# Aligning artificial intelligence with climate change mitigation

Lynn H Kaack, Priya L Donti, Emma Strubell, George Kamiya, Felix Creutzig, David Rolnick

# Aligning artificial intelligence with climate change mitigation

Lynn H. Kaack[1,2,3*], Priya L. Donti[4,5], Emma Strubell[4], George Kamiya[6], Felix Creutzig[7,8], and David Rolnick[9,10]

[1]Data Science Lab, Hertie School
[2]Energy and Technology Policy Group, Department of Humanities, Social and Political Sciences, ETH Zurich
[3]Institute of Science, Technology, and Policy, ETH Zurich
[4]School of Computer Science, Carnegie Mellon University
[5]Department of Engineering & Public Policy, Carnegie Mellon University
[6]International Energy Agency
[7]Mercator Research Institute on Global Commons and Climate Change, Berlin
[8]Sustainability Economics of Human Settlements, Technical University Berlin
[9]School of Computer Science, McGill University
[10]Mila – Quebec AI Institute
[*]*Corresponding author*

October 6, 2021

**Abstract**

Assessing and shaping the effects of artificial intelligence (AI) and machine learning (ML) on climate change mitigation demands a concerted effort across research, policy, and industry. However, there is great uncertainty regarding how ML may affect present and future greenhouse gas (GHG) emissions. This is owed in part to insufficient characterization of the different mechanisms through which such emissions impacts may occur, posing difficulties in measuring and forecasting them. We therefore introduce a systematic framework for describing ML's effects on GHG emissions, comprising three categories: (A) compute-related impacts, (B) immediate impacts of applying ML, and (C) system-level impacts. Using this framework, we assess and prioritize research and data needs for impact assessment and scenario analysis, and identify important policy levers.

## Introduction

As artificial intelligence (AI) and in particular machine learning (ML) are increasingly deployed across society [1], there has been a surge in interest in understanding the effects ML may have on climate action [2–4]. To explicitly and consistently account for ML in long-term climate and energy projections, and in the design of appropriate policies, the research community needs to develop a holistic and operational understanding of the different ways in which ML can positively and negatively impact climate change mitigation and adaptation strategies. In particular, those impacts that are easiest to measure are likely not those with the largest effects. This can lead to challenges in terms of estimating macro-scale effects, picking up on underlying dynamics and trends, and prioritizing actions to align ML with climate strategies. To aid in addressing these challenges, we present a systematic framework (Figure 1) for categorizing the different kinds of impacts of ML on global greenhouse gas (GHG) emissions — through compute-related impacts, the immediate impacts of ML applications, and the system-level changes ML induces.
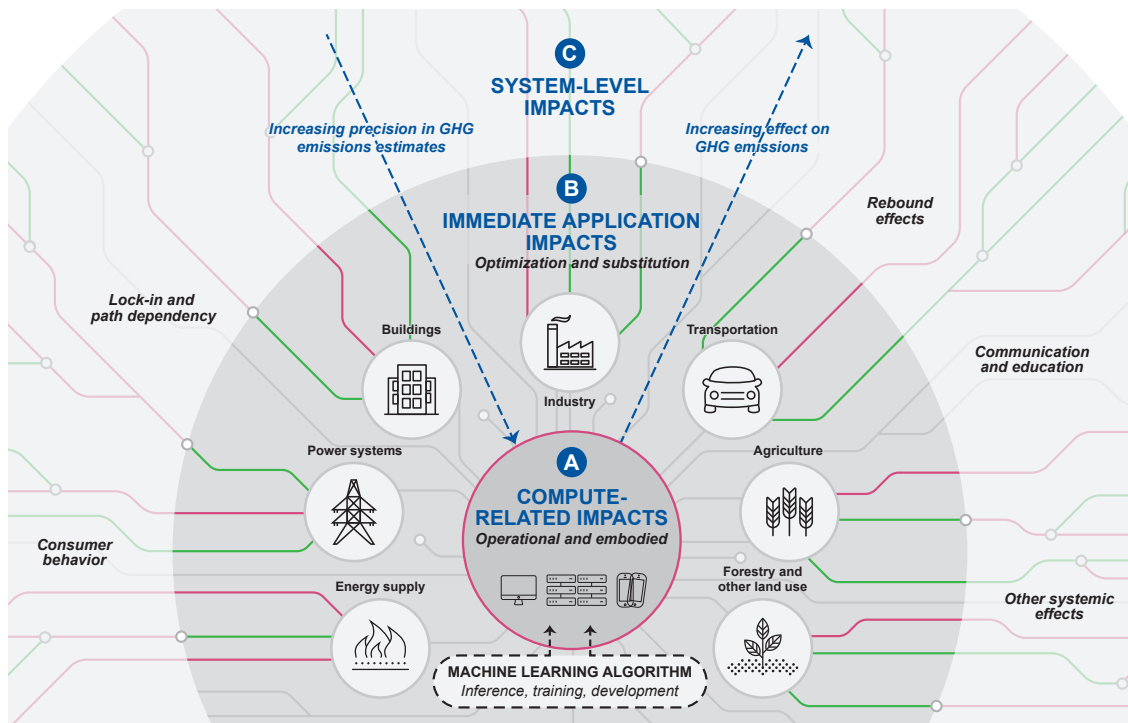
Figure 1: A framework for assessing the greenhouse gas (GHG) emissions impacts of machine learning. We distinguish between three categories (A, B, and C) with different kinds of potential emissions impacts, estimation uncertainties, and associated decarbonization levers. Green denotes effects relating to reductions in GHG emissions, and magenta to increases in emissions.

Recent work has discussed different aspects relating to this framework, describing applications of ML for tackling climate change [5], applications of ML that increase emissions [6, 7], and the energy consumption of ML through software and hardware [7–9]. A few pieces have engaged with both the positive and negative effects of ML on the climate [2, 3, 10–14], but no work has explicitly provided an overview of the different mechanisms by which ML may impact emissions. By presenting a unified framework of these mechanisms, we intend to provide a starting point for research, policy-making, and organizational action aiming to better align ML with climate change strategies.

Related literature on assessing the impacts of information and communications technologies (ICT) has often distinguished between the energy- and hardware-related GHG emissions of ICT ("direct" impacts) and the emissions impacts of ICT's applications ("indirect" impacts) [15–17]. Our framework in Figure 1 similarly distinguishes between the compute-related GHG emissions of ML, and the emissions reductions and increases resulting from applications of ML. Considering that ML encompasses a particularly novel and transformative set of software and analytics approaches with nuanced downstream effects, our framework covers three main aspects. The first involves the GHG emissions resulting from compute, caused by both the electricity used for ML computations and the embodied emissions associated with computing hardware. The second involves the "immediate" GHG emissions effects tied to the short-term outcomes of applications of ML. The third involves the structural or "system-level" GHG effects induced by these applications. Drawing a clear line between these latter two application-level effects is rather difficult, with different classifications available throughout the literature (see, e.g., Horner et al. [16]); our distinction is adapted from Hilty and Aebischer [15], and, while imperfect, is important for framing the discussion of ML's overall impacts and associated levers. We report quantitative assessments where available and where we believe these assessments are representative, and we discuss the current state of research on impact assessment. We then use our framework to propose a roadmap for assessing and forecasting impacts, and discuss approaches for shaping the impacts of ML. In terms of scope, our framework predominantly focuses on algorithm-related impacts, and omits impacts relating to data collection and management, ICT, and digitalization more broadly [2, 16, 18–20].

