# Carbon Aware Transformers Through Joint Model-Hardware Optimization

Irene Wang $^{1,2,*}$ , Newsha Ardalani $^1$ , Mostafa Elhoushi $^1$ , Daniel Jiang $^1$ , Samuel Hsia $^1$ , Ekin Sumbul $^1$ , Divya Mahajan $^2$ , Carole-Jean Wu $^1$ , Bilge Acun $^1$ 

The rapid growth of machine learning (ML) systems necessitates a more comprehensive evaluation of their environmental impact, particularly their carbon footprint, which comprises operational carbon from training and inference execution and embodied carbon from hardware manufacturing and its entire life-cycle. Despite the increasing importance of embodied emissions, there is a lack of tools and frameworks to holistically quantify and optimize the total carbon footprint of ML systems. To address this, we propose CATransformers, a carbon-aware architecture search framework that enables sustainability-driven co-optimization of ML models and hardware architectures. By incorporating both operational and embodied carbon metrics into early design space exploration of domain-specific hardware accelerators, CATransformers demonstrates that optimizing for carbon yields design choices distinct from those optimized solely for latency or energy efficiency. We apply our framework to multi-modal CLIP-based models, producing CarbonCLIP, a family of CLIP models achieving up to 17% reduction in total carbon emissions while maintaining accuracy and latency compared to state-of-the-art edge small CLIP baselines. This work underscores the need for holistic optimization methods to design high-performance, environmentally sustainable AI systems.

**Date:** May 12, 2025

Correspondence: Irene Wang at irene.wang@gatech.edu & Bilge Acun at acun@meta.com

Code: https://github.com/facebookresearch/CATransformers

Meta

#### 1 Introduction

As machine learning (ML) systems become more widespread across various industries, it's crucial to take a closer examination of their carbon footprint and find strategies to mitigate it across the stack. This requires taking a holistic view and considering both operational carbon, incurred from the energy use during model training and inference, and embodied carbon, associated with the manufacturing and life-cycle of hardware Gupta et al. (2022). This work tackles the question: "How does incorporating carbon footprint metrics into optimization workflows influence the design of ML models and hardware architectures?"

A fundamental challenge in designing sustainable ML systems is that model and hardware decisions are deeply interdependent. Effective co-optimization serves two key purposes. First, in the early stages of hardware accelerator design, identifying efficient architectures for executing target models is crucial. By co-optimizing model architectures alongside domain-specific hardware accelerators, we can determine optimal configurations that improve efficiency and reduce emissions while maintaining accuracy. Second, optimizing for total carbon requires jointly considering model and hardware architecture, where model execution on a particular hardware determines the operational carbon and the hardware area and design determines the embodied carbon.

This is particularly crucial for multi-modal models like CLIP Radford et al. (2021), which are computationally intensive and resource-heavy. However, existing hardware-aware neural architecture search (NAS) methods Wang et al. (2020); Zhou et al. (2022) primarily optimize for latency and energy, neglecting embodied carbon and failing to address the heterogeneous nature of multi-modal models. Unlike uni-modal architectures, multi-modal systems introduce distinct computational bottlenecks—e.g., vision transformers rely on memory bandwidth and parallelism Marino et al. (2023), whereas text transformers require sequential processing Jain et al. (2024). Balancing these diverse workloads for both efficiency and sustainability expands the design space,

<sup>&</sup>lt;sup>1</sup>FAIR at Meta, <sup>2</sup>Georgia Institute of Technology

<sup>\*</sup>Work done at Meta

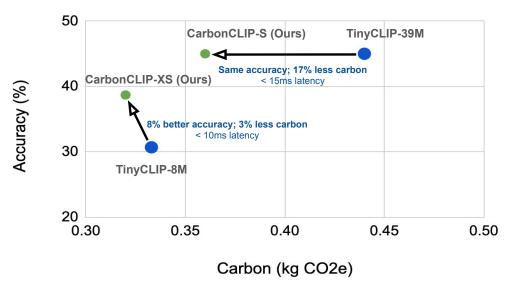


Figure 1 CarbonCLIP models achieves lower carbon footprint and higher accuracy compared to baseline CLIP models.

making carbon-aware co-optimization even more challenging. Effective co-optimization requires accurate quantification of both operational and embodied carbon, yet existing tools lack the capability to assess total carbon footprint—particularly for multi-modal models on custom accelerators designs.

To address these challenges, we propose CATransformers, a carbon-aware neural and hardware architecture search framework that incorporates both operational and embodied carbon into the optimization process to design systems for edge inference-only devices. Our key insight is that traditional optimization objectives—such as minimizing latency or energy—often lead to design choices that do not necessarily minimize total carbon emissions. We demonstrate that carbon-aware optimization yields different architectural trade-offs compared to conventional approaches, necessitating new tools and methodologies.

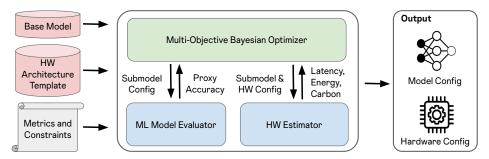
CATransformers, systematically explores the joint design space of model architectures and hardware accelerators using multi-objective Bayesian optimization. It consists of three key components: A *Multi-Objective Bayesian Optimizer* that balances trade-offs between accuracy, latency, energy, and carbon footprint, an *ML Model Evaluator* that efficiently explores model architecture variations via importance-based pruning and fine-tuning strategies, and a *Hardware Estimator* that quantifies latency, energy, and carbon emissions using open-source profiling tools.

Applying CATransformers to multi-modal CLIP-based models, we derive CarbonCLIP, a family of models that achieve up to 17% lower total carbon emissions compared to edge-deployed CLIP baselines, while maintaining comparable accuracy and latency. Figure 1 highlights that hardware-aware model co-optimization can produce models that are both low-carbon and high-accuracy, outperforming baselines that optimize only hardware configurations.

Beyond reducing the carbon footprint, our hardware-model co-optimization framework, CATransformers, enables early-stage design space exploration for next-generation ML accelerators by revealing fundamental trade-offs between hardware area, energy consumption, and carbon efficiency. Our findings show that: Carbon-optimized models-hardware combinations reduce embodied emissions but trade off latency (up to  $2.4 \times$  longer) compared to latency-optimized configurations. Energy and joint latency-carbon optimization mitigates these trade-offs, leading to a 19-20% reduction in carbon footprint without significantly increasing latency. Optimizing for latency alone leads to  $2 \times$  larger hardware architectures, increasing embodied carbon, while energy-optimized designs strike a balance between compute efficiency and sustainability.

Our key contributions include:

- 1. **Insights and Analysis:** Empirical insights revealing how carbon-aware optimization alters model-hardware trade-offs compared to traditional metrics.
- 2. Quantification Framework: A quantitative tool-chain to estimate latency, energy, and total carbon



**Figure 2** Overview of the CATransformers framework. The Bayesian Optimizer iteratively explores the search space by obtaining accuracy, carbon impact, and latency estimates for specific hardware and model architecture combinations from the evaluation modules, and outputs optimized Model and Hardware configuration combinations.

emissions during early design space exploration for custom hardware accelerators.

- 3. Carbon-Aware Co-optimization: A carbon-aware co-optimization framework, CarbonNaaS, using multi-objective Bayesian optimization to target Pareto-optimal trade-offs in accuracy, latency, energy and carbon.
- 4. Sustainable Multi-Modal Models: Using our simulation framework, CarbonNaaS, we demonstrate 17% reduction potential in total carbon emissions with CarbonCLIP while maintaining accuracy and latency compared to edge CLIP model baselines.

By integrating carbon footprint metrics directly into ML system design, this work provides a scalable pathway toward environmentally responsible AI, balancing performance with sustainability.

## 2 Background and Related Works

Hardware Accelerator Search: Specialized hardware accelerators have been developed to efficiently execute deep learning workloads, featuring tensor cores for matrix operations Jouppi et al. (2023); Mahajan et al. (2016); Park et al. (2017); Chen et al. (2016) and vector cores for point-wise operations and activation functions Jouppi et al. (2017); Ghodrati et al. (2024). Previous works have developed hardware accelerator search frameworks for deep learning models Adnan et al. (2024); Wang et al. (2024); Zhang et al. (2022); Sakhuja et al. (2023). These frameworks optimize accelerator configurations based on throughput and energy metrics but do not address the co-optimization of hardware and neural architectures or consider carbon footprint.

