28, prop. 4]) *The expected average result  $\Psi(\pi, M, p)$  can be rewritten as follows: in the particular case when  $\tilde{h}$  only gives discrete values:*

$$\Psi(\pi, M, p) = \sum_{h=0}^{\infty} p(h) \Psi_h(\pi, M, p) \quad (6)$$

where  $p(h)$  is a discrete probability function for eq. 6. Note that equations 4, 5 and 6 are generalisations, respectively, of equations 2, 3 and 1.

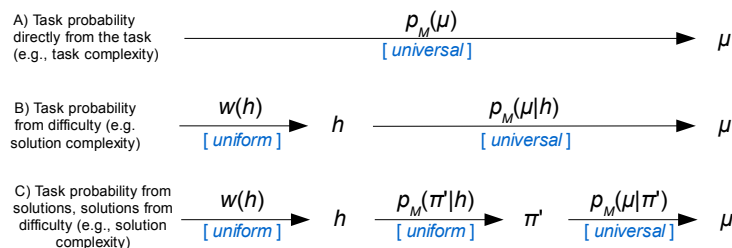
### 3 Difficulty functions

Before setting an appropriate measure of difficulty based on the policy, in this section we will analyse which properties a difficulty function may have.

The decomposition in previous section suggests that we could try to fix a proper measure of difficulty first and then think about a meaningful distribution  $p(h)$ . Once this is settled, we could try to find a distribution for all environments of that difficulty  $p(\mu|h)$ . In other words, once we determine how relevant a difficulty is we ask which tasks to take for that difficulty. This is the spirit of the  $C$ -test [32,16] as seen in eq. 1. In fact, we perhaps we do not need a  $p(h)$  that decays dramatically, as it is expectable to see performance to decrease for increasing difficulty, as in Figure 2.

To distinguish  $p(h)$  and  $p(\mu|h)$  we will denote the former with  $w$  and the latter with  $p_M$ . We will use any distribution or even a measure (not summing up to one, for reasons that we will see later on) as a subscript for  $\Psi$ . For instance, we will use the following notation  $\Psi_{\mathcal{U}(h_{min}, h_{max})}(\pi, M, p_M)$ , where  $\mathcal{U}(a, b)$  represents a uniform distribution between  $a$  and  $b$ . For instance, we can have two agents  $\pi_A$  and  $\pi_B$  such that  $\Psi_{\mathcal{U}(1,10)}(\pi_A) > \Psi_{\mathcal{U}(1,10)}(\pi_B)$  but  $\Psi_{\mathcal{U}(11,20)}(\pi_A) < \Psi_{\mathcal{U}(11,20)}(\pi_B)$ . We will use the notation  $\Psi_{\oplus}(\pi, M, p_M)$  when  $w(h) = 1$  (note that this is not the uniform distribution for discrete  $h$ ), which means that the partial aggregations for each difficulty are just added. In other words,  $\Psi_{\oplus}(\pi, M, p_M) \triangleq \sum_{h=0}^{\infty} \Psi_h(\pi, M, p_M)$  for discrete difficulty functions. We will explore whether this (area under the agent response curve) is bounded.

Figure 3 shows approach A, which has already been mentioned, while approaches B and C will be seen in sections 4.1 and 4.2 respectively.



**Fig. 3.** Three approaches to aggregate the results for a set of tasks. Top (A) shows the classical approach of choosing a probability for the task, according to the properties of the task. Middle (B) shows the approach where we arrange tasks by difficulty, and the notion of difficulty is derived from the properties of the policy. Bottom (C) shows a variation of B where we derive acceptable policies for a given difficulty and then generate tasks for each policy. Between square brackets some choices we examine in this paper.

When we aggregate environments with different scales on  $R$  and different difficulties, we may have that an agent focusses on a few environments with high difficulty while another focusses on many more environments with small responses. Agent response curves in [28], which are inspired by item response

curves in psychometrics (but inverting the view between agents and items), allow us to see how each agent performs for different degrees of difficulty. Looking at Figure 2 and similar agent response curves in psychometrics, we see that the notion of difficulty must be linked to  $R$ , i.e., how well the agents perform, and not about the complexity of the task, as in the previous section.

