

Accounting for the Neglected Dimensions of AI Progress

Fernando Martínez-Plumed

DSIC, Universitat Politècnica de València, Spain

fmartinez@dsic.upv.es

Shahar Avin

CSER, University of Cambridge, UK

sa478@cam.ac.uk

Miles Brundage

FHI, University of Oxford, UK

miles.brundage@philosophy.ox.ac.uk

Allan Dafoe

FHI, University of Oxford, UK

allan.dafoe@governance.ai

Seán Ó hÉigearthaigh

CSER, University of Cambridge, UK

so348@cam.ac.uk

José Hernández-Orallo

DSIC, Universitat Politècnica de València, Spain

jorallo@dsic.upv.es

June 5, 2018

Abstract

We analyze and reframe AI progress. In addition to the prevailing metrics of performance, we highlight the usually neglected costs paid in the development and deployment of a system, including: data, expert knowledge, human oversight, software resources, computing cycles, hardware and network facilities, development time, etc. These costs are paid throughout the life cycle of an AI system, fall differentially on different individuals, and vary in magnitude depending on the replicability and generality of the AI solution. The multidimensional performance and cost space can be collapsed to a single utility metric for a user with transitive and complete preferences. Even absent a single utility function, AI advances can be generically assessed by whether they expand the Pareto (optimal) surface. We explore a subset of these neglected dimensions using the two case studies of Alpha* and ALE. This broadened conception of progress in AI should lead to novel ways of measuring success in AI, and can help set milestones for future progress.

1 Introduction

Metrics of scientific progress can play an outsized role in the perception of a field and in the allocation of its resources. By contrast, that which goes unmeasured is often neglected. We argue for a more general accounting of progress in AI, so as to better map attention and metrics to scientific achievement.

The prevailing approach to assessing AI progress consists of measuring *performance*, such as the raw or normalized score in a game, ELO rating, error rate, accuracy, and so forth, often plotted over time to evaluate temporal progress [11, 32]. Performance, however, does not exactly correspond with social value or scientific progress in AI. Misalignment between what is measured and what is desired can lead to misallocation of energy and resources. Specifically, excessive effort is likely to go towards achieving novel performance benchmarks, and insufficient effort towards progress on other dimensions relevant to social value, economic value, and scientific progress, such as compute efficiency, data efficiency, novelty, replicability, autonomy, and generality.

This does not mean that quantitative assessment and benchmarks should be abandoned. On the contrary, we need more and better measurement [18]: measurement which is more comprehensive, general, and focused on the cost function of the ultimate beneficiaries. Ultimately, in assessing progress we would like to weight all the resources that users (or receivers) of a technology require to achieve their goals. For instance, to what extent does progress on a particular metric of performance in machine translation map on to user’s satisfaction? Does the progress also correspond to a reduction in cost per translation, or in time for execution? If a paper develops a new technique, how easily can new algorithms and applications integrate and benefit from it?

In general, users seek the benefits of high performance (at a set of tasks), while they seek to minimize the costs of developing and deploying their system. Sensitivity to costs is true for individual consumers, firms and developers, as well as other scientists. Some kinds of hidden costs can appear during development, when an application is produced, when reproduced at a large scale, or when adapted to other domains. Some future costs will be born by future developers or scientists, sometimes referred to as “technical debt” or “research debt”. Other costs may be spread more broadly, and are thus harder to account for. As in other sectors, there are externalities from AI development and deployment which are important to be aware of; among the negative externalities are environmental footprints, user privacy, skill atrophy (e.g., the Google effect), opacity in decision making, etc. Attention to, and ideally measurement of, these externalities is beneficial, as it is a first step towards internalizing them.

In this paper we consider this wide range of costs, though we focus on the costs born by future developers, such as the costs in computation, data, knowledge, software, human attention, and calendar time. We will identify how costs are distributed depending on the stage in which they are incurred, the number of times they are replicated, and the actor covering each cost. These dimensions should be integral to the measurement of AI progress.

The *estimation* of these dimensions is fraught with difficulties. To what extent are performance benchmarks actually representative of the target problem domain? To what extent are solutions overly specialized for the performance benchmark, as opposed to being more general, thus shaping the costs of adapting the solution to an adjacent problem domain? To what extent are solutions more reproducible by other teams, due to the availability of software and datasets? As an illustration of these difficulties and how they can be overcome, we will analyze several case studies where we evaluate performance alongside these other dimensions. As a result, we overhaul the notion of progress in these domains.

Our paper makes several contributions. First, we offer the most detailed and formal analysis to date of the dimensions of AI progress. While previous work has attempted to quantify progress in the performance of a specific system, we more fully account for the resources required and the generality of solutions. Second, in so doing we surface neglected dimensions of AI progress that may be optimized more directly. Third, we offer a novel framing under Pareto optimality for assessing performance and costs of an AI system, which suggests a more principled approach to forecasting future developments in AI, with myriad applications for policy, ethical, and economic analysis, and better research portfolio optimization within the field of AI itself.