Hardware and Neural Architecture Co-optimization: While previous works on co-optimizing hardware and neural architecture search (NAS) focus on operational metrics such as latency Zhou et al. (2022); Choi et al. (2021); Jiang et al. (2020); Lin et al. (2021), none addresses carbon footprint as a primary optimization goal, especially for multi-modal models. The ones that do use carbon footprint as a primary optimization objective, do not optimize jointly with hardware and neural architecture Elgamal et al. (2023, 2025); Zhao and Guo (2023); Gupta et al. (2023). Our research directly addresses this gap.

CLIP Models and Edge Variants: Recent advancements in multi-modal AI have led to the development of models that understand and correlate text and images, such as CLIP Radford et al. (2021) models. These models feature a Transformer-based text encoder and an image encoder based on ResNet He et al. (2016) or Vision Transformer (ViT) Dosovitskiy et al. (2021), and are trained on massive datasets (400 million to 2.5 billion images) to learn visual concepts and their textual associations Schuhmann et al. (2021); Xu et al. (2024); Schuhmann et al. (2022); Gadre et al. (2023). This enables applications in zero-shot classification, image generation, and multi-modal retrieval.

Several works have adapted CLIP models for edge computing by improving training strategies, pruning, and introducing new architectures Wu et al. (2023); Lin et al. (2024); Shi et al. (2023); Vasu et al. (2024). However, these efforts focus on optimizing accuracy and latency, overlooking the carbon footprint. Our framework optimizes ViT-based CLIP models to reduce both the carbon footprint and maintain comparable accuracy and latency.

#### 2.1 The Carbon Footprint of AI Systems

While operational carbon has been extensively studied, embodied carbon remains under-explored Acun et al. (2023); Zhao and Guo (2023); Lacoste et al. (2019); Wu et al. (2022); Zhao et al. (2024); Li et al. (2024). Recent developments of frameworks like ACT Gupta et al. (2022), IMEC.netzero IMEC (2025), and LLMCarbon Faiz et al. (2024) enable the modeling of embodied carbon in ML systems, but a comprehensive solution for accurately quantifying total carbon footprint—both operational and embodied—is still lacking. Moreover, joint co-optimization techniques to minimize ML systems' carbon footprint during custom hardware design remain unaddressed, representing a critical gap in sustainable AI.

#### 3 Framework Overview

In this section, we introduce CATransformers, a carbon-aware architecture search framework for sustainability-driven co-optimization of ML models and hardware architectures. As shown in Figure 2, CATransformers takes three inputs: (1) a base ML model, (2) a hardware architecture template, and (3) optimization metrics and constraints, which define the hardware and software search space. The framework consists of three core components: a multi-objective optimizer and two evaluation modules—an ML model evaluator and a hardware estimator. Below, we describe each component in detail.

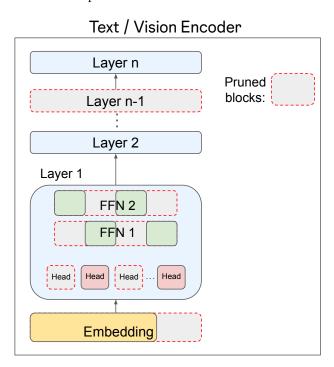


Figure 3 Overview of the dimensions pruned for the encoder. Each layer within a transformer is pruned to the same dimensions, and the Text and Vision encoders separately.

#### 3.1 CATransformers Inputs

Base Model: The base model is a large, pre-trained model that serves as the foundation for pruned models generated by the framework. It determines the overall architecture, shape, and functionality of the optimized models. This work focuses on CLIP-based architectures, but the framework can extend to other models. Pruning is performed along multiple dimensions, including the number of layers, feedforward network size, attention heads, and embedding dimension (as illustrated in Figure 3).

Hardware Architecture Template: The template (Figure 4) defines the accelerator's components and search parameters, based on prior academic and industry designs Jouppi et al. (2017); Adnan et al. (2024); Zhang et al. (2022); Wang et al. (2024). It consists of tensor cores with Processing Elements (PEs) arranged in X and Y dimensions to accelerate GEMM operations. Vector processing units handle element-wise computations,

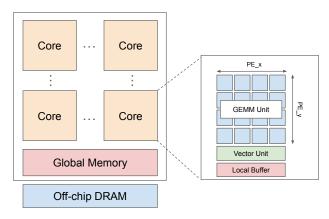


Figure 4 Architecture template based on prior architectures and exploration techniques Jouppi et al. (2017); Wang et al. (2024) for dense deep learning acceleration.

Table 1 Architecture Design Space Parameters.

Parameter Description	Notation	Potential Values
Design Space ( $S$ )		
Number of Cores	TC	1 to 4 powers of 2
PE Array X dim	$PE_{\mathbf{X}}$	1 to 256 powers of 2
PE Array Y dim	$PE_{\mathbf{Y}}$	1 to 256 powers of 2
Global Buffer Size	GLB	1 to 8 MB powers of 2
Local Buffer Size	L2	256 KB to 4 MB powers of 2
Local Bandwidth	$L2_{ m bw}$	1  to  256  words/cycle
Vector Unit width	$V_{ m pe}$	$= pe_{\mathrm{X}}$
Fixed Parameters		
Global Bandwidth	$GLB_{\mathrm{bw}}$	256  words/cycle
Off-chip DRAM Size	HBM	1 GB
Technology	Tech	22 nm
BitWidth	B	8
Maximum TOPS	$T_{max}$	20 TOPS
Frequency	f	500 MHz

while each core has a local buffer for data reuse. All cores share a global SRAM-based buffer and an off-chip memory for storing model parameters and activations. Inline with existing edge inference accelerators Sumbul et al. (2022); Wu et al. (2024b), the global SRAM is designed as scratchpad memory for efficient multi-core serving. We summarize the hardware design space parameters in Table 1. Increasing cores and compute units enhances performance but raises area and energy costs. Expanding on-chip memory reduces off-chip accesses but adds area and energy overhead. The optimal hardware configuration depends on model architecture, size, and performance constraints.

#### 3.2 ML Model Evaluator

The ML model evaluator estimates the accuracy of a given model architecture using OpenCLIP Ilharco et al. (2021).

**Pruning**: The evaluator first prunes the pre-trained model along multiple dimensions, handling the text and vision encoders separately. We adopt pruning strategies from prior work Lin et al. (2024) to reduce the Feed Forward Network dimension, Number of Attention Heads, and Number of Layers based on block

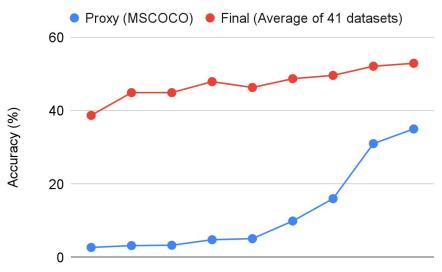


Figure 5 Correlation between the proxy accuracy, fine-tuned and evaluated on MS COCO and the final model accuracy evaluated across 41 datasets.

importance. Additionally, the Embedding Dimension is pruned using techniques from HAT Wang et al. (2020). Each transformer layer is pruned uniformly. Appendix A provides further details on pruning strategies and observations.

Fine-tuning for Accuracy Proxy: Directly evaluating pruned models without training leads to excessively low accuracies. However, fully training each model in the search space is computationally prohibitive, as CLIP training can take hours or even days Xu et al. (2024). To efficiently approximate accuracy while preserving ranking consistency, we fine-tune each pruned model on the MS COCO dataset Lin et al. (2015) (400K samples). As shown in Figure 5, fine-tuning maintains a high Spearman's rank correlation coefficient (0.98) with the final post-pruning training accuracy, ensuring reliable accuracy ranking. We report the mean top-1 recall accuracy on MS COCO as the accuracy proxy for each model.

#### 3.3 Hardware Estimator

The Hardware Estimator module calculates the total carbon footprint (the sum of embodied and operational) and inference latency for a given model and hardware design. Our framework integrates existing libraries to estimate operator latency, energy, and hardware area Wu et al. (2019); Olyaiy et al. (2023). To quantify carbon cost, we use ACT Gupta et al. (2022) for embodied carbon estimation and Electricity Maps Maps (2025) for operational carbon per inference. The operational carbon is then scaled over the hardware's lifetime. A detailed overview of the tool-chain integration is provided in Appendix B.