Another option is what is done in [22], as  $\bar{h}(\mu) \triangleq \min_{\pi: \mathbb{E}[R(\pi, \mu)] = R_{max}(\mu)} L(\pi)$  where  $R_{max}(\mu) = \max_{\pi} \mathbb{E}[R(\pi, \mu)]$ . However,  $R_{max}$  may be hard to calculate and even if it can be effectively calculated, any minor mistake or inefficiency in a very good agent will prevent the agent from reaching the optimal result, leading to a notion of difficulty linked to the complexity of the ‘perfect’ policy. In [22], a ‘tolerance value’ is considered and, instead of one policy, difficulty is linked to the probability of finding a policy under this tolerance.

We are going to consider this tolerance  $\epsilon$  of acceptability.

$$\mathbb{A}^{[\epsilon]}(\pi, \mu) \triangleq \mathbf{1}(\mathbb{E}[R(\pi, \mu)] \geq 1 - \epsilon) \quad (7)$$

This returns 1 if the expected response is above  $1 - \epsilon$  and 0 otherwise. If  $\mathbb{A}^{[\epsilon]}(\pi, \mu) = 1$  we say that  $\pi$  is  $\epsilon$ -acceptable. With this, we binarise responses. One can argue that we could have just defined a binary  $R$ , but it is important to clarify that it is not the same to have tolerance for each single  $R$  (or a binarised  $R$ ) than to have a tolerance for the expected value  $\mathbb{E}[R]$ . The tolerance on the expected value allows the agent to have variability in their results (e.g., stochastic agents) provided the expected value is higher than the tolerance. Finally, even if we will be very explicit about the value of  $\epsilon$ , and changing it will change the difficulty value of any environment, it is important to say that this value is not so relevant. The reason is that for any environment we can build any other environment where the responses are transformed by any function. In fact, we could actually consider one fixed threshold, such as 0.5, always.

And now we can just define a new version of eq. 5 using this new function:

$$\Psi_h^{[\epsilon]}(\pi, M, p_M) \triangleq \sum_{\mu \in M, \bar{h}(\mu) = h} p_M(\mu|h) \cdot \mathbb{A}^{[\epsilon]}(\pi, \mu) \quad (8)$$

We can just rewrite equations 6 accordingly:

$$\Psi_w^{[\epsilon]}(\pi, M, p_M) = \sum_{h=0}^{\infty} w(h) \Psi_h^{[\epsilon]}(\pi, M, p_M) \quad (9)$$

Given the above, we are now ready for a few properties about difficulty functions.

**Definition 2.** *A difficulty function  $\bar{h}$  is strongly bounded in  $M$  if for every  $\pi$  there is a difficulty  $h$  such that for every  $\mu \in M : \bar{h}(\mu) \geq h$  we have  $\mathbb{A}^{[\epsilon]}(\pi, \mu) = 0$ .*

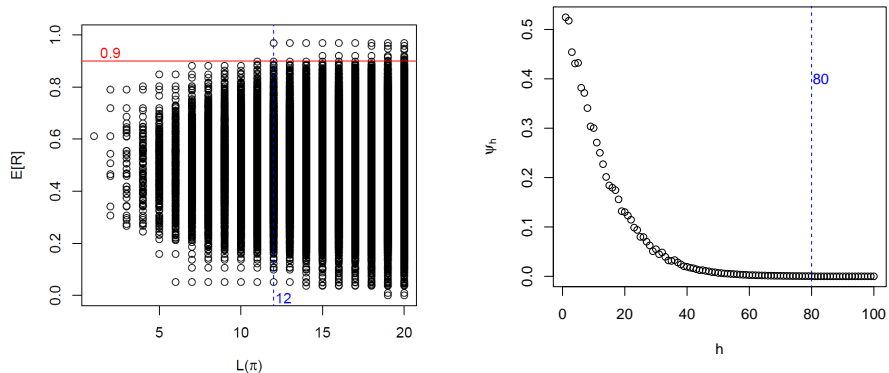
Now we choose the difficulty function in terms of  $\epsilon$ -acceptability, i.e.:

$$\bar{h}^{[\epsilon]}(\mu) \triangleq \min\{L(\pi) : \mathbb{E}[R(\pi, \mu)] \geq 1 - \epsilon\} = \min\{L(\pi) : \mathbb{A}^{[\epsilon]}(\pi, \mu) = 1\} \quad (10)$$