2 Background

There was a time that benchmarks were unusual in AI, but today almost every area of AI has its own benchmarks and competitions [17]. Most researchers accept these challenges and invest great effort in improving on these metrics of performance. Indeed, many reports about AI progress include summaries of these benchmarks [11, 32]. We discuss here four issue areas arising from excessive focus on performance: *representativeness*, *specialization*, *reproducibility* and *resources*.

Regarding *representativeness*, many benchmarks and competitions are used in AI, but they vary in how representative they are of the fundamental problems in their respective subfields [17, 19]. For instance, it has recently been recognized that the Winograd Schema challenge only partially represents commonsense reasoning. As a reaction, challenges in AI are realigned to achieve more and better automation [13, 7]), or the aspiration of more human-like AI [23, 27]. A deeper concern is that most benchmarks are not really fostering the basic scientific advances needed to move the field forward, be they theoretical advances, explanatory insights, or tools to facilitate other work. This issue of non-representativeness is partly addressed through the review process, and requirements such as controlling the percentage of papers in different areas [31].

The second issue, *specialization*, is related to representativeness. When a benchmark or competition becomes the target, researchers will have incentives to overly specialize their systems to performance according to that benchmark, at the cost of other features of their system, such as generalizability. If we had a satisfactory metric of generality then we could use that as a benchmark, but it remains an open question how best to operationalize generality [18], balancing between putting all the distribution mass on a few tasks [25]—and not really being general—or distributing it in a block-uniform way—facing the no free lunch theorems [38].

A third issue is *reproducibility*, and the wider notion of replicability. In AI this was usually understood as requiring the sharing of data and code, but the concept is becoming richer [10, 4, 16]. Indeed, we must distinguish between specifically reproducing the results, and replicating the findings with some variations [39]. Several initiatives have been proposed to facilitate (or even require) a wider replicability. For instance, with the “open leaderboards” [36], participants have to upload their code so that other participants can make modifications and submit another proposal.

Finally, users are generally sensitive to the resource cost of developing and deploying an AI system, which performance benchmarks rarely explicitly take into account. Much AI progress is said to be attributable to advances in computational power [29]. However, it is not straightforward to quantify what exactly can be attributed to software progress, hardware progress or several other resources [6, 14]. Accordingly, perhaps it is more effective to just measure the so-called “end-to-end performance”, including computational time and quality of the models, such as the recent DAWNBench for deep learning [9]. Other resources, such as data, are at least as important, especially in machine learning¹. But it seems subjective to determine what input is seen positively or negatively, or even considered as cheating: too much data (supervised or unsupervised), too much knowledge (constraints, rules or bias), enriched input [5], etc. The question depends mostly on the cost of the resource. Human resources (“human computation”) are also common in AI to increase performance or generality (but at the cost of autonomy).

Overall, there are many resources involved but, at the moment, there is no integrated framework taking into account all of them. Related approaches involve the ideas of utility functions, Pareto-optimal analysis and, most especially, cost-sensitive learning [12]. [37] identifies costs related to inputs and outputs in classification (errors, instability, attributes, labeling, actioning) data (cases), computation and human preprocessing. In this paper, we offer a general statement of this idea, applied to AI progress.

In the end, when assessing AI progress in a comprehensive way, one should consider the whole life cycle of research, innovation, production, and reproduction. Notions such as technical or research debt are becoming more recognized, as they incorporate some costs that are not perceived at early stages of the process but appear later on, when the technology or product is put into practice [30, 16, 28].

3 Components and integration

In this section, we flesh out a comprehensive list of dimensions that are required for an “AI system” to work. We use the term “system” in a flexible way, including an agent, an algorithm, a product, etc., proposed in a research paper or by a company. Given the fuzzy contours of AI, human automation is usually recognized as a goal for AI. However, it is actually difficult to distinguish when reports and forecasts about “automation” [13, 7] are assuming conditions such as “at a reasonable cost”, “to a high degree of automation”, etc., versus “full automation at whatever cost”. The estimated probability of automation for a given task might change completely depending on these conditions. In the end, automation is important, but it is the efficiency of the whole system what matters, including any “human computation” involved. This view of efficiency links us directly to the resources involved in an AI system and their associated costs.

¹See <https://sites.google.com/site/dataefficientml/bibliography> for a bibliography on data-efficient ML.

Table 1 shows the resources we identified as frequently involved in developing and deploying AI systems. These resources have fuzzy boundaries and are often fungible with each other. For instance, the distinction between data and knowledge is not always clear, and hardware and software may be highly intertwined. Human resources are typically considered under “manipulation”, but can appear in relation to the other resources (e.g., labeled data and teaching a robot might be assigned to r_d and r_m respectively). This is not a problem, as long as all the resources are identified.