# Compute-related impacts

In order to assess direct compute-related impacts, we take two different perspectives. The first is a *bottom-up* perspective aimed at assessing the energy use of individual ML models, capturing aspects of the use, development, and design of these models. The second is a *top-down* perspective aimed at estimating the total global GHG emissions associated with ML workloads, capturing both the sourcing of the electricity used to power computations as well as embodied emissions from materials extraction and manufacturing.

## ML model development and deployment

Creating and running an ML model uses computing power and therefore energy, with the amount varying dramatically between different algorithms and different stages in the ML model life cycle. While many models used in practice are relatively small and can be trained and run on a laptop (such as linear classifiers or decision trees), state-of-the art performance on more complex tasks is often achieved with very large models, typically using deep learning. The size of the largest deep learning models (measured in number of parameters), and likely the size of the average model, is growing rapidly, leading to much larger demand for computing resources [21–23].

To illustrate how ML models differ so drastically in the energy they consume and better understand approaches to reduce their energy consumption, it is necessary to take a deeper dive into the life cycle stages of an ML model: model inference (or use), model training, and model development and tuning. Model *inference* describes the stage where the model is in use in the world: for instance, given new inputs, such as

images, it labels those inputs (e.g., identifies whether an image is of a cat or a dog) according to a function that it has learned. The goal of the model *training* stage is to learn the underlying function that, e.g., maps from inputs to labels, by analyzing a dataset to choose a set of parameters that define the function. During model *development and tuning*, a researcher will typically train many different model variants on different datasets, in order to devise a variant that works best in the given problem setting.

Figure 2 ("Bottom-up") provides a schematic overview of relative energy requirements and frequency of each stage of the ML model life cycle. Model inference is the least energy-intensive process in the ML model life cycle, but it is likely to occur the most frequently. For instance, classifying toxic comments [24, 25] or the contents of images [26] on social media requires little power each time a model is used, but may be used on the order of billions of times a day. Also larger models, such as Google's machine translation system, may process more than 100 billion words per day [27]. Those computing requirements may add up: Amazon estimates that up to 90% of financial costs for ML workloads in production are due to inference [28] (similar numbers for energy consumption are not available and this ratio might not be the same for energy). The training stage may require many passes over the dataset, often denoted as *epochs*, with each epoch performing full model inference on each example, as well as computing updates to correct the model's prediction for future iterations. In the case of deep learning, for example, this means that each epoch typically requires about three times as much computation as inference itself [29]. Training an ML model is thus more energy-intensive than using it, but is done much less frequently. Hazelwood et al. [25] report that ML models in Facebook's datacenters are re-trained anywhere from hourly to multi-monthly. The most energy-intensive stage of the ML model life cycle is model development, which requires training many different models. Modern ML models that use neural networks are particularly energy-intensive in the development phase as they have many more possible model configurations than their predecessors, and it is not well understood how those configurations should be set to perform well on a given dataset, except through trial and error experimentation and validation (*hyperparameter search*), often involving thousands of training runs. For instance, the GHG emissions associated with developing certain large, cutting-edge models can be comparable to, e.g., the lifetime carbon emissions of a car [8], though such computationally intensive processes are performed rarely and by the fewest entities.

The computational requirements of ML models are often described in FLOPs: FLoating point OPerations, or the number of additions and multiplications of scalar values required to obtain a result. The precise mapping from FLOPs to energy draw is hardware- and algorithm-dependent, but more FLOPs generally correspond to higher energy use. ResNet-50 [30], a popular deep learning model for image classification, requires about 4 billion FLOPs (and 65 milliseconds) to map a 224x224 pixel input image to a label, with an error rate of 24.6% [31]. A less computationally efficient version, ResNet-152, requires about 11 billion FLOPs (and 150 milliseconds) per image, and obtains only a slightly better error rate of 23.0%. This case illustrates a trade-off in energy efficient ML: Is it worth the more than 2.5x increase in FLOPs, and corresponding energy, to reduce the error rate by 1.6%? Will the benefits outweigh the costs from both an emissions and a broader societal perspective [32]?

Software tools for measuring ML model energy use [33] and carbon emissions [34, 35] are already available, metrics for reporting model accuracy as a function of computational budget have been proposed [36, 37], and benchmarks measuring training and inference efficiency have been established [38, 39]. Yet, such reporting is still not standard for researchers and ML software maintainers. Standardized reporting is essential for including efficiency considerations during model development and for making energy consumption a criterion when choosing between different ML approaches in practice.

As larger neural network models have become more prominent in certain areas of machine learning, research into improving the efficiency of ML models has started to increase [26, 40–45], and now also discusses implications of compressing models on broader performance characteristics [46]. However, the vast majority of ML research and development still focuses on improving model accuracy, rather than balancing accuracy

and energy efficiency [9].

## Computing infrastructure

The global ICT sector — consisting of all data centers, data transmission networks, and connected devices — accounted for around 700 MtCO2e in 2020, equivalent to around 1.4% of global GHG emissions [47, 48]. Around two-thirds of the sector's emissions come from operational energy use, with the remainder resulting from materials extraction, manufacturing, transportation, and end-of-life [48].

Only a fraction of emissions from the ICT sector is attributable to AI and ML (Figure 2 "Top-Down"), but its exact share is not known due to challenges in boundary definition and a lack of data and established methodology. Based on the limited information available, we estimate that the majority of ML-related workloads today are likely taking place in cloud and hyperscale data centers, with a smaller share occurring on distributed devices such as personal computers. Cloud and hyperscale data centers account for 0.1-0.2% of global GHG emissions [49–51], and it is likely that less than a quarter of their workloads and traffic are currently ML-related based on estimates for Infrastructure-as-a-Service and Platform-as-a-Service [52] and IP traffic related to big data [53]. Over the coming years, edge devices such as smartphones are also expected to handle an increasing volume of inference tasks to reduce latency and dependence on network connectivity [54], with uncertain effects on overall energy use and emissions.

While the amount of compute needed for each of the largest ML training runs is growing rapidly [22], it is uncertain how quickly overall ML-related energy use in data centers is increasing. For example, Facebook's overall data center energy use increased rapidly over the past few years (+40% per year) [55] while compute demand for ML training (e.g., +150% per year [56]) and inference (e.g., +105% per year [57]) have grown even faster. At the same time, by some measures, Facebook's operational GHG footprint (accounting for renewable energy purchases) halved between 2015 and 2019 [55] due in part to energy efficiency improvements and increased renewable electricity procurement.

Energy efficiency has played a central role in limiting data center energy demand growth more generally. Between 2010 and 2018, global data center energy use rose by only 6%, despite a 550% increase in workloads and compute instances [49]. The rapid growth in demand for data center services has been offset by efficiency improvements in servers, storage, networking, and infrastructure, as well as a shift away from smaller, less efficient data centers to large cloud and hyperscale data centers [18, 58], which have higher virtualization, more efficient cooling, and greater use of specialized "AI accelerator" hardware such as application-specific integrated circuits (ASICs) and graphics processing units (GPUs). For instance, a 2017 study found that Google's custom ASIC, the Tensor Processing Unit (TPU), was on average 30-80 times more energy efficient than a contemporary CPUs or GPUs [59]. However, the use of GPUs and ASICs for ML applications could drastically increase the power density of data center racks, which may in turn require liquid cooling technologies and increase water use. Although energy use across all data centers has been flat over the past decade, energy use by large data centers has grown by around 20% annually, and this trend is expected to continue [50]. Limiting overall data center energy demand growth over the next decade will therefore require even stronger energy efficiency improvements. For instance, operators can increase utilization and virtualization to maximize the energy efficiency of existing hardware and infrastructure, while replacing hardware when advisable from a life cycle perspective with the most efficient option. Companies and governments will also need to invest in research, development, and demonstration (RD&D) for efficient next-generation computing and communications technologies [49].