We can say a few words about the cases where a truly random agent gives an acceptable policy for an environment. If this is the case, we intuitively consider

the environment easy. So, in terms, of  $L$ , we consider truly random agents to be simple, which is more reasonable than considering them of infinite difficulty, and goes well with our consideration of stochastic agents and environments.

Figure 4 (left) shows the distribution of response according to  $L(\pi)$ , but setting  $\epsilon = 0.9$ . We see that the simplest  $\epsilon$ -acceptable policy has  $L = 12$ .



**Fig. 4.** Left: an illustrative distribution of responses of a population of agents for a single environment. If we set the threshold at  $0.9 = 1 - \epsilon$ , the simplest policy above this threshold is of ‘complexity’  $h = 12$ . Right: an illustrative distribution of the result of  $\Psi_h$  considering a population of environments for a single agent, an *agent response curve*. There is no task  $\pi$  of difficulty above 80 for which  $\mathbb{E}[R(\pi, \mu)] \geq 1 - \epsilon$ , i.e., there is no task for which  $\pi$  is  $\epsilon$ -acceptable, so  $\Psi_h$  is 0 from 80 on. If we were using the definition of  $h$  as for Eq. 10, this 80 would be  $L(\pi)$ . Also note on the right plot that all ‘heaven’ tasks (good results independently of what the agent does) are at  $h = 1$ , while all ‘hell’ tasks (bad response independently of what the agent does) are at  $h = \infty$ .

With the difficulty function in eq. 10 we have:

**Proposition 2.** *The difficulty function  $h^{[\epsilon]}$  in eq. 10 is strongly bounded.*

*Proof.* For every policy  $\pi$ , if a task  $\mu$  has a difficulty  $h^{[\epsilon]}(\mu) > L(\pi)$ , it means that  $\pi$  is not  $\epsilon$ -acceptable, because otherwise the difficulty would be  $L(\pi)$  and not  $h$ . Consequently,  $\mathbb{A}^{[\epsilon]}(\pi, \mu) = 0$  for all  $\mu$  of difficulty  $h^{[\epsilon]}(\mu) > L(\pi)$ . It is sufficient to take  $h > L(\pi)$  for every  $\pi$  to see that  $h$  is strongly bounded.

This is what we see in Fig. 4 (right), where  $L(\pi) = 80$ . With  $h^{[\epsilon]}$  in eq. 10, we can ensure that the values are going to be 0 from  $h = 80$  on.

This may not be the case for other difficulty functions. We can imagine a situation where the curve never converges to zero. For instance, if the difficulty function is decoupled from resources (length and/or steps) of the acceptable policies or we do not use the notion of  $\epsilon$ -acceptability then we cannot avoid that a very simple policy could eventually score well in a problem with very high difficulty. This would be counter-intuitive, as if there is a simple policy for a difficult problem, the latter should not be considered difficult any more.



## 4 Difficulty-conditional task probabilities

In the previous sections we have focussed on  $w(h)$  and whether it is necessary or not. We have seen difficulty functions where just aggregating  $\Psi_h$  without  $w(h)$  (or  $w(h) = 1$ ) leads to a  $\Psi_{\oplus}(\pi, M, p_M)$  that is bounded. The question now is how to choose the conditional probability  $p_M(\mu|h)$ . In the  $C$ -test, eq. 1, this was chosen as a uniform distribution. However, this is not possible in an interactive scenario if we consider all possible tasks, as the number of tasks for which there is an acceptable policy  $\pi$  of  $L(\pi) = n$  can be infinite. Even if we cannot set a uniform distribution, we want a choice of  $p_M(\mu|h)$  that keeps the task diversity (unless there is any special bias to choose the tasks).