	Description	Example
r_d	<i>Data</i> : All kinds of data (unsupervised, supervised, queries, measurements).	A self-driving car needs on-line traffic information.
r_k	<i>Knowledge</i> : Rules, constraints, bias, utility functions, etc., that are required.	A spam filter requires the cost matrix from the user.
r_s	<i>Software</i> : Main algorithm, associated libraries, operating system, etc.	A planner uses a SAT solver.
r_h	<i>Hardware</i> : Computer hardware, sensors, actuators, motors, batteries, etc.	A drone needs a 3D radar for operation.
r_m	<i>Manipulation</i> : Manual (human-operated) intervention through assistance	A robot needs to be manually re-calibrated.
r_c	<i>Computation</i> : Computational resources (CPU, GPU usage) of all the components	A nearest neighbor classifier computes all distances.
r_n	<i>Network</i> : Communication resources (Internet, swarm synchronisation, distribution).	An automated delivery system connects all drones.
r_t	<i>Time</i> : Calendar (physical) time needed: waiting/night times, iteration cycles.	A PA requires cyclical data (weeks) to find patterns.

Table 1: Resources that are frequently needed by AI systems.

It is appealing to collapse the benefits and costs of an AI system to a single metric. For any given user with rational (transitive and complete) preferences, their preferences can be represented using a utility function. A firm’s utility function, for example, might correspond to risk-adjusted expected profit. A user’s utility function might be harder to quantify, but is generically increasing in the performance of the system and decreasing in the costs of the system. Denote a performance vector, ψ , for a given problem, which is often a unidimensional quantitative score (such as the error), but could also have several components. A utility function maps performance and all associated resources to a single dimension:

$$U(\psi, \bar{r}) = U(\psi, r_d, r_k, r_s, r_h, r_m, r_c, r_n, r_t) \rightarrow u \quad (1)$$

In some cases this is an additively separable function, such that $U(\psi, \bar{r}) = B(\psi) - \sum_x C_x(r_x)$, with the first term accounting for the benefit according to the performance of the system minus the costs produced by the use of resources (note that the cost functions C_x are different for each resource). For economic applications, we might be able to separate the utility function into performance generating revenue (in dollars), and resources imposing costs (in dollars).

In many cases, we are not able to collapse performance and costs into a single metric, perhaps because the utility function is not known or varies across a population of users. Still, we can productively examine the relative performance and costs of different systems. For any number of dimensions, we can assess the Pareto-optimal surface, as we do in Fig. 1 for two indicators (we explore this further in section 5). We may want to focus on one dimension of costs, such as economic costs or energy costs (as per the “carbon footprint”). For example, Fig. 1 shows algorithms and architectures according to their MNIST prediction error and power consumption, revealing that most solutions are not on the Pareto surface on these dimensions, with notable exceptions, such as some ASIC architectures, which focus on efficiency in terms of chip space, speed and “energy footprint” [8].

4 The full range of accounting

The benefits and costs of developing and deploying an AI system are not incurred only once, but throughout the many uses, reuses, and follow-on contributions. Some costs are born exclusively during the initial conception and development, while others recur with each adaptation to a new application, or even each application to a particular user. In general, the total resource burden should be accounted for according to the whole cycle of the AI system.

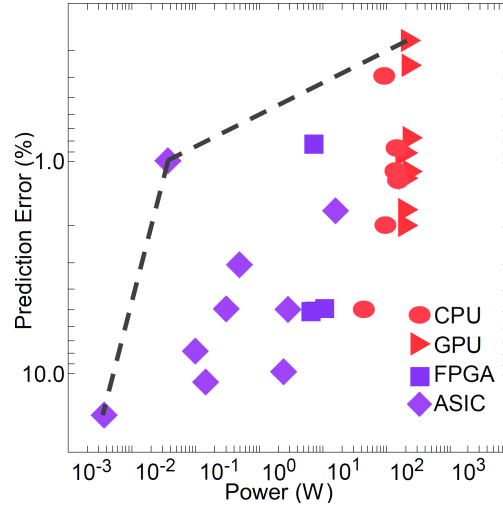


Figure 1: Performance for MNIST [24], for 22 papers, compared to power consumption (data from [29]). The Pareto front is also shown (we will discuss whether the points can actually be joined by straight segments in section 5).

Fig. 2 shows how the dimensions we identified can become relevant at different stages of the life cycle of an AI system. Consider a new algorithm for voice recognition. Apart from all the human thinking, there will be a great effort in terms of failed experiments, different libraries used, users testing the early systems, etc. If a company takes these ideas and builds a prototype, the tests, software, hardware, and compute will concentrate on production. When the system is reproduced (installed or shipped) to users, additional resource costs will be incurred. Further, if the idea can be adapted for other applications (e.g., adapting a voice recognition system to other languages), depending on its generality and reproducibility, the initial contribution can provide further value, at some further adaptation cost including the need for new corpora, training, semantic knowledge, etc.