Some of the largest data center operators are now purchasing as much renewable electricity as they consume on a global annual basis [50], however, this does not guarantee that their data centers are actually fully powered by renewable sources all the time. More ambitious approaches to low-carbon electricity include shifting flexible workloads to times of the day (or locations) with higher shares of renewables generation [60]
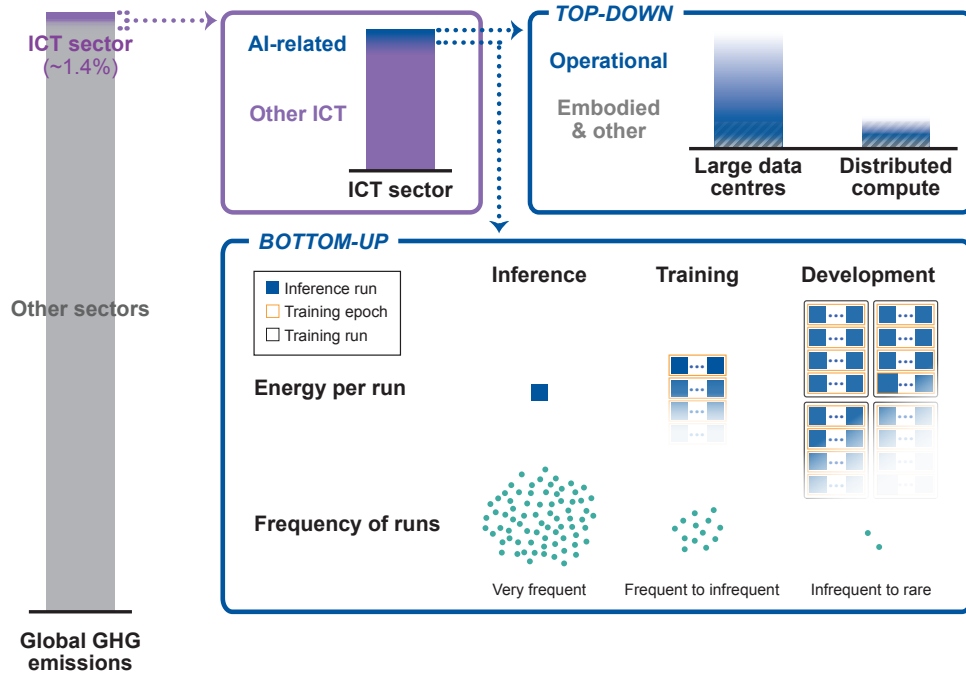
Figure 2: Compute-related global greenhouse gas (GHG) emissions impacts of machine learning (ML). The information and communications technology (ICT) sector accounts for around 1.4% of GHG emissions today, of which ML likely accounts for a small but unknown share (indicated by shading). Compute-related impacts of ML can be assessed from both top-down and bottom-up perspectives. Top-down: The majority of ML-related GHG emissions likely come from compute loads in large data centers, with a smaller share from distributed compute (e.g., personal computers, smartphones); these GHG emissions result from both operational energy use from computation and from other phases of the hardware life cycle (including embodied emissions). Bottom-up: The amount of energy needed for an ML model throughout different stages of the model life cycle differs based on problem setting and usage patterns.

and replacing on-site diesel generators with battery storage.

Computing hardware and infrastructure is also responsible for "embodied" emissions from raw materials extraction and manufacturing as well as emissions from transportation and end-of-life. For decentralized computing (e.g., desktops, laptops, smartphones), embodied emissions account for 40-80% of devices' life cycle GHG emissions, while for data centers this is typically less than 10% [48, 61–63]. Servers in large data centers are typically replaced every 3-4 years, which can result in higher operational efficiency [49, 64]; however, shorter lifespans could also increase the share of life cycle emissions from manufacturing, which can be mitigated by reusing servers and equipment (such as older GPUs for inference). As data centers become increasingly efficient and powered by clean electricity, the relative importance of emissions from non-operational life cycle phases will grow – particularly embodied emissions in computing hardware and data center building construction [65].

## Immediate application impacts

The broad applicability of ML algorithms means that they can be used both in applications that alleviate bottlenecks in addressing climate change, and in applications that may counteract climate action. In Rolnick et al. [5], we describe a number of settings in which ML can enable or accelerate climate change mitigation and adaptation strategies. As shown in Figure 3, these applications span areas such as energy, transportation,
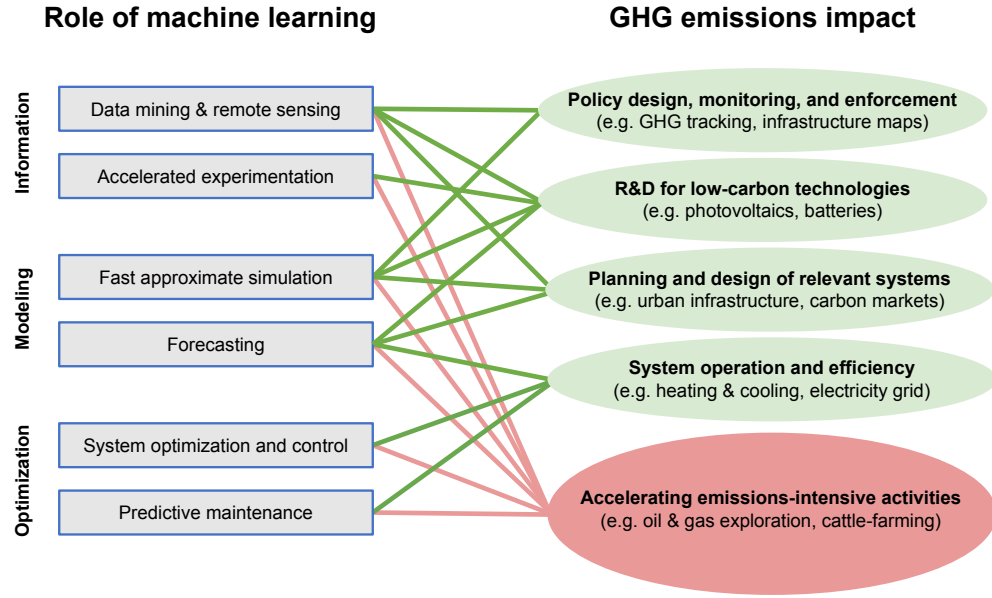
Figure 3: Immediate application impacts of machine learning (ML). ML applications are grouped by their functional role (left) and the associated greenhouse gas emissions impacts (right). ML can both reduce emissions (indicated in green) and increase emissions (red on bottom right). This figure differentiates ML applications for addressing climate change in more detail based on the findings in Rolnick et al. [5]; however, the net effect of those applications addressing climate change vs. those accelerating emissions-intensive industries is unclear.

and land use. For example, via data mining and remote sensing, ML can be used to translate raw data such as text documents or satellite imagery into usable insights for RD&D, policy-making, and systems planning – e.g., by tracking GHG emissions, fusing datasets, or gathering information on building efficiency characteristics. By accelerating the search for experimental parameters in scientific discovery, ML can help in designing next-generation batteries and solar cells. By learning from time series, ML can forecast renewable power production, crop yields, and transportation demands. By controlling and improving the operational efficiency of complex systems, e.g., industrial heating and cooling systems, ML can help to save resources and energy. ML can also be used to speed up time-intensive physics-based simulations, e.g., for urban planning or climate modeling. Predictive maintenance approaches leveraging ML can also help climate change mitigation if they are applied to low-carbon systems, where they can improve efficiency, reduce costs, and build resilience.