### 4.1 Task probability depends on difficulty

The first thing we can do is to assume  $p(\mu|h)$  in eq. 8 as  $p(\mu|h) = \frac{2^{-K(\mu)}}{\nu(h)}$  if  $\hbar^{[\epsilon]}(\mu) = h$  and 0 otherwise, where  $\nu(h)$  is a normalisation term to make the mass of the distribution equal to 1, which can be formulated as  $\nu(h) = \sum_{\mu: \hbar^{[\epsilon]}(\mu)=h} 2^{-K(\mu)}$ .

And now we have:

$$\Psi_h^{[\epsilon]}(\pi, M, p_M) = \sum_{\mu \in M, \hbar^{[\epsilon]}(\mu)=h} p_M(\mu|h) \cdot \mathbb{A}^{[\epsilon]}(\pi, \mu) = \frac{1}{\nu(h)} \sum_{\mu \in M, \hbar^{[\epsilon]}(\mu)=h} 2^{-K(\mu)} \cdot \mathbb{A}^{[\epsilon]}(\pi, \mu)$$

From here, we can plug it into eq. 9 for the discrete case:

$$\Psi_w^{[\epsilon]}(\pi, M, p_M) = \sum_{h=0}^{\infty} w(h) \frac{1}{\nu(h)} \sum_{\mu \in M, \hbar^{[\epsilon]}(\mu)=h} 2^{-K(\mu)} \cdot \mathbb{A}^{[\epsilon]}(\pi, \mu) \quad (11)$$

Note that the above is going to be bounded independently of the difficulty function if  $w$  is a probability distribution. Also notice that  $\frac{1}{\nu(h)}$  is on the outer sum, and that  $\nu(h)$  is lower than 1, so the normalisation term is actually greater than 1.

And if we use any of the difficulty functions in equations 10 we can choose  $w(h) = 1$  and  $\Psi_{\oplus}^{[\epsilon]}(\pi, M, p_M)$  is bounded.

### 4.2 Task probability depends on the policy probability

One of things of the use of equation 10 is that the number of acceptable policies per difficulty is finite. This is what happened in the  $C$ -test and that is the reason why a uniform distribution could be used for the inner sum. We could try to decompose the inner sum by using the policy and get the probability of the task given the policy.

The interpretation would be as follows: for each difficulty value we aggregate all the acceptable policies with size equal to that difficulty uniformly and for each of these policies all the environments where each policy is acceptable with a universal distribution. This extra complication with respect to eq. 11 can only be justified if we generate environments and agents and we check them as we populate *Pairs*, as a way of constructing a test more easily.

## 5 Using computational steps

As we mentioned in the introduction, the  $C$ -test [32,16] used Levin’s  $Kt$  instead of  $K$ . Apart from being computable,  $Kt$  is related to Levin’s search. This makes its choice much more appropriate for a measure of difficulty. However, when working with interactive tasks and with stochastic tasks and agents, we have that the number of computational steps must be calculated as  $\mathbb{E}[S(\pi, \mu)]$ . For simplicity, up to this point, we have considered just  $K$ . Now we are briefly exploring the inclusion of the computational steps. In section 2 we defined  $LS$ . Now we define a version that accounts for the tolerance  $\epsilon$  as follows:

$$LS^{[\epsilon]}(\pi, \mu) \triangleq \mathbb{E}[LS(\pi, \mu)] \text{ if } \mathbb{A}^{[\epsilon]}(\pi, \mu) = 1 \text{ and } \infty \text{ otherwise}$$

and we define a new difficulty function that considers computational steps:

$$h^{[\epsilon]}(\mu) \triangleq \min_{\pi} LS^{[\epsilon]}(\pi, \mu)$$

This difficulty function is not bounded, as  $LS$  depends on  $\mu$ , and we can always find a very short policy that takes an enormous amount of steps for a task with very high difficulty. This is an acceptable policy, but does not reduce the difficulty of the task, so it can always score non-zero beyond any limit.

This means that for this difficulty function we would need to use equation eq. 9 with an appropriate  $w(h)$  (e.g., a small decay or a uniform interval of difficulties we are interested in).