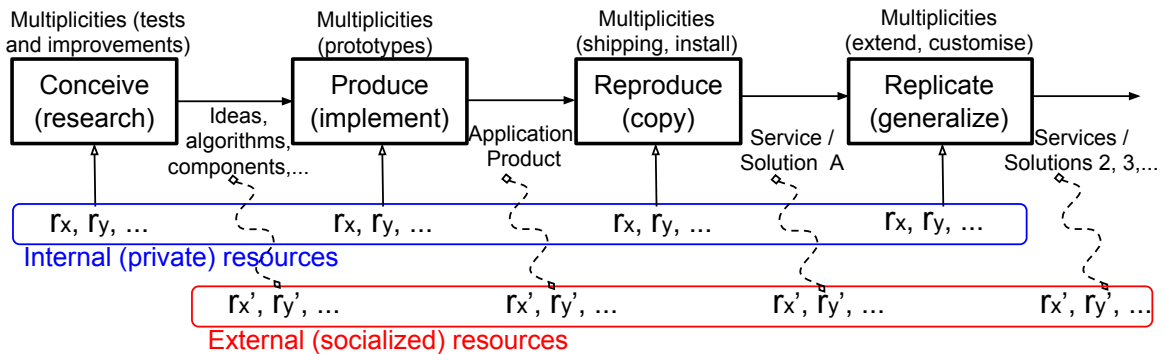


Figure 2: Illustrative representation of stages of the AI system life cycle where resources might be required.

At each stage of the life cycle, the contribution may be deployed a multiplicity of times (represented above the boxes in Fig. 2). The total value of the contribution thus needs to take into account the scale of its deployment. For instance, some early speech recognition systems were pre-trained once (the *system cost*, denoted by C , covering the “conceive” and “produce” stages in Fig. 2) and then adapted to thousands of users, with extra hours of customization per user (the *application cost*, denoted by C^j with j indexing each of the n applications, or users, covering the “reproduce” and “replicate” stages). More recent general speech recognition systems do not need such customization. Consequently, the application cost C^j is lower per user. In both cases, the total cost C is $C + \sum_{j=1}^n C^j$. As the number of applications increases, the average cost will converge to the average application cost as the system cost is amortized. For this reason, for contributions that have many possible applications, it is worth paying additional system costs so as to make the contribution more general, adaptable, and reusable, and thereby bring down the application costs. Since AI contributions often have broad potential applicability, contributions that are general, adaptable, and reusable are likely to

have especially high utility.

Fig. 2 not only covers direct “internal” costs (r_x, r_y, \dots) but also some external “debts” or “socialization” costs (r'_x, r'_y, \dots). For instance, automated customer service systems (call centers) clearly were not a Pareto improvement relative to previous systems, even though they may be a profit maximizing improvement. Companies reduce their labor costs for customer service by substituting in phone-trees and voice recognition, but in the process impose time, frustration, and other costs onto the customer. Some navigators and personal assistants can make users more dependent on them, atrophying some capabilities or leading to a simplification of language. In other words, the user adapts to the AI system, and assumes part of the effort or cost. In general, technological innovation both involves developing technology to fit a given conception of the task, and adapting conceptions of the task to fit the capabilities of technology. In the process of adapting work processes, customer expectations, relationship norms, and even urban design to what is technologically convenient, there can be consequences for society that are not internalized by the designers and deployers of these systems. This footprint of AI is not usually acknowledged in benchmarking.

From the previous sections, we conclude that the contribution of an AI development should, in principle, be given a full accounting of the costs and benefits, across the contribution’s full life cycle. The current emphasis on targeting and reporting performance benchmarks, however, poses an obstacle to a full accounting. Reproducibility and replicability are two traditional tools for addressing this. More precisely:

- *Specific reproducibility* refers to whether the *same result* can be obtained from the same conditions and procedures. In AI, this requires that all the necessary code and data are given. This also assumes the same cost functions as well: $\sum_{j=1}^n \sum_x C_x^j(r_x^j) = n \sum_x C_x(r_x)$.
- *General replicability* will check whether the AI technique can be *applied to other problems*, a set of n tasks, applications, or users indexed by j , with an overall cost $\sum_{j=1}^n \sum_x C_x^j(r_x^j)$ that must consider the adaptation effort, with different resources r_x^j and cost functions C_x^j per user.

Especially for replicability, we can experiment with different hardware architectures, change some of the software and get different computational costs, apart from different performance. That means that the partial results for each B^j and $C_x^j(r_x^j)$ might be different, but we still have something replicable with similar utility. A clear example of this notion of replicability is “approximate computation” in deep learning, where one can get much smaller computational costs without a significant change in accuracy [29].

5 Exploring the Pareto-front of AI research

Corporations, governments, startups, NGOs, personal users, and contemporary and future AI researchers are the intended recipients, or *receivers*, of the AI technologies being developed, and they each have different preferences, resources and constraints, or in other words different operating characteristics. The familiar concept of the ROC curve plots true positive rates (TPR) and false positive rates (FPR) for binary classifiers, and emphasizes the importance of comparing multi-dimensional surfaces, rather than single metrics.

For instance, Fig. 3 (left) just shows a single metric, performance, as a function of time. This plot does not explain what the cluster of attempts after 2014 really contribute, when they have more error than the already obtained human level. Other dimensions are neglected in this plot, limiting insight about progress. In the next section we will see other domains where some of the resources are actually put as dimensions.