While ML is often seen as a "futuristic" technology, most of these applications are possible with current ML techniques, and many are already being deployed [1, 5]. In addition, areas of cutting-edge ML research such as interpretable and probabilistic machine learning [66, 67], physics-integrated machine learning [68], and transfer learning [69] can both enable new applications and better support integration within existing systems. In order to support the development and deployment of this kind of work, it will be crucial to facilitate interdisciplinary and applied research via science policy, advance the technological readiness of applications through RD&D programs, and adapt current regulatory environments to mitigate bottlenecks to deployment in relevant sectors and industries. This includes targeted funding and research programs, testbeds and demonstration projects, public procurement programs, and relevant data management initiatives.

As a general-purpose tool, ML has also been applied in ways that may make climate goals harder to achieve. One such effect is when ML is used to decrease the cost of emissions-intensive activities, thereby potentially increasing their consumption. For example, ML can accelerate oil and gas exploration and

extraction by decreasing production costs and boosting reserves [6], which could in turn lead to greater use of fossil fuels. Likewise, ML is used in the "Internet of Cows" to help manage livestock at scale [70], which can increase cattle-farming, an activity already responsible for an estimated 15% of GHG emissions [71]. Sometimes, essentially the same ML algorithm has the potential to be used both in ways that help climate action and in ways that hinder it. For example, an optimization algorithm that aims to minimize energy consumption (e.g., reduce the amount of transportation fuel used) would use essentially the same methodology as one that minimizes total cost, but the latter could potentially increase energy consumption if other cost aspects (e.g., labor costs) dominate. A potential approach to reduce or avoid the emissions increases associated with such applications is to require ML solutions providers to account for and report the emissions impacts of the applications they support, even if only at the level of order-of-magnitude or qualitative assessments where more detailed numbers are infeasible to obtain.

The total immediate impact of ML applications on GHG emissions is extremely difficult to estimate due to the lack of data on the deployment rate of ML, the diversity of application areas, and the lack of procedures to appropriately attribute emissions effects to the use of ML algorithms. While there exist some scientific reviews within isolated fields or sectors, the only attempts to provide overall numbers are from ML solutions providers in the private sector [72–74] (in particular, these studies are not peer-reviewed and do not disclose all methodology). We also note that ML can be used with the motivation to elevate the profile of sustainability-related activities in corporations in a way that could provide a false impression of overall organizational sustainability [10].

## System-level impacts

While the previous section describes ML applications that are directly beneficial or detrimental to climate change mitigation, many societal ML applications may not have clear immediate impacts on climate change. However, many of these applications can have broader societal implications beyond their immediate impact, and these system-level effects can influence GHG emissions both positively and negatively. Though these kinds of impacts may be hard to quantify, they have the potential to outweigh immediate application impacts and are extremely important to consider when evaluating ML use cases.

One pathway to system-level impacts occurs when ML enables changes to a technology that in turn affect the ways in which that technology is used. For example, rebound effects can occur when ML increases the efficiency of a service. While improved efficiency may result in lower GHG emissions per use, a decrease in cost may lead to increased consumption of the same or another good. This can eat into GHG benefits from efficiency gains or even counteract them [75]. Such rebound effects can be direct, for example by allowing a manufacturing plant to use ML-enabled efficiency gains to increase production of the same goods, thereby (partially) negating emissions savings. Even larger impacts can be expected from more structural types of rebound effects [10], which occur for example with ML-enabled autonomous driving. Specifically, autonomous vehicles can improve fuel efficiency, but they may also lead to higher rates of individualized vehicle travel, potentially increasing overall energy use and emissions if autonomous vehicles are not shared and/or electrified [76–79].

Given ML's role as an accelerator of technological development, it may also induce path dependencies that affect climate change mitigation. For instance, the phenomenon of "lock-in" refers to a scenario where a particular technology reaches market first and prevents competitors from entering the market [80]. Depending on how it is applied, ML may end up entrenching the role of a potentially inferior technology in a way that prevents other, e.g., low-carbon technologies, from entering the market. For example, the adoption of autonomous vehicles may ingrain the role of trucks and private cars as the dominant means for transportation, instead of enabling infrastructure and space for less emissions-intensive rail, public transit, and micromobility options [77]. On the other hand, ML may help break path-dependency effects or create a

first-mover advantage for a technology that is beneficial for the climate. The potential effects of ML on such path-dependencies in the context of climate change mitigation should be carefully analyzed.

Another avenue to system-level impact occurs when ML technologies influence broader lifestyle changes across society, for example by changing the demand for goods and services. A likely negative example here is in advertising, where ML algorithms such as recommender systems can be used to increase consumption of goods and services with embodied GHG emissions. Given that ML is fueled by data, its use could also incentivize increasingly large data infrastructures, which can come with their own carbon footprint and systemic implications. Various other paradigm-changing applications of ML have highly unclear effects from a climate perspective – for example, in automatic translation tools, virtual assistants, and augmented/virtual reality.

These examples demonstrate how important it is to assess the impacts of an ML application at the system level, rather than only estimate marginal effects, and to design public policy to shape system effects. Such policy levers include requiring climate impacts to be considered within regulations surrounding ML-driven emerging technologies [81], and implementing carbon pricing or other mechanisms to incentivize GHG emissions reductions and avoid rebound effects when ML is applied for efficiency.

# A roadmap for assessing and forecasting impacts

Above, we have discussed the extent to which it is currently possible to estimate the GHG emissions associated with ML. However, holistic and realistic predictions of ML's impact across several areas of our framework will require new reporting standards, more data collection, novel measurement methodologies, and new approaches for developing forecasts and scenarios. Moreover, given the heterogeneous nature of the capabilities and impacts of different digitalization technologies, ML and other forms of data analytics warrant separate consideration within impact assessment frameworks for digital technologies. Such efforts could, for example, build upon and extend existing methods and standards for life-cycle assessment of ICT (e.g., [82, 83]) to devise approaches that take that heterogeneity into account. We call on the academic fields of life-cycle analysis, industrial ecology, and others to actively grapple with this task, and have designed the framework provided in this paper (Figure 1) to lay the groundwork for such methodological development.

To estimate compute-related impacts (Category A of our framework), better access to information will be central. For example, while it is relatively straightforward to estimate the compute-related GHG emissions resulting from individual runs of AI systems, the usage patterns in practice are typically opaque. Practitioners could disclose relevant information such as specifics on computing power needed for system development, training, fine-tuning, and inference at appropriate temporal resolutions, as well as information about the type, time, and location of computing infrastructure used. Also informative are specifics about the model type and size; training requirements for model development (or pre-trained models used); frequency of training, re-training, and fine-tuning; and average number of inference uses per unit of time. From a top-down perspective, another important datapoint is the share of the total compute load in data centers that can be attributed to ML, ideally distinguished by the relevant model life-cycle stages. This information would allow for a top-down estimate of global compute-related impacts and underlying dynamics, but is currently not made public by data center operators.