#### 3.4 Multi-Objective Optimization

The Multi-Objective Optimizer iteratively explores both model and hardware search spaces, taking the ML model and hardware template as inputs. It utilizes the Ax platform Meta Platforms (2024) and BoTorch Balandat et al. (2020) to perform multi-objective Bayesian optimization via the qNEHVI algorithm Daulton et al. (2021). The optimization targets maximizing accuracy while minimizing latency, energy, and total carbon (embodied + operational). It explores hardware parameters (Table 1) and model dimensions (Figure 3) within a compute constraint in terms of Tera Operations per Second (TOPS) based on publicly available edge accelerators. Our framework supports four optimization modes, all subject to a number of compute constraint: (1) Accuracy & Total Carbon (given a latency constraint), (2) Accuracy & Latency, (3) Accuracy & Energy, and (4) Accuracy, Total Carbon, and Latency.

#### 3.5 Outputs

CATransformers outputs a combination of model and hardware configurations, which when used together will provide efficiency improvements. After Multi-Objective Optimizer module has identified carbon efficient model architectures that are pruned versions of our base model, we fine-tune our models to reach their final accuracy. Based on preliminary studies, our pruned models require significantly fewer training steps compared to what is used to pre-train the model in prior works Radford et al. (2021); Xu et al. (2024); Wu et al. (2023) to recover the accuracy loss from pruning. We train our CarbonCLIP models with the MetaCLIP Xu et al. (2024) dataset with 2 epochs, using only 40% of the training steps needed for pre-training in previous studies Radford et al. (2021); Xu et al. (2024); Wu et al. (2023).

## 4 Experimental Settings

Model: We use CATransformers to optimize the CLIP-ViT-B/16 architecture, pre-trained on DataComp-1B Gadre et al. (2023), for various metrics and use cases. We generate CarbonCLIP model and hardware configurations and compare them with (1) CLIP ViT-B/16 and (2) TinyCLIP, a state-of-the-art small CLIP model. Additionally, for baseline comparisons, we perform a hardware architecture search with fixed model parameters for two sets of models: (1) CLIP architectures (ViT-B/16, ViT-L/14, ViT-H/14) and (2) the TinyCLIP family of models.

Hardware: Accelerator area, operator latency, and energy are estimated using open-source tool-chains (Section 3.3), assuming a 22nm process technology for 8-bit integer operations. We validate our latency estimates using SCALE-Sim Samajdar et al. (2018, 2020) for QKV projection operators in CLIP ViT-B/16 and CLIP ViT-B/32. Analytical estimates, widely used for design space exploration Parashar et al. (2019); Wu et al. (2019); Ghodrati et al. (2024); Olyaiy et al. (2023), achieve up to 95% accuracy Parashar et al. (2019); Wu et al. (2019) and were on average within 13% of SCALE-Sim's cycle-accurate latency results. Unlike Sunstone, which maps tensor computations for general dataflows involving spatial arrays, SCALE-Sim specifically targets systolic-array-like architectures. Therefore, for fair comparisons, we validate against similar architectures previously evaluated in the tool (square arrays of dimensions 16x16, 32x32, and 64x64) that can be simulated using SCALE-Sim Samajdar et al. (2018).

For operational carbon footprint estimates, we use the U.S. California grid's carbon intensity and the energy use of the accelerator per inference. We then scale the operational carbon footprint to a 3-year hardware lifespan, assuming one inference per second, consistent with mobile device usage patterns Apple Inc. (2021); NSYS Group (2024); Harmony Healthcare IT (2025). For embodied carbon, we assume the Taiwan grid.

**Execution Setup:** Optimization runs on a single node with 8 V100-SXM2-16GB GPUs and 80 CPU cores, performing 100 Bayesian optimization trials. Post-pruning training uses 224 GPUs with a batch size of 128 per GPU and a learning rate of  $5 \times 10^{-4}$ . CarbonCLIP models are trained for 2 epochs on the MetaCLIP-2.5B dataset Xu et al. (2024) with distillation from the base model.

#### 5 Evaluation Results

#### 5.1 HW Architecture Search with Fixed Model

To understand how we can optimize the total carbon footprint of machine learning models, we first analyze the carbon footprint of state-of-the-art CLIP architectures. To establish a fair baseline, we leverage components of CATransformers to perform a hardware architecture search under a fixed model architecture. We optimize for both carbon footprint and latency of inference systems within a constraint of 20 TOPS compute performance. The framework generates a Pareto frontier of results, illustrating the optimal trade-offs between carbon and latency optimizations. We present the hardware architecture configurations with the minimum carbon footprint and minimum latency from the optimization for each baseline model in Table 2.

From this analysis, we make two key observations:

**Takeaway 1:** Optimizing for carbon footprint results in smaller hardware architectures with fewer compute and memory resources, which leads to reduced area and carbon footprint but often results in higher latency.

Table 2 Preliminary Hardware Architecture Search for fixed baselines CLIP model architectures.

	Total			Minimum	Carbon		Minimum Latency					
Model Architecture	Params	Carbon	Latency	F	Iardware Arc	hitecture	Carbon	Latency	Hardware Architecture			
	(M)	(kgCO2e)	(ms)	# Cores	Core	Memory Config	(kgCO2e)	(ms)	# Cores	Core	Memory Config	
				# Cores	Dimension	{Local, Global}			# Cores	Dimension	{Local, Global}	
CLIP-B/16	149	0.54	18.5	1 (256,8)		64 KB, 2MB	0.69	5.4	2	(256,16)	256KB, 4MB	
CLIP-L/14	427	1.43	68.7	2	(128,16)	128 KB, 2MB	1.76	66.4	4	(64,64)	256KB, 4MB	
CLIP-H/14	986	1.92	71.0	1	(128,32)	128KB, 4MB	2.60	70.2	4	(256,4)	512KB, 4MB	
TinyCLIP-8M/16	41	0.34	3.0	2	(256,4)	64KB, 2MB	0.56	1.3	1	(256,64)	256KB, 8MB	
TinyCLIP-39M/16	83	0.44	9.4	1	(256,8)	64KB, 2MB	0.59	2.2	4	(256,16)	128KB, 4MB	
TinyCLIP-40M/32	84	0.37	8.6	1 (32,32)		64KB, 2MB	0.46	1.1	4	(128,32)	64KB, 2MB	
TinyCLIP-61M/32	115	0.39	9.7	1 (128,8)		64KB, 2MB	0.49	1.4	4	(128,32)	64KB, 2MB	

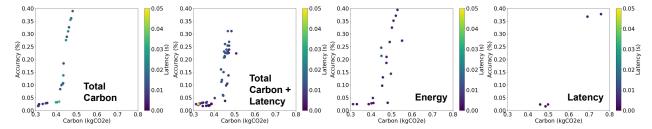


Figure 6 Pareto frontiers for different optimization modes under a 20 TOPS compute constraint.

For example, for the CLIP ViT-B/16 architecture, optimizing for carbon yields a hardware architecture consisting of a single core accelerator with 4K PEs, 64KB of local memory, and 2MB of global memory. In contrast, when optimizing for latency, the selected hardware architecture features two tensor cores, each with 8K PEs, 256KB of local memory, and a 4MB global memory. The carbon-optimized configuration has a 22% lower total carbon footprint but over 3 times longer latency per inference. Note that since the total energy consumption is proportional to latency (i.e. energy = power  $\times$  delay), the total latency can not be made arbitrarily long when using a smaller area hardware configuration. This is because excessively long latencies would not only fail to meet realistic performance targets but also start to significantly impact operational carbon costs again. These trade-offs highlight the need for a careful balance between carbon footprint and latency requirements when designing AI systems.

**Takeaway 2:** It is important to tailor the hardware to the specific characteristics of the model architecture.

For example, TinyCLIP-61M/32, despite having more parameters, can achieve comparable or even lower latency than TinyCLIP-39M/16 when paired with tailored hardware configuration, also while achieving lower carbon footprint. This is due to differences in parameters, embedding dimensions, and patch sizes, which significantly impact model execution patterns and resulting hardware configurations. These findings emphasize the interdependence of model architecture and hardware design, highlighting the need for co-optimization to enhance resource utilization, performance, and environmental sustainability.

#### 5.2 Joint Model and Hardware Architecture Search Using Different Metrics

In this section, we use CATransformers to employ a joint model and the hardware architecture search with carbon footprint as a central design metric, alongside traditional metrics like accuracy and latency. We compare across each optimization modes under a 20 TOPS compute constraint, reflecting a setup comparable to publicly available edge accelerators HAILO (2025); Nvidia (2025).

**Takeaway 3:** Optimizing for total carbon footprint yields model and hardware architectures with the lowest overall carbon impact, but at the cost of increased latency. In contrast, optimizing for energy consumption strikes a balance between latency and carbon footprint.