Before analyzing the case studies, we have to understand how to build and work with the Pareto front. When resources are included, the analysis of optimal Pareto surfaces might be slightly different than the traditional triangulization approach. When showing performance metrics such as TPR and FPR for two models, any point in between can be obtained by interpolation, connecting any two points by a straight segment. However, we should note that these points require the implementation of both models. While some of the resources can be interpolated, others (e.g., software) will simply sum up, and the points between two other points will not be achievable with a straight line, but by an axis-parallel route.

For instance, Fig. 3 (right) shows performance against one particular resource. For each method, A, B, C, D, and E, the numbers represent the extremes when varying their parameters. E1 represents a random (or baseline) model. Assuming no interpolation is possible, the Pareto front here is shown in blue. Methods C and B can be discarded, as they do not reach anywhere that cannot already be achieved with A, E and D.

The diversity of receivers and the number of dimensions suggest that a single utility metric is simplistic — different receivers would have different subjective utilities for different dimensions. This operating condition

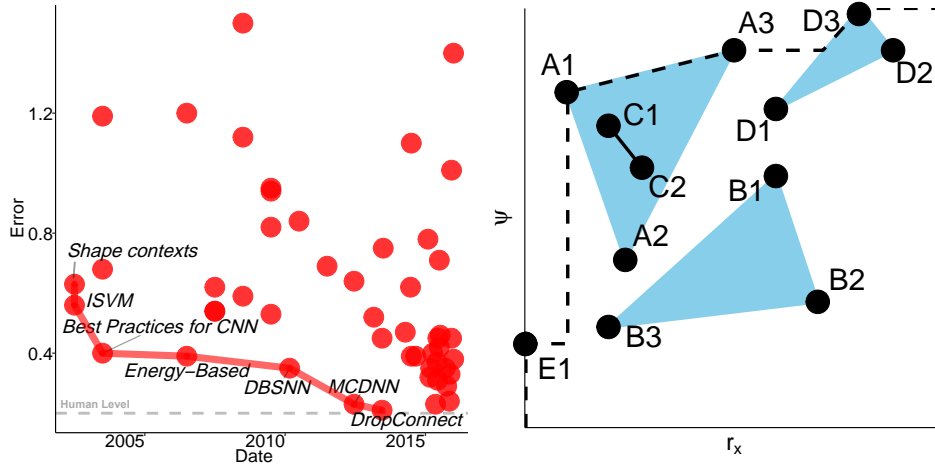


Figure 3: Left: Performance for the MNIST benchmark (data from EFF). Right: A schematic representation of techniques A, B, C, D, E, with variants, the areas they cover, and the resulting Pareto front.

translates into a vector, or gradient, in the multidimensional space. For example, large technology corporations may gain significant utility from a discovery that allows modest speed-ups in exchange for significantly increased compute demands, whereas individual researchers, personal users and startups may find little value in such a discovery. Conversely, the existence of real recipients whose preferences can be known in advance allows us to prioritize exploration of those configurations. From the above, we derive a few criteria to identify progress events:

- *Improving the Pareto front for a known group of recipients* (A1, A3 or D3 in Fig. 3, right). This would include all-else-being-equal improvements in performance, but also reductions in computation, data, manipulation or other resources in Table 1. This would not, however, consider extreme regions no recipient assigns value to.
- *Covering a location slightly under the Pareto front with more flexibility* (B3 in Fig. 3, right). Instead of reaching some areas by combining existing approaches, a new technique can reach there easily with a trade-off between its own parameters, allowing more receivers to easily find their subjectively optimal trade-offs.
- *Covering a location slightly under the Pareto front with more diversity* (C in Fig. 3, right, if it is very different from A). The current dominant technique or paradigm can push the Pareto-optimal front for some time, but slightly suboptimal approaches, especially if they are radically different (i.e., alternative “research programs”), should not be discarded because they may lead to potential improvement in the Pareto-optimal front if the current paradigm stalls.

Receivers can be incentivized to generate and communicate their gradients (though in some cases, countervailing considerations may exist such as commercial secrecy). It is also in the interests of discoverers to show the recipients benefited by their discovery. Brokers of such information (peer-review, surveys, competitions, etc.) are in a position to meet the incentives (and gradients) of both researchers and recipients by ensuring such discoveries are properly rewarded.

6 Case studies

In this section we will examine two representative case studies of progress in AI: Alpha* and ALE.

Alpha* refers to a series of papers and associated techniques by DeepMind to play board games. We analyzed the whole series, from AlphaGo [33] (including the Fan and Lee versions, used against Fan Hui and Lee Sedol, respectively, and its latest version, AlphaGo Master, which won 60 straight online games against professional Go players), AlphaGo Zero [35] (a version created without using data from human games) and AlphaZero [34] (which uses an approach similar to AlphaGo Zero to master not just Go, but also chess and shogi).

	AlphaGo Fan	AlphaGo Lee	AlphaGo Master	AlphaGo Zero	AlphaZero
r_d (Data)	✓	✓	✓	✓	✓
r_k (Knowledge)	○	○	○	○	○
r_s (Software)	○	×	×	○	×
r_h (Hardware)	×	×	×	×	×
r_m (Manipulation)	✓	✓	✓	✓	✓
r_c (Computation)	✓	○	○	✓	○
r_n (Network)	-	-	-	-	-
r_t (Time)	-	-	-	-	-
ψ (Performance)	✓	✓	✓	✓	○

Table 2: Dimensions (resources and performance) reported in the Alpha* papers. Systems from [33, 35, 34].