There are currently limited quantitative estimates available about the immediate impacts of ML applications (Category B). The lack of established methodology poses a central bottleneck here. Research and practice need to establish how to estimate the marginal and counterfactual benefit that ML could have if introduced in established processes, including distinguishing between use cases that would not exist without ML vs. those where ML provides improvements to an existing use case. For such efforts, it will be important to develop a more fine-grained taxonomy of ML systems and application areas, that can help to generalize beyond single case studies and also help stakeholders weigh costs and benefits of new projects *a priori*. Espe-

cially regarding ML application effects, also obtaining better data will be difficult yet crucial, considering the potentially large magnitude and uncertainty around those developments. To estimate impacts more broadly and systematically, reviewing, synthesizing, and generalizing case studies will be important, and where data cannot be easily obtained, approaches such as stakeholder surveys or expert elicitation might help to fill gaps.

Perhaps the most important yet most difficult to assess are the system-level impacts (Category C). ML is a fast-growing enabling technology that has the potential to affect present and future societal and technological trajectories and thus needs to be appropriately accounted for in forecasting and scenario analysis. ML can influence many input factors of climate and energy system models, such as efficiency, production and consumption rates, learning rates, resource constraints, financial assumptions, etc., which makes ML a "wild card" that could introduce large transformations in different ways. How to appropriately factor that uncertainty into climate and energy system models is yet to be established. Importantly, ML builds on digital infrastructure, yet the impact assessment of digitalization is itself at an early stage (especially when it comes to estimating the impacts of how digital technologies are applied) [82, 83]. Energy and climate models, such as energy system models developed by the International Energy Agency (IEA) and U.S. Information Administration (EIA) or the Shared Socioeconomic Pathways used for the IPCC, generally do not explicitly or systematically account for digitalization, let alone the effects introduced by ML. One exception is perhaps the inclusion of autonomous vehicles in scenarios, e.g. by the EIA [79]. Our framework can be used as a starting point, and is sufficiently general to provide a comprehensive framing for incorporating current and future ML effects within scenario analysis.

## Approaches for aligning ML with climate change mitigation

Given ML's multi-faceted relationship with climate change, many different kinds of approaches from the public and private sectors are needed to shape its impacts. This will require progress in both climate policy and AI policy, coupled with algorithmic and hardware innovation and the development of adequate impact assessment methodologies. Table 1 provides an overview of a number of these strategies.

While not addressing ML explicitly, general climate policy approaches, such as carbon pricing, may be effective in driving the development and use of ML in a manner that is aligned with climate change mitigation. Science policy approaches that foster low-carbon technologies may also facilitate uses of ML that enable or improve these technologies (though they may not necessarily address ML-specific barriers).

To address more technology-specific opportunities and risks, it will be important for climate change to become a major consideration within AI innovation and deployment policies. This includes (a) promoting the research, development, and deployment of ML applications that are beneficial to the climate, (b) requiring transparency and accountability for those use cases that could increase emissions or otherwise counteract climate change goals, as well as on computational energy use, and (c) employing climate-cognizant technology assessment for ML use cases that are not traditionally within the realm of climate policy, but where decisions today may have large implications for future climate impacts. Many of the associated policy approaches in Table 1 can be developed and implemented starting today.

Further, mandating emissions measurement and reporting for ML use cases – considering the impacts of both compute *and* applications – can enable these emissions to be regulated via climate policy approaches, and further shape the design of targeted policies. Such reporting requirements, however, need to be carefully designed based on what is feasible to estimate given the state of the measurement research and bureaucratic burden, based on an understanding of where top-down measurement might suffice to inform regulatory approaches, and with an eye towards preventing strategic behavior such as the "hiding" of emissions in cloud compute servers [84]. Climate-related reporting for ML-based systems can potentially be more easily implemented where other AI reporting requirements are planned or in place (such as proposed in the EU [85]).

Table 1: Levers to reduce the greenhouse gas (GHG) emissions impacts associated with machine learning (ML) compute and applications.

| Lever type | Compute-related (algorithm, infrastructure) | Application-related (immediate, systemic) |
|---|---|---|
| **Public sector** | | |
| **Economic instruments** | ● *Implement economy-wide or sector-specific carbon pricing to incentivize emissions reductions and mitigate rebound effects* | |
| **Research, development & demonstration (RD&D)** | ● Support research in energy-efficient ML<br>● Support RD&D in energy-efficient, specialized, and low-resource hardware<br>● Support RD&D in data center operational efficiency | ● Support interdisciplinary and applied ML research for climate-relevant applications of ML<br>◑ Provide mechanisms to advance the technological readiness of climate-beneficial AI applications (e.g. testbeds, demonstration projects, public procurement programs) |
| **Regulation** | ○ Employ a climate-cognizant technology assessment lens within AI strategies and when regulating ML-driven emerging technologies<br>● *Implement clean electricity mandates (e.g., low-carbon portfolio standards)*<br>● Implement efficiency standards for data center hardware and infrastructure | ● *Employ regulatory approaches to constrain sector-specific GHG emissions*<br>◑ Reduce deployment barriers in relevant sectors and industries for AI applications that are beneficial to the climate |
| **Best practices and standards** | ◑ Develop interoperability standards for commercial ML approaches to prevent lock-in to particular solutions providers and facilitate a decentralized solutions provider space<br>● Develop and implement standardized metrics for evaluating model efficacy that include energy efficiency | ● Implement data governance standards that spur impactful work and are mindful of privacy and ownership<br>● Require meaningful civic and stakeholder engagement in scoping, developing, and deploying ML-driven projects<br>○ Develop best practices and systematic approaches to weigh benefits and costs for ML applications |
| **Monitoring and reporting** | ◑ Develop measurement methodologies and guidance to estimate and report ML-related GHG emissions<br>◑ Mandate appropriate life-cycle transparency and reporting of GHG emissions for ML use cases, including both compute and application-related impacts | |
| **Capacity-building** | ◑ Build in-house public-sector capacity in ML to facilitate governance and deployment<br>● Promote ML education and literacy among climate-relevant entities and in the public sector | ◑ Incentivize ML workforce shifts towards climate-oriented entities (e.g., via placement programs) |
| **Private sector** | | |
| **Corporate climate action** | ◑ Adopt organizational carbon pricing strategies that account for both compute- and application-related emissions (e.g., Scope 1, 2, and 3 emissions, including from cloud compute, as well as from products and services)<br>● Reduce wasteful model re-training and execution<br>● Make energy efficiency a central criterion in evaluating model efficacy<br>● Reduce GHG emissions across supply chains and product life cycle (including embodied emissions)<br>● Maximize energy efficiency in data centers and support related RD&D<br>● Shift compute load to geographies and times with lower carbon-intensity of the grid<br>● Purchase low-carbon electricity and invest in energy technologies to decarbonize the grid<br>○ Develop standardized ML platforms to facilitate rapid company-wide adoption of energy efficiency improvements | ● Adjust business models to avoid ML applications that drive GHG emissions increases<br>● Encourage ML applications that drive GHG emissions reductions<br>◑ Measure and engage in voluntary reporting of the emissions impacts of ML products and services |

**Legend**

- ● Policies that are ready to implement or already exist
- ◑ Policies that can be developed today
- ○ More analysis needed to develop policies
- *General climate policy levers*

Lastly, we note that ML expertise today is often concentrated among a limited set of actors, raising potential challenges with respect to the governance and implementation of ML in the context of climate change. For instance, the use of ML in certain contexts may yield or exacerbate societal inequities, e.g., by widening the digital divide [32, 86], through algorithmic bias [87], or by shifting power from public to large private entities by virtue of who controls relevant data or intellectual capital. Strategies to address such gaps include strengthening small and medium-sized ML solutions providers, developing incentives such as placement programs and dedicated education to shift the ML workforce towards public and climate-relevant entities, developing interoperability standards to prevent lock-in to particular solutions providers, and developing best practices for when state-of-the-art ML models vs. other (simpler) alternatives should be used. Developing meaningful civic engagement processes for the scoping, design, and deployment of projects (and associated data collection and provision efforts) will also be critical to ensuring that ML approaches are both effective and avoid potential pitfalls [88].