Figure 6 presents the Pareto frontiers for latency-only, carbon-only, energy-only, and latency+carbon. Each data point in the figure represents a model and hardware architecture configuration, with accuracy, carbon footprint, and latency represented on the y-axis, x-axis, and color map, respectively. Each experiment is repeated three times for consistency, and accuracy is estimated using the MS COCO dataset. When latency is not an optimization target, a maximum latency constraint of 50ms is enforced Žádník et al. (2022) to ensure

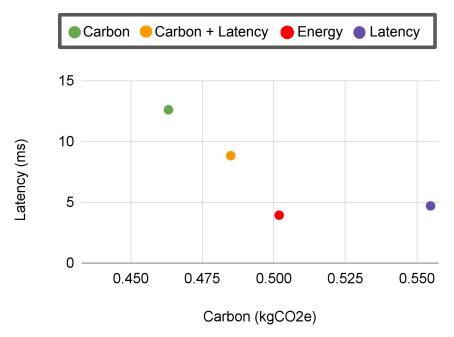


Figure 7 ISO-Accuracy plot illustrating the trade-off between latency and carbon footprint for different optimization modes, achieved at 31% accuracy.

realistic specifications. To provide further comparison of the selected configurations by each optimization mode, Figure 7 illustrates the latency-carbon trade-off at an iso-accuracy point of 31% (in terms of proxy accuracy before fine-tuning). Appendix C provides details into the model and hardware configurations at more accuracy levels.

Comparing optimization modes, we find that carbon optimization significantly reduces carbon footprint by 24% but latency is  $2.4\times$  higher compared to latency-optimized models. Joint carbon and latency optimization mitigates this trade-off, achieving a 20% carbon reduction with only  $1.8\times$  higher latency. Energy optimization strikes a balance, yielding a 19% carbon reduction without increasing latency, effectively balancing sustainability and performance.

**Takeaway 4:** Although energy is not directly linked to latency, energy optimization often leads to lower latency.

Interestingly, this occurs because energy minimization targets both power efficiency and computational duration, indirectly optimizing latency. In contrast, carbon optimization prioritizes minimizing hardware area to account for embodied carbon, leading to a weaker emphasis on latency reduction.

**Takeaway 5:** Optimizing for latency-only favors more compute resources and larger local and global memory compared to carbon-only optimizations.

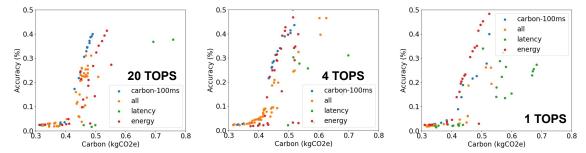
This results in hardware architectures of  $2 \times$  larger area, favoring more compute resources and larger local and global memory compared to carbon-only optimizations. In addition, latency-only optimizations often select models with larger number of parameters, which is compensated by a larger architecture to maintain a lower latency. In contrast, energy-only optimizations on average select larger architectures but instead choose smaller size models to reduce energy consumption by decreasing delay.

**Takeaway 6:** To maximize accuracy while minimizing total carbon footprint, The Bayesian Optimizer prioritizes pruning CLIP models in the following order: FFN dimension and number of attention heads, then the number of layers, and finally the embedding dimension. Text encoders are pruned more aggressively than vision encoders.

We observe that, in general, text encoders are pruned more aggressively than vision encoders. This is because vision models have a greater impact on accuracy, as highlighted in Appendix A. We observed that pruning any dimension across the vision encoder resulted in more significant accuracy loss compared to its text encoder counterpart. This is because vision models process higher-dimensional inputs, more aggressive pruning

**Table 3** The hardware and model architecture properties of each variant of the CarbonCLIP family. Hardware configurations are specified as:  $\{TC, PE_x, PE_y, L2, L2_{bw}, GLB\}$ . Text and Vision encoders are specified as:  $\{Num Layers, FFN Dim, Hidden Dim, Num Heads\}$ .

Name	Carbon	Latency	Hardware	M	Avg. Accuracy		
Name	(kgCO2e)	(ms)	Configuration	Text Encoder	Vision Encoder	Params (M)	over 41 datasets
				Configuration	Configuration	rarams (M)	
CarbonCLIP-XS	0.32	7.1	{1, 256, 4, 64, 32, 2}	{6, 1024, 284, 4}	{6, 1536, 576, 6}	41	38.7
CarbonCLIP-S	0.36	12.0	{1, 256, 4, 64, 64, 2}	{6, 1024, 512, 6}	{8, 1920, 672, 6}	63	45.0
CarbonCLIP-M	0.39	19.7	{1, 256, 4, 64, 128,2 }	{8, 1536, 512, 6}	{9, 2304, 672, 6}	79	47.9
CarbonCLIP-L	0.42	13.7	{1, 256, 8, 64, 128, 2}	{6, 1280, 384, 5}	{10, 2688, 768, 7}	83	48.7
CarbonCLIP-XL	0.49	19.1	{1, 256, 8, 64, 128, 2}	{12, 2048, 512, 4}	{12, 3072, 768, 5}	123	52.0
TinyCLIP-8M/16	0.34	3.0	{2, 256, 4, 64, 32, 2}	{3, 1024, 256, 4}	{10, 1024, 256, 4}	41	30.7
TinyCLIP-39M/16	0.44	9.4	{1, 256, 8, 64, 128, 2}	{6, 2048, 512, 8}	{12, 2048, 512, 8}	83	45.0
CLIP-B/16 - DataComp	0.54	18.5	{1, 256, 8, 64, 128, 2}	{12, 2048, 512, 8}	{12, 3072, 768, 12}	149	53.2



**Figure 8** Pareto frontiers for different optimization modes with different compute design constraints. Tighter computational constraints yield similar outcomes when optimizing for total carbon emissions and optimizing for energy consumption (operational carbon).

disrupts its ability produce embeddings that align well with the text encoder, leading to degraded retrieval and classification accuracy.

#### 5.3 Accuracy Evaluation

We compare CarbonCLIP models with TinyCLIP and the CLIP-ViT-B/16 baseline (pretrained on DataComp-1B). The evaluation considers the most carbon-efficient configurations of each baseline model. Table 3 presents the CarbonCLIP family of models. We select models at various sizes from the Pareto frontiers shown in Figure 6, perform post-pruning training, and evaluate their performance on all 41 zero-shot evaluation benchmarks in the CLIP benchmark LAION-AI (2022).

The table details the model and hardware configurations of CarbonCLIP-XL to CarbonCLIP-XS (largest to smallest models), and the corresponding carbon footprint, latency, and average accuracy. Detailed breakdown of the accuracy for each dataset is in Appendix E. We further extend the CarbonCLIP family to CLIP-B/32 architecture to compare to TinyCLIP's B/32 baselines in Appendix D.

CarbonCLIP family of models achieve better average performance at various carbon footprint levels. Notably, CarbonCLIP-XL achieves baseline-level accuracy with an 10% reduction in carbon footprint. CarbonCLIP-XS achieves an 8% increase in accuracy with a 3% reduction in carbon footprint compared to TinyCLIP-8M/16. CarbonCLIP-L, CarbonCLIP-M, and CarbonCLIP-S all achieve significant reductions in carbon footprint compared to TinyCLIP-39M/16, with CarbonCLIP-L achieving a 4% increase in accuracy and a 4.5% reduction in carbon footprint, CarbonCLIP-M achieving an 11% reduction in carbon footprint with a 3% decrease in accuracy, and CarbonCLIP-S achieving a 17% reduction in carbon footprint without any regression in accuracy.

In terms of hardware configuration, CarbonCLIP models select accelerators with cores of  $PE_x$  dimension 256 to align with the underlying operator dimensions of the CLIP ViT-B/16 architecture, with sequence length of 197 for the vision encoder. Smaller models select a total of 1024 PE units per core, whereas larger models select twice as many PEs to keep the latency of the task low. Due to the reduced size of the CarbonCLIP models, a 64KB local memory and 2MB global memory are sufficient to keep the core utilized.

#### 5.4 Evaluating for Different Compute Constraints

We evaluate various accelerators compute capability scenarios with CATransformers, inspired by commercially available edge devices. Figure 8 illustrates Pareto frontiers for three groups designed with maximum peak performance constraint: 20 TOPS HAILO (2025); Nvidia (2025), 4 TOPS Google (2020), and 1 TOPS Wu et al. (2024b); Intel (2025).