Table 2 shows whether the dimensions were reported in the papers (✓), only partially accounted for (○), not mentioned but relevant (×) and not applicable (−). Many dimensions are relevant for the analysis: the data, the knowledge, the software, the hardware, manipulation, computation and, of course, performance, etc. However, only some of them are provided, which makes a comprehensive comparison of the whole space difficult. Still, we will represent three dimensions: performance (in ELO ranking, which can only be partially estimated for AlphaZero), computational resources (using the equivalence: $1 \text{ TPU}_{v2} \simeq 3 \text{ TPU}_{v1} \simeq 36 \text{ GPU} \simeq 180 \text{ CPU}$ [21]) and human manipulation resources (as represented quantitatively by the ELO ranking of the player or players the system learns from). Other dimensions (like knowledge² about Go, software, etc.) are not included because of insufficient information from some papers.

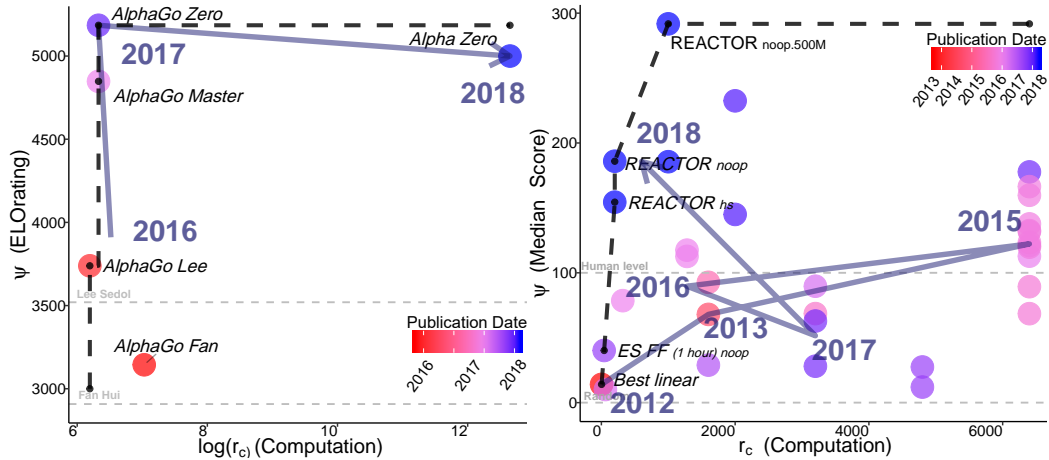


Figure 4: Multidimensional utility space for Alpha* (left) and ALE (right). Research gradient evolution from 2013 to 2018 represented with a segmented gray arrow. The Pareto front (dashed black) does not include other resources (software, and humans used for training) that duplicate for connecting segments.

What we see in Fig. 4 (left) is that the Pareto front at the moment is represented by AlphaGo Lee and AlphaGo Zero. AlphaGo Fan is discarded because AlphaGo Zero needs less compute, no manipulation and gets better performance.

Why is AlphaZero a breakthrough if it is not Pareto optimal? The answer is generality. AlphaGo* only solved one task (Go) and AlphaZero can solve several tasks. Finally, if we look chronologically at the plot, we see that the main gradient that has been followed has been performance.

The second case study is ALE [3], a collection of Atari games that has become popular for the evaluation of general-purpose RL algorithms learning from screen shots. We selected all the papers (systems) from

²We have the constructed features: stones to be captured or escaped, legal moves, ‘liberties’, etc. While this knowledge is crucial, there is no cost for a new match (reproduction), but the adaptation of AlphaZero to other games (replication), may be important.

EFF’s AI Progress Measurement Project [11] and the papers introducing the Rainbow [20] and REACTOR agents [15].

	Sarsa	Best Linear	DQN best	NatureDQN	Gorila	DQN _{noop} & hs	DUEL _{noop} & hs	DDQN _{tuned} hs	PRIOR _{hs} & noop	P. DUEL _{hs} & noop	AC3 _{LSTM, FF & FFid}	DDQN _{Pop-Art} noop	AC3 _{CTS}	SARSA _e & f-EB	TRPO _{hash}	DQN _{CTS} & PixelCNN	C51 _{noop}	ES FF (1h) noop	RAINBOW	REACTOR
r_d	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
r_k	○	○	×	✓	×	○	×	○	○	○	○	○	×	×	○	○	○	×	✓	✓
r_s	×	×	×	✓	×	×	×	×	×	×	×	×	×	×	×	×	×	×	✓	×
r_h	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×
r_m	×	✓	×	✓	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×	×
r_c	○	○	○	○	○	○	○	○	○	○	✓	○	○	○	○	○	○	○	✓	✓
r_n	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
r_t	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
ψ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	○	✓	✓	✓	✓

Table 3: Same as Table 2 for the ALE papers (from EFF [11] and [15, 20]).