ML's ultimate effect on the climate is far from predestined, and societal decisions will play a large role in shaping its overall impacts. This will require a holistic portfolio of approaches across policy, industry, and academia to incentivize uses of ML that support climate change strategies while mitigating the impacts of use cases that may counteract climate change goals. Most importantly, society cannot wait to act: with the rapidly growing prevalence of ML and the increasing urgency of climate change, now presents a critical window of opportunity to shape ML's impacts for decades to come.

# Acknowledgements

# Author Information

## Contributions

P.L.D., L.H.K., and D.R. conceived the idea for this manuscript. All authors wrote and edited the manuscript text and figures, with primary contributions from E.S. and G.K. for the section on compute-related impacts, from D.R., F.C., and L.H.K. for the sections on application-related impacts, from P.L.D. for the section on shaping ML's impacts, and from L.H.K. for the introduction, roadmap for assessing impacts, and overall conceptual framing.

## Disclaimer

Opinions are G.K.'s own and do not reflect those of the OECD, the IEA or their Member countries.

# References

[1] Daniel Zhang, S Mishra, E Brynjolfsson, J Etchemendy, D Ganguli, B Grosz, T Lyons, J Manyika, JC Niebles, M Sellitto, et al. Artificial intelligence index report 2021. Technical report, Tech. rep., AI Index Steering Committee, Human-Centered AI Institute . . . , 2021.

[2] The Royal Society. Digital technology and the planet: Harnessing computing to achieve net zero. Technical report, The Royal Society, 2020. URL `https://royalsociety.org/-/media/policy/projects/digital-technology-and-the-planet/digital-technology-and-the-planet-report.pdf`.

[3] Lynn H Kaack, Priya L Donti, Emma Strubell, and David Rolnick. Artificial Intelligence and Climate Change: Opportunities, considerations, and policy levers to align AI with climate change goals, 2020. URL `https://eu.boell.org/en/2020/12/03/artificial-intelligence-and-climate-change`.

[4] World Economic Forum. Harnessing Artificial Intelligence to Accelerate the Energy Transition, 2021. URL `https://www.weforum.org/whitepapers/harnessing-artificial-intelligence-to-accelerate-the-energy-transition`.

[5] David Rolnick, Priya L Donti, Lynn H Kaack, Kelly Kochanski, Alexandre Lacoste, Kris Sankaran, Andrew Slavin Ross, Nikola Milojevic-Dupont, Natasha Jaques, Anna Waldman-Brown, et al. Tackling climate change with machine learning, 2019. URL `https://arxiv.org/abs/1906.05433`. arXiv:1906.05433.

[6] Greenpeace. Oil in the Cloud: How Tech Companies are Helping Big Oil Profit from Climate Destruction, 2019. URL `https://www.greenpeace.org/usa/reports/oil-in-the-cloud/`.

[7] Roel Dobbe and Meredith Whittaker. AI and climate change: How they're connected, and what we can do about it. *AI Now Institute, Medium, October*, 17, 2019. URL `https://medium.com/@AINowInstitute/ai-and-climate-change-how-theyre-connected-and-what-we-can-do-about-it-6aa8d0f5b32c`.

[8] Emma Strubell, Ananya Ganesh, and Andrew McCallum. Energy and policy considerations for deep learning in NLP. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3645–3650, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1355.

[9] Roy Schwartz, Jesse Dodge, Noah A. Smith, and Oren Etzioni. Green AI. *CoRR*, abs/1907.10597, 2019. URL `http://arxiv.org/abs/1907.10597`.

[10] Peter Dauvergne. Is artificial intelligence greening global supply chains? Exposing the political economy of environmental costs. *Review of International Political Economy*, (0):1–23. doi: 10.1080/09692290.2020.1814381.

[11] Mark Coeckelbergh. AI for climate: freedom, justice, and other ethical and political challenges. *AI and Ethics*, pages 1–6, 2020.

[12] Henry Gunther and Julietta Rose. Governing AI: The Importance of Environmentally Sustainable and Equitable Innovation. *The Environmental Law Reporter*, 10888, 2020.

[13] Amy L Stein. Artificial intelligence and climate change. *Yale J. on Reg.*, 37:890, 2020.

[14] Josh Cowls, Andreas Tsamados, Mariarosaria Taddeo, and Luciano Floridi. The AI gambit—leveraging artificial intelligence to combat climate change: Opportunities, challenges, and recommendations. 2021. URL `https://ssrn.com/abstract=3804983`. Available at SSRN.

[15] Lorenz M Hilty and Bernard Aebischer. ICT for Sustainability: An Emerging Research Field. In Lorenz M Hilty and Bernard Aebischer, editors, *ICT Innovations for Sustainability*, pages 3–36. Springer, 2015.

[16] Nathaniel C Horner, Arman Shehabi, and Inês L Azevedo. Known unknowns: Indirect energy effects of information and communication technology. *Environmental Research Letters*, 11(10):103001, 2016.

[17] Johanna Pohl, Lorenz M Hilty, and Matthias Finkbeiner. How LCA contributes to the environmental assessment of higher order effects of ICT application: A review of different approaches. *Journal of Cleaner Production*, 219:698–712, 2019.

[18] International Energy Agency. Digitalization & Energy. Technical report, OECD/IEA, Paris, 2017. URL `https://www.iea.org/reports/digitalisation-and-energy`.

[19] Varun Sivaram, Stephen D Comello, David G Victor, Lidija Sekaric, Ben Hertz-Hagel, Peter Fox-Penner, Rohit T Aggarwalla, Kyle Bradbury, Sunil Garg, Erfan Ibrahim, Jesse Scott, John O'Leary, Richard Kauffman, and Hiang Kwee Ho. *Digital Decarbonization Promoting Digital Innovations to Advance Clean Energy Systems*. 2018. URL `https://www.cfr.org/report/digital-decarbonization`.

[20] Charlie Wilson, Laurie Kerr, Frances Sprei, Emilie Vrain, and Mark Wilson. Potential climate benefits of digital consumer innovations. *Annual Review of Environment and Resources*, 45:113–144, 2020.

[21] Alfredo Canziani, Adam Paszke, and Eugenio Culurciello. An analysis of deep neural network models for practical applications, 2017. URL `https://arxiv.org/abs/1605.07678`. arXiv:1605.07678.

[22] OpenAI. AI and Compute, 2018. URL `https://openai.com/blog/ai-and-compute/`.

[23] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. 2021. URL `https://arxiv.org/abs/2108.07258`. arXiv:2108.07258.

[24] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain, April 2017. Association for Computational Linguistics. URL `https://www.aclweb.org/anthology/E17-2068`.

[25] K. Hazelwood, S. Bird, D. Brooks, S. Chintala, U. Diril, D. Dzhulgakov, M. Fawzy, B. Jia, Y. Jia, A. Kalro, J. Law, K. Lee, J. Lu, P. Noordhuis, M. Smelyanskiy, L. Xiong, and X. Wang. Applied Machine Learning at Facebook: A Datacenter Infrastructure Perspective. In *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 620–629, 2018. doi: 10.1109/HPCA.2018.00059.

[26] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, April 2018.

[27] Barak Turovsky. Ten years of Google Translate, 2016. URL `https://blog.google/products/translate/ten-years-of-google-translate/`.