**Takeaway 7:** As hardware architectures shrink due to tighter computational constraints, optimizing for energy consumption (operational carbon) yields outcomes that minimizes both total carbon and inference latency.

As the architecture size decreases due to tighter computational constraints, the outcomes of optimizing for carbon emissions become increasingly similar to those of optimizing for energy consumption (operational carbon). This is because the embodied carbon becomes a less significant portion of the total carbon footprint as the architecture shrinks, causing energy or operational carbon to dominate the overall environmental impact of the ML model. We present the trade-off between embodied and operational carbon for CarbonCLIP models in Appendix F. Note that different grid intensities and expected hardware lifetime can also affect the ratio of embodied and operational carbon, affecting the optimization results.

For smaller architecture areas, despite the inherent reduction in embodied carbon and power consumption, there is a noticeable trend of increased total carbon emissions. This is primarily because smaller architectures have less compute and memory resources at their disposal to accelerate the model efficiently, leading to a rise in delay, and thus in operational energy use. This trend underscores the critical importance of hardware and model architecture co-design to minimize the total carbon impact of ML systems.

We further extended the study to evaluate the impact of latency constraints on model and hardware configurations, as well as their carbon footprints in Appendix G.

#### 6 Conclusion

In this work, we introduced CATransformers, a solution for co-optimizing domain-specific accelerators and model architectures to minimize carbon footprint. By integrating both operational and embodied carbon metrics, our framework enables environmentally conscious design for edge computing environments. We demonstrated its effectiveness with CLIP-based models, achieving significant carbon reductions without sacrificing performance. This work fills a critical gap in sustainable AI deployment and provides a foundation for future advancements in carbon-aware ML design Wu et al. (2024a). Future directions to improve the work include: expanding the framework to support a wider range of model architectures and training scenarios in a data-center setting.

## Impact Statement

This paper introduces a framework that integrates carbon-aware neural architecture search and hardware co-optimization to advance sustainable machine learning. It aims to significantly reduce the carbon footprint of deploying and running ML models, particularly in resource-constrained edge environments. By optimizing both operational and embodied carbon in a holistic way, this work paves the way for environmentally responsible AI systems, focusing on multi-modal models like CLIP. The framework promotes efficient deployment strategies while addressing compute, memory, and energy constraints, fostering collaboration across the machine learning, hardware, and sustainable computing communities.

## Acknowledgements

We would like to thank Bernie Beckerman, David Eriksson, and Max Balandat for their expertise and assistance in building the optimization platform with Ax, Igor Fedorov for sharing insights on pruning models, Daniel Li and Hu Xu for providing valuable guidance on training CLIP models and working with the MetaCLIP dataset. We also would like thank Kim Hazelwood and Kristen Lauter for supporting this work.

#### References

- Bilge Acun, Benjamin Lee, Fiodar Kazhamiaka, Aditya Sundarrajan, Kiwan Maeng, Manoj Chakkaravarthy, David Brooks, and Carole-Jean Wu. Carbon dependencies in datacenter design and management. *SIGENERGY Energy Inform. Rev.*, 3(3):21–26, October 2023. doi: 10.1145/3630614.3630619. https://doi.org/10.1145/3630614.3630619.
- Muhammad Adnan, Amar Phanishayee, Janardhan Kulkarni, Prashant J. Nair, and Divya Mahajan. Workload-aware hardware accelerator mining for distributed deep learning training, 2024. https://arxiv.org/abs/2404.14632.
- Apple Inc. Apple environmental progress report 2021. https://www.apple.com/environment/pdf/Apple\_Environmental\_Progress\_Report\_2021.pdf, 2021.
- Maximilian Balandat, Brian Karrer, Daniel Jiang, Samuel Daulton, Ben Letham, Andrew G Wilson, and Eytan Bakshy. Botorch: A framework for efficient monte-carlo bayesian optimization. *Advances in neural information processing systems*, 33:21524–21538, 2020.
- Yu-Hsin Chen, Tushar Krishna, Joel S Emer, and Vivienne Sze. Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE journal of solid-state circuits*, 52(1):127–138, 2016.
- Kanghyun Choi, Deokki Hong, Hojae Yoon, Joonsang Yu, Youngsok Kim, and Jinho Lee. Dance: Differentiable accelerator/network co-exploration. In 2021 58th ACM/IEEE Design Automation Conference (DAC), pages 337–342, 2021. doi: 10.1109/DAC18074.2021.9586121.
- Benoit Courty, Victor Schmidt, Sasha Luccioni, Goyal-Kamal, MarionCoutarel, Boris Feld, Jérémy Lecourt, Liam-Connell, Amine Saboni, Inimaz, supatomic, Mathilde Léval, Luis Blanche, Alexis Cruveiller, ouminasara, Franklin Zhao, Aditya Joshi, Alexis Bogroff, Hugues de Lavoreille, Niko Laskaris, Edoardo Abati, Douglas Blank, Ziyao Wang, Armin Catovic, Marc Alencon, Michał Stęchły, Christian Bauer, Lucas Otávio N. de Araújo, JPW, and MinervaBooks. mlco2/codecarbon: v2.4.1, May 2024. https://doi.org/10.5281/zenodo.11171501.
- Samuel Daulton, Maximilian Balandat, and Eytan Bakshy. Parallel bayesian optimization of multiple noisy objectives with expected hypervolume improvement. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, NIPS '21, Red Hook, NY, USA, 2021. Curran Associates Inc. ISBN 9781713845393.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. https://openreview.net/forum?id=YicbFdNTTy.
- Mariam Elgamal, Doug Carmean, Elnaz Ansari, Okay Zed, Ramesh Peri, Srilatha Manne, Udit Gupta, Gu-Yeon Wei, David Brooks, Gage Hills, et al. Carbon-efficient design optimization for computing systems. In *Proceedings of the 2nd Workshop on Sustainable Computer Systems*, pages 1–7, 2023.
- Mariam Elgamal, Doug Carmean, Elnaz Ansari, Okay Zed, Ramesh Peri, Srilatha Manne, Udit Gupta, Gu-Yeon Wei, David Brooks, Gage Hills, and Carole-Jean Wu. CORDOBA: Carbon-Efficient Optimization Framework for Computing Systems . In 2025 IEEE International Symposium on High Performance Computer Architecture (HPCA), 2025.
- Ahmad Faiz, Sotaro Kaneda, Ruhan Wang, Rita Osi, Prateek Sharma, Fan Chen, and Lei Jiang. Llmcarbon: Modeling the end-to-end carbon footprint of large language models, 2024. https://arxiv.org/abs/2309.14393.
- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, Eyal Orgad, Rahim Entezari, Giannis Daras, Sarah Pratt, Vivek Ramanujan, Yonatan Bitton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Krishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. Datacomp: In search of the next generation of multimodal datasets, 2023. https://arxiv.org/abs/2304.14108.
- Soroush Ghodrati, Sean Kinzer, Hanyang Xu, Rohan Mahapatra, Byung Hoon Ahn, Dong Kai Wang, Lavanya Karthikeyan, Amir Yazdanbakhsh, Jongse Park, Nam Sung Kim, and Hadi Esmaeilzadeh. Tandem processor: Grappling with emerging operators in neural networks. In ASPLOS, 2024.
- Google. Helping you bring local ai to applications from prototype to production, 2020. https://coral.ai/products/.
- Udit Gupta, Mariam Elgamal, Gage Hills, Gu-Yeon Wei, Hsien-Hsin S. Lee, David Brooks, and Carole-Jean Wu. Act: designing sustainable computer systems with an architectural carbon modeling tool. In *Proceedings of*