Table 3 shows what information we found about the resources and performance. Again, many dimensions are relevant, but only a few are systematically reported: data, computation and performance. Fig. 4 (right) represents computation and performance. Computation time (whenever the authors do not provide this information explicitly) is roughly estimated from the kind of approach used, whether it is follow-up work, the training setting used, etc., or from figures in more recent papers, which make explicit comparisons between them and the state of the art [20, 15]. What we see in Fig. 4 (right) is a current Pareto front dominated by REACTOR variants, ES FF and Best Linear. The research gradient (in gray) has changed over the years, with some disregard of compute initially and more concern in efficiency recently.³

For this benchmark, it is common to find “learning curves” in the papers (e.g., [26]), which show performance varying on the number of episodes. This is clearly the r_d (data) but it also influences directly on computation. These learning curves give information of full regions of the multidimensional space, as we saw in Fig. 2.

Finally, for some papers, the comparison was not possible (e.g., due to different subsets of games). It is important to note, however, that some approaches based on genetic programming [22] and on planning [2] are valuable in terms of diversity.

7 Conclusions

The interest in more comprehensive evaluation protocols, going beyond performance alone, is represented by some of the references we included in section 2 on cost-sensitive learning, reproducibility, generality, data-efficiency and computational costs. However, in order to rigorously evaluate a novel contribution to AI progress more broadly, we need a more formal analysis. This is done by an explicit enumeration of all the dimensions (as represented by Table 1) and their integration into utility functions or their representation in a multidimensional space, with a clear delimitation of the extent of accounting. This is what happened in cost-sensitive learning more than 15 years ago [12, 37], leading to a wide range of techniques that covered different operating conditions. While all these costs are nowadays integrated into the measures of performance, many

³ The computation times shown in Fig. 4 (left) include both training and deployment (system and application costs). Hence, a model that is half way between models A and B (choosing between them with equal probability), denoted by \overline{AB} , has performance $\psi(\overline{AB}) = 0.5\psi(A) + 0.5\psi(B)$, but has a computational cost of $r_c(\overline{AB}) = r_c(A) + r_c(B)$. This is why the Pareto front in Fig. 4 (left) has parallel segments, as in Fig. 3 (right). However, in Fig. 4 (right), we can have A train and play for half of the ALE games and B train and play for the rest. As we average for the whole set of games, we can actually have $r_c(\overline{AB}) = 0.5r_c(A) + 0.5r_c(B)$, at least if there is no transfer effort between games. This is why the Pareto front on the right is shown with direct straight segments.

other resources are not, as we have surfaced here. We hope this paper can launch the study of “cost-sensitive AI”. Within this framework, we make a series of recommendations:

- Benchmarks and competitions should be defined in terms of a more comprehensive utility function, considering as many dimensions as possible, or recognize the value of all contributions that have any of the positive effects on the Pareto front identified in Section 5, in short or long terms.
- Papers presenting or evaluating algorithms should generally try to report the whole region they cover, and how to navigate the region by modifying parameters or resources. There are many partial examples nowadays: learning curves, plots comparing the number of models vs. performance, planning performance vs. lookahead, etc.
- These utility functions and multidimensional spaces must also be seen in terms of replicability, for variants of the problems and at different stages of the AI life cycle. The multiplicities are more difficult to plot graphically, but we can still define operating conditions depending on the adaptation (or transfer) effort for m problems, or n users.

Frequently, we will not be able to say that one technique is ‘better’ than another: they just cover different regions of the multidimensional space. It is the receiver who will choose the system that best fits their needs. Having a representation of the Pareto front may hugely facilitate this choice for other researchers and industry, as simply as moving the gradient until touching the Pareto surface. Also, small players in AI could focus on those areas that require less resources and still contribute to the Pareto front or to diversity. Finally, the Pareto surface can help detect some societal risks, especially if we see that a powerful capability in AI can be achieved with very few resources, becoming available to malicious actors.

This view of the operating condition as a gradient may suggest clever approaches to push the front for some resources, as gradient descent is increasingly being used at a meta-level [1]. In general, we hope this paper will help change perceptions, promote more general and versatile techniques, highlight the trade-offs, and raise awareness of the overall “AI footprint”, well beyond performance.