[28] Paras Jain, Xiangxi Mo, Ajay Jain, Alexey Tumanov, Joseph E. Gonzalez, and Ion Stoica. The OoO VLIW JIT Compiler for GPU Inference, 2019. URL `https://arxiv.org/abs/1901.10008`. arXiv:1901.10008.

[29] Angela H. Jiang, Daniel L.-K. Wong, Giulio Zhou, D. Andersen, J. Dean, G. Ganger, Gauri Joshi, M. Kaminsky, Michael A. Kozuch, Zachary Chase Lipton, and P. Pillai. Accelerating deep learning by focusing on the biggest losers, 2019. URL https://arxiv.org/abs/1910.00762. arXiv:1910.00762.

[30] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. doi: 10.1109/ CVPR.2016.90.

[31] Samuel Albanie. convnet-burden: Estimates of memory consumption and flop counts for various convolutional neural networks. URL https://github.com/albanie/convnet-burden.

[32] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *FAccT*, 2021.

[33] Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. Towards the systematic reporting of the energy and carbon footprints of machine learning. *Journal of Machine Learning Research*, 21(248):1–43, 2020.

[34] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. ICML Workshop on Challenges in Deploying and monitoring Machine Learning Systems, December 2019. URL https://arxiv.org/abs/1910.09700. arXiv:1910.09700.

[35] Lasse F. Wolff Anthony, Benjamin Kanding, and Raghavendra Selvan. Carbontracker: Tracking and predicting the carbon footprint of training deep learning models. ICML Workshop on Challenges in Deploying and monitoring Machine Learning Systems, July 2020. URL https://arxiv.org/abs/2007. 03051. arXiv:2007.03051.

[36] Ermao Cai, Da-Cheng Juan, Dimitrios Stamoulis, and Diana Marculescu. NeuralPower: Predict and deploy energy-efficient convolutional neural networks. In *The 9th Asian Conference on Machine Learning (ACML)*, 2017.

[37] Jesse Dodge, Suchin Gururangan, Dallas Card, Roy Schwartz, and Noah A. Smith. Show your work: Improved reporting of experimental results. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2185–2194, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1224.

[38] Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, David Brooks, Dehao Chen, Debo Dutta, Udit Gupta, Kim Hazelwood, Andy Hock, Xinyuan Huang, Daniel Kang, David Kanter, Naveen Kumar, Jeffery Liao, Deepak Narayanan, Tayo Oguntebi, Gennady Pekhimenko, Lillian Pentecost, Vijay Janapa Reddi, Taylor Robie, Tom St John, Carole-Jean Wu, Lingjie Xu, Cliff Young, and Matei Zaharia. Mlperf training benchmark. In I. Dhillon, D. Papailiopoulos, and V. Sze, editors, *Proceedings of Machine Learning and Systems*, volume 2, pages 336–349, 2020.

[39] V. J. Reddi, C. Cheng, D. Kanter, P. Mattson, G. Schmuelling, C. Wu, B. Anderson, M. Breughe, M. Charlebois, W. Chou, R. Chukka, C. Coleman, S. Davis, P. Deng, G. Diamos, J. Duke, D. Fick, J. S. Gardner, I. Hubara, S. Idgunji, T. B. Jablin, J. Jiao, T. S. John, P. Kanwar, D. Lee, J. Liao, A. Lokhmotov, F. Massa, P. Meng, P. Micikevicius, C. Osborne, G. Pekhimenko, A. T. R. Rajan, D. Sequeira, A. Sirasao, F. Sun, H. Tang, M. Thomson, F. Wei, E. Wu, L. Xu, K. Yamada, B. Yu, G. Yuan, A. Zhong, P. Zhang, and Y. Zhou. MLPerf Inference Benchmark. In *2020 ACM/IEEE*

*47th Annual International Symposium on Computer Architecture (ISCA)*, pages 446–459, 2020. doi: 10.1109/ISCA45697.2020.00045.

[40] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In *NeurIPS Deep Learning Workshop*, 2014.

[41] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. In *International Conference on Learning Representations*, 2016.

[42] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *Journal of Machine Learning Research*, 18(187):1–30, 2018. URL `http://jmlr.org/papers/v18/16-456.html`.

[43] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 46–54, 2020.

[44] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. In *International Conference on Learning Representations*, 2020.

[45] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. {GS}hard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021.

[46] Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. Characterizing and mitigating bias in compact models. In *ICML Workshop on Human Interpretability in Machine Learning (WHI)*, 2020.

[47] International Telecommunication Union. Greenhouse gas emissions trajectories for the information and communication technology sector compatible with the UNFCCC Paris Agreement. *ITU-T*, 2020. URL `http://handle.itu.int/11.1002/1000/14084`.

[48] Jens Malmodin and Dag Lundén. The energy and carbon footprint of the global ICT and E&M sectors 2010–2015. *Sustainability*, 10(9):3027, 2018.

[49] Eric Masanet, Arman Shehabi, Nuoa Lei, Sarah Smith, and Jonathan Koomey. Recalibrating global data center energy-use estimates. *Science*, 367(6481):984–986, 2020. ISSN 0036-8075. doi: 10.1126/science.aba3758.

[50] International Energy Agency. Data Centres and Data Transmission Networks, 2020. URL `https://www.iea.org/reports/data-centres-and-data-transmission-networks`.

[51] Francesca Montevecchi, Therese Stickler, Ralph Hintemann, and Simon Hinterholzer. Energy-efficient Cloud Computing Technologies and Policies for an Eco-friendly Cloud Market, 2020. URL `https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=71330`.

[52] Cisco. Cisco Global Cloud Index: Forecast and Methodology, 2016–2021, 2018. URL `https://www.cisco.com/c/en/us/solutions/collateral/service-provider/global-cloud-index-gci/white-paper-c11-738085.pdf`.

[53] Cisco. Cisco's Global Cloud Index Study: Acceleration of the Multicloud Era, 2018. URL `https://blogs.cisco.com/news/acceleration-of-multicloud-era`.

[54] Carole-Jean Wu, David Brooks, Kevin Chen, Douglas Chen, Sy Choudhury, Marat Dukhan, Kim Hazelwood, Eldad Isaac, Yangqing Jia, Bill Jia, et al. Machine learning at Facebook: Understanding inference at the edge. In *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 331–344. IEEE, 2019.

[55] Facebook. Facebook Sustainability Data 2019, 2020. URL `https://sustainability.fb.com/asset/fb_sustainability-data-disclosure-2019/`.

[56] Maxim Naumov, John Kim, Dheevatsa Mudigere, Srinivas Sridharan, Xiaodong Wang, Whitney Zhao, Serhat Yilmaz, Changkyu Kim, Hector Yuen, Mustafa Ozdal, Krishnakumar Nair, Isabel Gao, Bor-Yiing Su, Jiyan Yang, and Mikhail Smelyanskiy. Deep Learning Training in Facebook Data Centers: Design of Scale-up and Scale-out Systems, 2020. URL `https://arxiv.org/abs/2003.09518`. arXiv:2003.09518.

[57] Jongsoo Park, Maxim Naumov, Protonu Basu, Summer Deng, Aravind Kalaiah, Daya Khudia, James Law, Parth Malani, Andrey Malevich, Satish Nadathur, Juan Pino, Martin Schatz, Alexander Sidorov, Viswanath Sivakumar, Andrew Tulloch, Xiaodong Wang, Yiming Wu, Hector Yuen, Utku Diril, Dmytro Dzhulgakov, Kim Hazelwood, Bill Jia, Yangqing Jia, Lin Qiao, Vijay Rao, Nadav Rotem, Sungjoo Yoo, and Mikhail Smelyanskiy. Deep Learning Inference in Facebook Data Centers: Characterization, Performance Optimizations and Hardware Implications, 2018. URL `https://arxiv.org/abs/1811.09886`. arXiv:1811.09886.