- the 49th Annual International Symposium on Computer Architecture, ISCA '22, page 784–799, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450386104. doi: 10.1145/3470496.3527408. https://doi.org/10.1145/3470496.3527408.
- Udit Gupta, Daniel R Jiang, Maximilian Balandat, and Carole-Jean Wu. Towards green, accurate, and efficient ai models through multi-objective optimization. In *ICLR 2023 Workshop on Tackling Climate Change with Machine Learning*, 2023. https://www.climatechange.ai/papers/iclr2023/55.
- HAILO. Hailo-8 ai accelerator, 2025. https://hailo.ai/products/ai-accelerators/hailo-8-ai-accelerator/.
- Harmony Healthcare IT. Are you addicted to your phone? american phone usage & screen time statistics. https://www.harmonyhit.com/phone-screen-time-statistics/, 2025.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian 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.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. https://doi.org/10.5281/zenodo.5143773. If you use this software, please cite it as below.
- IMEC. imec.netzero. "https://netzero.imec-int.com/", 2025.
- Intel. Intel movidius myriad x vision processing unit, 2025. https://www.intel.com/content/www/us/en/products/sku/204770/intel-movidius-myriad-x-vision-processing-unit-0gb/specifications.html.
- Akriti Jain, Saransh Sharma, Koyel Mukherjee, and Soumyabrata Pal. First: Finetuning router-selective transformers for input-adaptive latency reduction, 2024. https://arxiv.org/abs/2410.12513.
- Weiwen Jiang, Lei Yang, Edwin Hsing-Mean Sha, Qingfeng Zhuge, Shouzhen Gu, Sakyasingha Dasgupta, Yiyu Shi, and Jingtong Hu. Hardware/software co-exploration of neural architectures. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(12):4805–4815, 2020. doi: 10.1109/TCAD.2020.2986127.
- 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. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*, ISCA '17, page 1–12, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348928. doi: 10.1145/3079856.3080246.
- Norman P. Jouppi, George Kurian, Sheng Li, Peter Ma, Rahul Nagarajan, Lifeng Nai, Nishant Patil, Suvinay Subramanian, Andy Swing, Brian Towles, Cliff Young, Xiang Zhou, Zongwei Zhou, and David Patterson. Tpu v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings, 2023.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020. https://arxiv.org/abs/2001.08361.
- Hyoukjun Kwon, Prasanth Chatarasi, Michael Pellauer, Angshuman Parashar, Vivek Sarkar, and Tushar Krishna. Understanding Reuse, Performance, and Hardware Cost of DNN Dataflow: A Data-Centric Approach. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO '52, page 754–768, New York, NY, USA, 2019. Association for Computing Machinery.
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. arXiv preprint arXiv:1910.09700, 2019.
- LAION-AI. Clip benchmark, 2022. https://github.com/LAION-AI/CLIP\_benchmark.
- Baolin Li, Yankai Jiang, Vijay Gadepally, and Devesh Tiwari. Sprout: Green generative AI with carbon-efficient LLM inference. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21799–21813, Miami, Florida, USA,

- November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1215. https://aclanthology.org/2024.emnlp-main.1215/.
- Haokun Lin, Haoli Bai, Zhili Liu, Lu Hou, Muyi Sun, Linqi Song, Ying Wei, and Zhenan Sun. Mope-clip: Structured pruning for efficient vision-language models with module-wise pruning error metric. 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 27360–27370, 2024. https://api.semanticscholar.org/CorpusID:268363825.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. https://arxiv.org/abs/1405.0312.
- Yujun Lin, Mengtian Yang, and Song Han. Naas: Neural accelerator architecture search. In 2021 58th ACM/IEEE Design Automation Conference (DAC), pages 1051–1056, 2021. doi: 10.1109/DAC18074.2021.9586250.
- Divya Mahajan, Jongse Park, Emmanuel Amaro, Hardik Sharma, Amir Yazdanbakhsh, Joon Kyung Kim, and Hadi Esmaeilzadeh. Tabla: A unified template-based framework for accelerating statistical machine learning. In 2016 IEEE International Symposium on High Performance Computer Architecture (HPCA), pages 14–26. IEEE, 2016.
- Electricity Maps. The leading electricity grid api, 2025. https://www.electricitymaps.com/.
- Kyle Marino, Pengmiao Zhang, and Viktor K. Prasanna. ME- ViT: A Single-Load Memory-Efficient FPGA Accelerator for Vision Transformers. In 2023 IEEE 30th International Conference on High Performance Computing, Data, and Analytics (HiPC), pages 213–223, Los Alamitos, CA, USA, December 2023. IEEE Computer Society. doi: 10.1109/HiPC58850.2023.00039. https://doi.ieeecomputersociety.org/10.1109/HiPC58850.2023.00039.
- Inc. Meta Platforms. Adaptive experimentation platform, 2024. https://ax.dev/.
- NSYS Group. Average device lifespan: How long does a cell phone last? https://nsysgroup.com/blog/average-device-lifespan-how-long-does-a-cell-phone-last, 2024.
- Nvidia. Jetson orin nano, 2025. https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-orin/nano-super-developer-kit/.
- MohammadHossein Olyaiy, Christopher Ng, Alexandra Fedorova, and Mieszko Lis. Sunstone: A Scalable and Versatile Scheduler for Mapping Tensor Algebra on Spatial Accelerators. In 2023 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), 2023.
- Angshuman Parashar, Priyanka Raina, Yakun Sophia Shao, Yu-Hsin Chen, Victor A. Ying, Anurag Mukkara, Rangharajan Venkatesan, Brucek Khailany, Stephen W. Keckler, and Joel Emer. Timeloop: A Systematic Approach to DNN Accelerator Evaluation. In 2019 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), pages 304–315, 2019.
- Jongse Park, Hardik Sharma, Divya Mahajan, Joon Kyung Kim, Preston Olds, and Hadi Esmaeilzadeh. Scale-out acceleration for machine learning. In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 367–381, 2017.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. https://arxiv.org/abs/2103.00020.
- Chirag Sakhuja, Zhan Shi, and Calvin Lin. Leveraging domain information for the efficient automated design of deep learning accelerators. In 2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA), pages 287–301, 2023. doi: 10.1109/HPCA56546.2023.10071095.
- Ananda Samajdar, Yuhao Zhu, Paul Whatmough, Matthew Mattina, and Tushar Krishna. Scale-sim: Systolic cnn accelerator simulator. arXiv preprint arXiv:1811.02883, 2018.
- Ananda Samajdar, Jan Moritz Joseph, Yuhao Zhu, Paul Whatmough, Matthew Mattina, and Tushar Krishna. A systematic methodology for characterizing scalability of dnn accelerators using scale-sim. In 2020 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS), pages 58–68. IEEE, 2020.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. ArXiv, abs/2111.02114, 2021. https://api.semanticscholar.org/CorpusID:241033103.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine

- Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models, 2022. https://arxiv.org/abs/2210.08402.
- Dachuan Shi, Chaofan Tao, Ying Jin, Zhendong Yang, Chun Yuan, and Jiaqi Wang. UPop: Unified and progressive pruning for compressing vision-language transformers. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pages 31292–31311. PMLR, 2023.
- H. Ekin Sumbul, Tony F. Wu, Yuecheng Li, Syed Shakib Sarwar, William Koven, Eli Murphy-Trotzky, Xingxing Cai, Elnaz Ansari, Daniel H. Morris, Huichu Liu, Doyun Kim, Edith Beigne, Reality Labs, and Meta. System-level design and integration of a prototype ar/vr hardware featuring a custom low-power dnn accelerator chip in 7nm technology for codec avatars. In 2022 IEEE Custom Integrated Circuits Conference (CICC), pages 01–08, 2022. doi: 10.1109/CICC53496.2022.9772810.
- Pavan Kumar Anasosalu Vasu, Hadi Pouransari, Fartash Faghri, Raviteja Vemulapalli, and Oncel Tuzel. Mobileclip: Fast image-text models through multi-modal reinforced training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15963–15974, 2024.
- Jakub Žádník, Markku Mäkitalo, Jarno Vanne, and Pekka Jääskeläinen. Image and video coding techniques for ultra-low latency. *ACM Comput. Surv.*, 54(11s), September 2022. ISSN 0360-0300. doi: 10.1145/3512342. https://doi.org/10.1145/3512342.
- Hanrui Wang, Zhanghao Wu, Zhijian Liu, Han Cai, Ligeng Zhu, Chuang Gan, and Song Han. Hat: Hardware-aware transformers for efficient natural language processing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.686. http://dx.doi.org/10.18653/v1/2020.acl-main.686.
- Irene Wang, Jakub Tarnawski, Amar Phanishayee, and Divya Mahajan. Integrated hardware architecture and device placement search. In *International Conference on Machine Learning*, 2024.
- Carole-Jean Wu, Ramya Raghavendra, Udit Gupta, Bilge Acun, Newsha Ardalani, Kiwan Maeng, Gloria Chang, Fiona Aga, Jinshi Huang, Charles Bai, Michael Gschwind, Anurag Gupta, Myle Ott, Anastasia Melnikov, Salvatore Candido, David Brooks, Geeta Chauhan, Benjamin Lee, Hsien-Hsin Lee, Bugra Akyildiz, Maximilian Balandat, Joe Spisak, Ravi Jain, Mike Rabbat, and Kim Hazelwood. Sustainable ai: Environmental implications, challenges and opportunities. In D. Marculescu, Y. Chi, and C. Wu, editors, *Proceedings of Machine Learning and Systems*, volume 4, pages 795–813, 2022. https://proceedings.mlsys.org/paper\_files/paper/2022/file/462211f67c7d858f663355eff93b745e-Paper.pdf.
- Carole-Jean Wu, Bilge Acun, Ramya Raghavendra, and Kim Hazelwood. Beyond Efficiency: Scaling AI Sustainably .  $IEEE\ Micro,\ 44(05),\ 2024a.$
- Kan Wu, Houwen Peng, Zhenghong Zhou, Bin Xiao, Mengchen Liu, Lu Yuan, Hong Xuan, Michael Valenzuela, Xi (Stephen) Chen, Xinggang Wang, Hongyang Chao, and Han Hu. Tinyclip: Clip distillation via affinity mimicking and weight inheritance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 21970–21980, October 2023.
- Tony F. Wu, Huichu Liu, H. Ekin Sumbul, Lita Yang, Dipti Baheti, Jeremy Coriell, William Koven, Anu Krishnan, Mohit Mittal, Matheus Trevisan Moreira, Max Waugaman, Laurent Ye, and Edith Beigné. 11.2 a 3d integrated prototype system-on-chip for augmented reality applications using face-to-face wafer bonded 7nm logic at  $< 2\mu$ m pitch with up to 40% energy reduction at iso-area footprint. In 2024 IEEE International Solid-State Circuits Conference (ISSCC), volume 67, pages 210–212, 2024b. doi: 10.1109/ISSCC49657.2024.10454529.
- Yannan Nellie Wu, Joel S. Emer, and Vivienne Sze. Accelergy: An Architecture-Level Energy Estimation Methodology for Accelerator Designs. In 2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD), pages 1–8, 2019.
- Hu Xu, Saining Xie, Xiaoqing Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying CLIP data. In *The Twelfth International Conference on Learning Representations*, 2024. https://openreview.net/forum?id=5BCFlnfE1g.
- Dan Zhang, Safeen Huda, Ebrahim Songhori, Kartik Prabhu, Quoc Le, Anna Goldie, and Azalia Mirhoseini. A full-stack search technique for domain optimized deep learning accelerators. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, page 27–42, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392051.
- Yiyang Zhao and Tian Guo. Carbon-efficient neural architecture search. In *Proceedings of the 2nd Workshop on Sustainable Computer Systems*, HotCarbon '23, page 1–7. ACM, July 2023. doi: 10.1145/3604930.3605708. http://dx.doi.org/10.1145/3604930.3605708.