References

- [1] Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, and Nando de Freitas. Learning to learn by gradient descent by gradient descent. *CoRR*, abs/1606.04474, 2016.
- [2] Wilmer Bandres, Blai Bonet, and Hector Geffner. Planning with pixels in (almost) real time. *arXiv preprint arXiv:1801.03354*, 2018.
- [3] 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, jun 2013.
- [4] Fabio Bonsignorio and Angel P Del Pobil. Toward replicable and measurable robotics research. *IEEE Robotics & Aut. M.*, 22(3):32–35, 2015.
- [5] Nicolas Bougie and Ryutaro Ichise. Deep reinforcement learning boosted by external knowledge. *arXiv preprint arXiv:1712.04101*, 2017.
- [6] Miles Brundage. Modeling progress in ai. *The Workshops of the Thirtieth AAAI Conference on Artificial Intelligence AI, Ethics, and Society: Technical Report WS-16-02*, *arXiv preprint arXiv:1512.05849*, 2016.
- [7] Erik Brynjolfsson and Tom Mitchell. What can machine learning do? Workforce implications. *Science*, 358(6370):1530–1534, 2017.
- [8] Tianshi Chen, Zidong Du, Ninghui Sun, Jia Wang, Chengyong Wu, Yunji Chen, and Olivier Temam. Diannao: A small-footprint high-throughput accelerator for ubiquitous machine learning. In *Sigplan Not.*, volume 49, pages 269–284. ACM, 2014.
- [9] Cody Coleman, Deepak Narayanan, Daniel Kang, Tian Zhao, Jian Zhang, Luigi Nardi, Peter Bailis, Kunle Olukotun, Chris Ré, and Matei Zaharia. Dawnbench: An end-to-end deep learning benchmark and competition, 2017.

- [10] C. Drummond. Replicability is not reproducibility: nor is it good science. *Evaluation Methods for ML Ws at the 26th ICML, Montreal, Canada*, 2009.
- [11] P Eckersley and N Yomna. Measuring the progress of AI research, 2017.
- [12] Charles Elkan. The foundations of cost-sensitive learning. In *IJCAI*, volume 17, pages 973–8, 2001.
- [13] Carl Benedikt Frey and Michael A Osborne. The future of employment: how susceptible are jobs to computerisation? *Technological Forecasting and Social Change*, 114:254–280, 2017.
- [14] Katja Grace. Trends in algorithmic progress, 2017.
- [15] Audrunas Gruslys, Mohammad Gheshlaghi Azar, Marc G. Bellemare, and Rémi Munos. The reactor: A sample-efficient actor-critic architecture. *CoRR*, abs/1704.04651, 2017.
- [16] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. *arXiv preprint arXiv:1709.06560*, 2017.
- [17] José Hernández-Orallo. Evaluation in artificial intelligence: from task-oriented to ability-oriented measurement. *Artificial Intelligence Review*, 48(3):397–447, 2017.
- [18] José Hernández-Orallo. *The measure of all minds: evaluating natural and artificial intelligence*. Cambridge University Press, 2017.
- [19] Jose Hernández-Orallo, Marco Baroni, Jordi Bieger, Nader Chmait, David L Dowe, Katja Hofmann, Fernando Martínez-Plumed, Claes Strannegård, and Kristinn R Thórisson. A new ai evaluation cosmos: Ready to play the game? *AI Magazine, Association for the Advancement of Artificial Intelligence*, 2017.
- [20] Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Daniel Horgan, Bilal Piot, Mohammad Gheshlaghi Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. *CoRR*, abs/1710.02298, 2017.
- [21] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of a tensor processing unit. In *44th Intl Symposium on Computer Architecture*, pages 1–12. ACM, 2017.
- [22] Stephen Kelly and Malcolm I Heywood. Emergent tangled graph representations for atari game playing agents. In *European Conference on Genetic Programming*, pages 64–79. Springer, 2017.
- [23] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *BBS*, 40, 2017.
- [24] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998.
- [25] Shane Legg and Marcus Hutter. Universal intelligence: A definition of machine intelligence. *Minds and Machines*, 17(4):391–444, 2007.
- [26] Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, and Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open problems for general agents. *CoRR*, abs/1709.06009, 2017.
- [27] Gary Marcus. Deep learning: A critical appraisal. *arXiv preprint arXiv:1801.00631*, 2018.
- [28] Chris Olah and Shan Carter. Research debt. <https://distill.pub/2017/research-debt/>, 2017.
- [29] Brandon Reagen, Robert Adolf, Paul Whatmough, Gu-Yeon Wei, and David Brooks. Deep learning for computer architects. *SL on Comp. Architecture*, 12(4):1–123, 2017.
- [30] D Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-Francois Crespo, and Dan Dennison. Hidden technical debt in machine learning systems. In *NIPS*, pages 2503–2511, 2015.

- [31] Nihar B Shah, Behzad Tabibian, Krikamol Muandet, Isabelle Guyon, and Ulrike Von Luxburg. Design and analysis of the nips 2016 review process. *arXiv preprint arXiv:1708.09794*, 2017.
- [32] Yoav Shoham, Raymond Perrault, Erik Brynjolfsson, Jack Clark, and Calvin LeGassick. AI Index, 2017.
- [33] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [34] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017.
- [35] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550:354, 2017.
- [36] J. Spohrer. Opentech AI workshop, 2017.
- [37] Peter D Turney. Types of cost in inductive concept learning. *arXiv preprint cs/0212034*, 2002.
- [38] David H Wolpert. What the no free lunch theorems really mean; how to improve search algorithms. In *Santa fe Institute Working Paper*, page 12. 2012.
- [39] Rolf A Zwaan, Alexander Etz, Richard E Lucas, and M Brent Donnellan. Making replication mainstream. *Behavioral and Brain Sciences*, pages 1–50, 2017.