If the testing procedure established a limit on the number of steps (total or per transition) we would have this new difficulty function would be strongly bounded. Alternatively, we could reconsider the inclusion the computational steps in the notion of acceptability. In this case, the approach in section 4.2 could not be used, as the probability of  $\pi$  given  $h$  would also depend on  $\mu$ .

## 6 Discussion

We have gone from eq. 1 taken from  $C$ -test to eq. 9 when using discrete difficulty functions. We have seen that difficulties allow for a more detailed analysis of what happens for a given agent, depending on whether it succeeds at easy or difficult tasks. For some difficulty functions, we do not even need to determine the weight for each difficulty and just calculate the area, as an aggregated performance for all difficulties. Actually, given the characteristics of most of the difficulty functions we have seen, setting a maximum difficulty could be enough, for instance,  $h = 200$ , more than convenient, as it is really challenging to look into a space of policies of 200 bits (without prior knowledge). So, a finite range of difficulties would be enough for the evaluation of feasible agents.

The important thing is that now we do not need to set an a priori distribution for all tasks  $p(\mu)$ , but just a conditional distribution  $p(\mu|h)$ . Note that if we set a high  $h$  we have the freedom to find the simplest task that creates that difficulty. In fact, this is desirable, as when we try to create a difficult game with a small

number of rules. Actually, the choice of  $p(\mu|h)$  as a universal distribution still depends on the reference machine and can set most of the probability mass on smaller tasks, but as it is conditional on  $h$ , all trivial, dead or simply meaningless tasks have usually very extreme values of  $h$  (very low or infinite). That means that there is a range of *interesting* difficulties, discarding very small values of  $h$  and very large values of  $h$ . Figure 2 is a nice example of this, where only difficulties between 1 and 8 were used, and we see also that  $h = 1$  and  $h > 7$  are not really very discriminating. The bulk of the testing effort must be performed in this range.

## References

1. M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 06 2013.
2. J. Hernández-Orallo. Beyond the Turing Test. *J. Logic, Language & Information*, 9(4):447–466, 2000.
3. J. Hernández-Orallo. Computational measures of information gain and reinforcement in inference processes. *AI Communications*, 13(1):49–50, 2000.
4. J. Hernández-Orallo. On the computational measurement of intelligence factors. In A. Meystel, editor, *Performance metrics for intelligent systems workshop*, pages 1–8. National Institute of Standards and Technology, Gaithersburg, MD, 2000.
5. J. Hernández-Orallo. AI evaluation: past, present and future. *arXiv preprint arXiv:1408.6908*, 2014.
6. J. Hernández-Orallo. On environment difficulty and discriminating power. *Autonomous Agents and Multi-Agent Systems*, pages 1–53, 2014.
7. J. Hernández-Orallo and D. L. Dowe. Measuring universal intelligence: Towards an anytime intelligence test. *Artificial Intelligence*, 174(18):1508–1539, 2010.
8. J. Hernández-Orallo, D. L. Dowe, and M. V. Hernández-Lloreda. Universal psychometrics: Measuring cognitive abilities in the machine kingdom. *Cognitive Systems Research*, 27:5074, 2014.
9. J. Hernández-Orallo and N. Minaya-Collado. A formal definition of intelligence based on an intensional variant of Kolmogorov complexity. In *Proc. Intl Symposium of Engineering of Intelligent Systems (EIS'98)*, pages 146–163. ICSC Press, 1998.
10. B. Hibbard. Bias and no free lunch in formal measures of intelligence. *Journal of Artificial General Intelligence*, 1(1):54–61, 2009.
11. S. Legg and M. Hutter. Universal intelligence: A definition of machine intelligence. *Minds and Machines*, 17(4):391–444, 2007.
12. M. Li and P. Vitányi. *An introduction to Kolmogorov complexity and its applications (3rd ed.)*. Springer-Verlag, 2008.
13. T. Schaul. An extensible description language for video games. *Computational Intelligence and AI in Games, IEEE Transactions on*, PP(99):1–1, 2014.
14. R. J. Solomonoff. A formal theory of inductive inference. Part I. *Information and control*, 7(1):1–22, 1964.