Yiyang Zhao, Yunzhuo Liu, Bo Jiang, and Tian Guo. CE-NAS: An end-to-end carbon-efficient neural architecture search framework. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. https://openreview.net/forum?id=v6W55lCkhN.

Yanqi Zhou, Xuanyi Dong, Tianjian Meng, Mingxing Tan, Berkin Akin, Daiyi Peng, Amir Yazdanbakhsh, Da Huang, Ravi Narayanaswami, and James Laudon. Towards the co-design of neural networks and accelerators. In D. Marculescu, Y. Chi, and C. Wu, editors, *Proceedings of Machine Learning and Systems*, volume 4, pages 141–152, 2022. https://proceedings.mlsys.org/paper\_files/paper/2022/file/4c430a4d0a7de11e85fa5b076e7f1895-Paper.pdf.

## **Appendix**

## A Ablation Study: Pruning and Finetuning

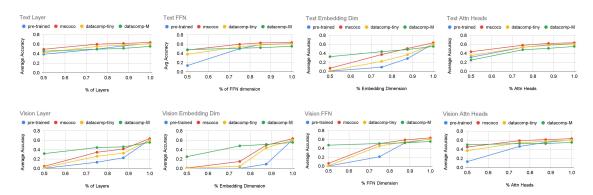


Figure 9 Pruning and fine-tuning results for each dimension of the Text and Vision encoders. In this ablation study, we studied the effect of pruning each dimension and the effect of fine-tuning on various datasets to find a good proxy for approximating the accuracy after training the pruned models on larger datasets. All models were evaluated against MS COCO, and our key observations are as follows:

We make some key observations in terms of the importance of each dimension to the overall accuracy of the model. First, accuracy drops significantly after pruning to 50% of any dimension, even after fine-tuning. Therefore, we confine the search space to a minimum of half of each dimension. Second, The vision model has a more significant impact on accuracy compared to the text model. Finally, pruning the embedding dimension has the most significant impact on accuracy among all pruned parameters followed by number of layers, then FFN dimensions and Number of attention heads.

Fine-tuning on MS COCO: Even with just a single epoch, fine-tuning on MS COCO significantly improves the accuracy of pruned models, making it a good and fast proxy for evaluating their overall potential. Fine-tuning on Datacomp-Tiny: Training on a small subset of a pre-training dataset (Datacomp-Tiny, a 400k subset of Datacomp-Medium unfiltered) also improves accuracy, albeit with lower overall accuracy compared to MS COCO. Fine-tuning on Datacomp-Medium: Using a general pre-training dataset (Datacomp-Medium unfiltered) reduces the variance in accuracy between models of different sizes, showing that smaller models can achieve comparable accuracy when trained with enough data. This insentifies our post-pruning training with a large and high quality MetaCLIP 2.5B Dataset.

In general, larger models will attain higher accuracy compared to models with fewer parameters when trained for the same number of steps. However, models with fewer parameters may still achieve comparable accuracy as larger models given more training steps Kaplan et al. (2020). Therefore, we fine-tune each model with the same computation FLOPS, allowing smaller models to train for more flops and recover their accuracy from more extensive pruning.

## **B** Estimation Tool-chain Integration

In this section we present in detail the integration of the tools used in the hardware estimator.

We use Accelergy Wu et al. (2019) to estimate the area and access energy of each component, and Sunstone Olyaiy et al. (2023) to estimate the per-operator latency and energy. For carbon estimations, we use ACT Gupta et al. (2022) to estimate the embodied carbon of the hardware architecture based on the area of the accelerator. Given the energy estimates from sunstone, Electricity Maps Maps (2025) is used to estimate the operational carbon of executing a single inference. We then scale the operational carbon to the total lifetime of the hardware architecture.

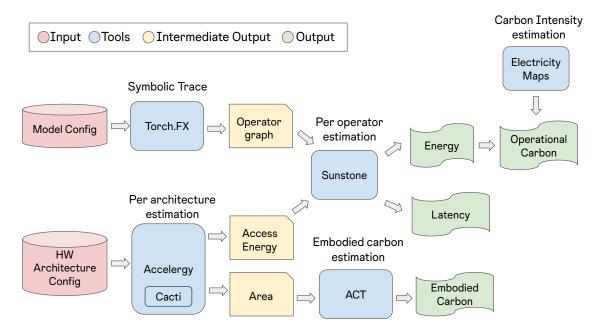


Figure 10 Integration of tools to estimate the total carbon footprint of the given hardware and model configuration.

## C ISO Accuracy

In this sections we present full details on different design points found by each optimization metric for a given accuracy point in Table 4. The key observations are discussed in Section 5.2

## D CLIP ViT B/32 Accuracy Evaluations

In this section, we present results that demonstrate the applicability of CATransformers to the CLIP ViT-B/32 architecture and compare its performance against the TinyCLIP baselines. (Table 5) Our findings show that CATransformers's optimization can be effectively generalized to other model architectures, yielding up to 5% and 8% reductions in carbon footprint while achieving higher accuracy and comparable latency compared to the TinyCLIP baselines, respectively.

#### E CLIP Benchmark Full Result

In this section we provide a detailed breakdown on the accuracy of each dataset for CarbonCLIP and the evaluated baselines in Table 6.

## F Carbon Footprint Breakdown

We provide a breakdown of the carbon footprint for each variant of the CarbonCLIP family model. As shown in Figure 11, as the model increases in size, the proportion of operational carbon in the overall carbon footprint of the model increases from 20% to over 40%. The CarbonCLIP-XL model has  $3 \times 1$  the number of parameters and almost  $3 \times 1$  the latency of CarbonCLIP-XS, but the selected hardware architecture only has double the number of compute PEs. Therefore, the operational carbon increases proportionally more than the increase in embodied carbon. This highlights need of co-optimizing the model and hardware architecture to maintain an intricate balance between operational and embodied carbon, keeping the overall carbon footprint of the system low. As such, the expected lifetime and source of power are also important factors that need to be taken into account during the optimization process.