[58] Arman Shehabi, Sarah Smith, Dale Sartor, Richard Brown, Magnus Herrlin, Jonathan Koomey, Eric Masanet, Nathaniel Horner, Inês Azevedo, and William Lintner. United States Data Center Energy Usage Report. 2016. URL `https://eta.lbl.gov/publications/united-states-data-center-energy`.

[59] Norman P. Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, Diemthu Le, Chris Leary, Zhuyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Horia Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, and Doe Hyun Yoon. In-datacenter performance analysis of a tensor processing unit. *SIGARCH Comput. Archit. News*, 45(2):1–12, June 2017. ISSN 0163-5964. doi: 10.1145/3140659.3080246.

[60] Ana Radovanovic. Our data centers now work harder when the sun shines and wind blows, Apr 2020. URL `https://blog.google/inside-google/infrastructure/data-centers-work-harder-sun-shines-wind-blows`.

[61] Beth Whitehead, Deborah Andrews, and Amip Shah. The life cycle assessment of a UK data centre. *The International Journal of Life Cycle Assessment*, 20(3):332–349, 2015.

[62] Eric Masanet, Arman Shehabi, and Jonathan Koomey. Characteristics of low-carbon data centres. *Nature Climate Change*, 3(7):627–630, 2013.

[63] Roland Hischier, Vlad C Coroama, Daniel Schien, and Mohammad Ahmadi Achachlouei. Grey energy and environmental impacts of ICT hardware. In *ICT Innovations for Sustainability*, pages 171–189. Springer, 2015.

[64] Luiz André Barroso, Jimmy Clidaras, and Urs Hölzle. The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines. *Synthesis Lectures on Computer Architecture*, 8(3):1–154, 2013.

[65] Udit Gupta, Young Geun Kim, Sylvia Lee, Jordan Tse, Hsien-Hsin S. Lee, Gu-Yeon Wei1, David Brooks, and Carole-Jean Wu. Chasing Carbon: The Elusive Environmental Footprint of Computing". In *IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2021.

[66] Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. Interpretable machine learning: Fundamental principles and 10 grand challenges, 2021. URL `https://arxiv.org/abs/2103.11251`. arXiv:2103.11251.

[67] Zoubin Ghahramani. Probabilistic machine learning and artificial intelligence. *Nature*, 521(7553):452–459, 2015.

[68] Jared Willard, Xiaowei Jia, Shaoming Xu, Michael Steinbach, and Vipin Kumar. Integrating physics-based modeling with machine learning: A survey. 2020. URL `https://arxiv.org/abs/2003.04919`. arXiv:2003.04919.

[69] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020.

[70] Amber Adams-Progar, Glenn A. Fink, Ely Walker, and Don Llewellyn. *Security and Privacy Issues in the Internet of Cows*, chapter 18, pages 375–398. John Wiley & Sons, Ltd, 2017. ISBN 9781119226079. doi: https://doi.org/10.1002/9781119226079.ch18.

[71] Pierre J Gerber, Henning Steinfeld, Benjamin Henderson, Anne Mottet, Carolyn Opio, Jeroen Dijkman, Allessandra Falcucci, Giuseppe Tempio, et al. *Tackling climate change through livestock: a global assessment of emissions and mitigation opportunities*. Food and Agriculture Organization of the United Nations (FAO), 2013.

[72] C Herweijer, B Combes, and J Gillham. How AI can enable a sustainable future. *Microsoft and PWC: London, UK*, 2018. URL `https://www.pwc.co.uk/services/sustainability-climate-change/insights/how-ai-future-can-enable-sustainable-future.html`.

[73] Capgemini. Climate AI: How artificial intelligence can power your climate action strategy, 2020. URL `https://www.capgemini.com/research/climate-ai/`.

[74] Charlotte Degot, Sylvain Duranton, Michel Frédeau, and Rich Hutchinson. Reduce Carbon and Costs with the Power of AI, January 2021. URL `https://www.bcg.com/en-us/publications/2021/ai-to-reduce-carbon-emissions`.

[75] Inês M.L. Azevedo. Consumer end-use energy efficiency and rebound effects. *Annual Review of Environment and Resources*, 39(1):393–418, 2014. doi: 10.1146/annurev-environ-021913-153558.

[76] James M. Anderson, Nidhi Kalra, Karlyn D. Stanley, Paul Sorensen, Constantine Samaras, and Tobi A. Oluwatola. *Autonomous Vehicle Technology: A Guide for Policymakers*. RAND Corporation, Santa Monica, CA, 2016. doi: 10.7249/RR443-2.

[77] Felix Creutzig, Martina Franzen, Rolf Moeckel, Dirk Heinrichs, Kai Nagel, Simon Nieland, and Helga Weisz. Leveraging digitalization for sustainability in urban transport. *Global Sustainability*, 2, 2019.

[78] Zia Wadud, Don MacKenzie, and Paul Leiby. Help or hindrance? The travel, energy and carbon impacts of highly automated vehicles. *Transportation Research Part A: Policy and Practice*, 86:1–18, 2016.

[79] Nicholas Chase, John Maples, and Mark Schipper. Autonomous Vehicles: Uncertainties and Energy Implications, May 2018. URL https://www.eia.gov/outlooks/aeo/av.php.

[80] W. Brian Arthur. Competing technologies, increasing returns, and lock-in by historical events. *The Economic Journal*, 99(394):116–131, 1989. ISSN 00130133, 14680297.

[81] E Cox, S Royston, and J Selby. Impact of non-energy policies on energy systems. *Retrieved from UK Energy Research Centre (UKERC): http://www. ukerc. ac. uk/asset/1B9BBB2F-B98C-4250-BEE5DE0F253EAD91*, 2016.

[82] Vlad C Coroamă, Pernilla Bergmark, Mattias Höjer, and Jens Malmodin. A methodology for assessing the environmental effects induced by ict services: Part i: Single services. In *Proceedings of the 7th International Conference on ICT for Sustainability*, pages 36–45, 2020.

[83] Pernilla Bergmark, Vlad C Coroamă, Mattias Höjer, and Craig Donovan. A methodology for assessing the environmental effects induced by ict services: Part ii: Multiple services and companies. In *Proceedings of the 7th International Conference on ICT for Sustainability*, pages 46–55, 2020.

[84] David Mytton. Hiding greenhouse gas emissions in the cloud. *Nature Climate Change*, 10(8):701–701, 2020. doi: 10.1038/s41558-020-0837-6.

[85] European Commission. Proposal for a Regulation of the European Parliament and of the Council Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act) and Amending Certain Union Legislative Acts, 2021. URL https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206.

[86] Martin Hilbert. Big data for development: A review of promises and challenges. *Development Policy Review*, 34(1):135–174, 2016. doi: https://doi.org/10.1111/dpr.12142.

[87] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. 2019. URL https://arxiv.org/abs/1908.09635. arXiv:1908.09635.

[88] Elizabeth Bondi, Lily Xu, Diana Acosta-Navas, and Jackson A Killian. Envisioning Communities: A Participatory Approach Towards AI for Social Good. 2021. URL https://arxiv.org/abs/2105.01774. arXiv:2105.01774.