**Table 4** The hardware and model architecture configuration found by each optimization metric at each accuracy point. Hardware configurations are specified in the format of:  $\{TC, pe_x, pe_y, L2, L2_{bw}, glb\}$ . Text and Vision encoders are specified in the format of  $\{\text{Num Layers, FFN Dim, Hidden Dim, Num Heads}\}$ 

Λ	0-4::4:	Combon	T =4======	II and and an	Model						
Accuracy $(+/-1\%)$	Optimization Metric	Carbon	Latency (ms)	Hardware	Configuration						
(+/- 170) Metric		(kgCO2e)	(ms)	Configuration	Text Encoder	Vision Encoder	Danama (M)				
					Configuration	Configuration	Params (M)				
	Carbon	0.46	12.6	{1, 256, 8, 64, 64, 2}	{9, 1536, 512, 6}	{12, 576, 768, 8}	104				
31%	Energy	0.50	3.9	$\{2, 256, 16, 128, 32, 2\}$	{7, 1536, 384, 8}	{12, 576, 768, 8}	101				
31/0	Latency	0.55	4.7	$\{2, 256, 16, 128, 32, 2\}$	{10, 1792, 384, 7}	{12, 672, 768, 11}	111				
	Carbon + Latency	0.48	8.8	$\{2, 256, 8, 64, 64, 2\}$	{9, 1280, 512, 6}	{11, 672, 768, 9}	105				
	Carbon	0.44	10.9	{1,256, 8, 64, 64 2}	{9, 1536, 512, 6}	{11, 672, 768,6}	95				
19.5%	Energy	0.48	3.5	$\{2, 256, 16, 128, 32, 2\}$	{6, 1536, 512, 8}	{12, 480, 768, 7}	90				
19.570	Latency	0.55	8.2	{4, 256, 4, 128, 64, 2}	{11, 1536, 384, 5}	{11, 576, 768, 12}	99				
	Carbon + Latency	0.45	7.3	$\{2, 256, 4, 64, 128, 2\}$	{9, 1024, 512, 5}	{11, 576, 768, 8}	94				
	Carbon	0.43	22.1	{1, 256, 4, 64, 64, 2}	{7, 1536, 6384 5}	{11, 576, 768, 8}	84				
13%	Energy	0.49	7.3	$\{2, 256, 8, 64, 64, 2\}$	{8, 1792, 448, 5}	{10, 576, 768, 7}	92				
13/0	Latency	0.54	15.9	$\{4, 256, 2, 128, 64, 2\}$	{12, 1792, 320, 5}	{11, 576, 768, 12}	96				
	Carbon + Latency	0.47	7.3	$\{1, 256, 16, 128, 64, 2\}$	{9, 2048, 384, 4}	{11, 3072, 768, 6}	98				
	Carbon	0.32	4.6	{1, 256, 8, 64, 64, 2}	{6, 1024, 256, 4}	{6, 1536, 384, 6}	27				
2.5%	Energy	0.33	1.8	$\{1, 256, 16, 128, 64, 2\}$	{6, 1024, 256, 6}	{6, 1536, 384, 6}	28				
2.370	Latency	0.46	1.3	$\{4, 256, 16, 128, 128, 2\}$	{6, 1024, 384, 8}	{9, 1536, 384, 4}	43				
	Carbon + Latency	0.31	5.1	$\{1, 256, 4, 64, 64, 2\}$	{6, 1024, 256, 4}	{6, 1536, 384, 6}	27				

**Table 5** The hardware and model architecture properties of each variant of the CarbonCLIP family. Hardware configurations are specified as:  $\{TC, PE_x, PE_y, L2, L2_{bw}, GLB\}$ . Text and Vision encoders are specified as:  $\{Num \text{ Layers}, FFN \text{ Dim}, \text{ Hidden Dim}, \text{ Num Heads}\}$ 

Name	Carbon	Latency	Hardware	M	Avg. Accuracy		
Name	(kgCO2e)	(ms) Configuration		Text Encoder	Vision Encoder	Params (M)	over 41 datasets
				Configuration	Configuration	Tarams (W)	
CarbonCLIP-32-S	0.35	7.3	{1, 64, 16, 64, 64, 2}	{6, 1536, 384, 5}	{10, 3072, 672, 12}	89	46.4
CarbonCLIP-32-M	0.36	15.3	{1, 64, 8, 64, 64, 2}	{8, 1280, 448, 7}	{11, 2688, 768, 8}	99	47.6
CarbonCLIP-32-L	0.38	9.1	{1, 128, 8, 64, 128, 2}	{7,2048,512,6}	{11,3072,768,10}	113	49.1
TinyCLIP-39M/32	0.37	3.0	{1, 32, 32, 64, 32, 2}	{6, 2048, 512, 8}	{12, 2048, 512, 8}	84	45.2
${\rm TinyCLIP\text{-}61M}/32$	0.39	9.4	{1, 128, 8, 64, 64, 2}	{9, 2048, 512, 8}	{12, 2560, 640, 10}	115	47.2
CLIP-B/32 - DataComp	0.42	15.1	{1, 32, 32,64,64,2}	{12, 2048, 512, 8}	{12, 3072, 768, 12}	144	51.1

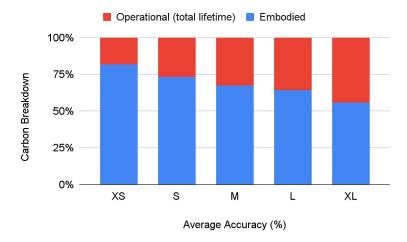


Figure 11 Operational and embodied carbon footprint breakdown for the CLIP-ViT-B/16 architecture.

 $\textbf{Table 6} \ \ \text{Results across all 41 evaluation benchmarks from CLIP Benchmark}$ 

D. L	-	CarbonC	CLIP B/	16 (our	s)	TinyCLI	P B/16	DataComp B/16	Carbo	nCLIP I	3/32 (ours)	TinyCLIP B/32		DataComp B/32
Dataset	XS	S	M	L	XL	39M/16	8M/16	ViT-B-16	S	M	L	40M/32	61M/32	ViT-B-32
cars	0.74	0.82	0.84	0.85	0.87	0.52	0.08	0.89	0.82	0.83	0.84	0.77	0.8	0.87
country211	0.11	0.14	0.16	0.17	0.2	0.18	0.12	0.22	0.15	0.16	0.17	0.13	0.15	0.18
fer2013	0.17	0.21	0.3	0.36	0.34	0.52	0.33	0.39	0.25	0.23	0.38	0.47	0.49	0.33
fgvc aircraft	0.12	0.2	0.23	0.23	0.29	0.15	0.07	0.3	0.2	0.21	0.23	0.14	0.18	0.25
flickr30k	0.55	0.66	0.7	0.7	0.76	0.76	0.52	0.76	0.65	0.67	0.68	0.68	0.71	0.7
flickr8k	0.52	0.62	0.65	0.64	0.7	0.71	0.5	0.7	0.62	0.64	0.66	0.63	0.66	0.65
gtsrb	0.25	0.4	0.46	0.5	0.55	0.32	0.11	0.55	0.47	0.49	0.52	0.38	0.3	0.52
imagenet-a	0.11	0.19	0.26	0.32	0.39	0.33	0.15	0.48	0.22	0.22	0.24	0.17	0.21	0.3
imagenet-o	0.55	0.55	0.51	0.44	0.44	0.49	0.4	0.43	0.49	0.49	0.5	0.52	0.51	0.5
imagenet-r	0.55	0.66	0.73	0.77	0.82	0.7	0.3	0.84	0.72	0.74	0.75	0.7	0.73	0.78
imagenet1k	0.51	0.6	0.63	0.65	0.7	0.63	0.41	0.74	0.62	0.64	0.65	0.6	0.62	0.69
imagenet sketch	0.35	0.45	0.5	0.52	0.57	0.4	0.1	0.6	0.49	0.51	0.52	0.47	0.5	0.57
imagenetv2	0.43	0.53	0.55	0.58	0.62	0.56	0.35	0.66	0.54	0.55	0.58	0.51	0.54	0.61
mnist	0.28	0.58	0.7	0.66	0.7	0.37	0.1	0.76	0.69	0.75	0.71	0.51	0.6	0.81
mscoco captions	0.33	0.41	0.44	0.44	0.49	0.47	0.29	0.49	0.41	0.43	0.44	0.41	0.45	0.45
objectnet	0.37	0.46	0.51	0.54	0.6	0.43	0.19	0.64	0.47	0.5	0.51	0.41	0.44	0.55
renderedsst2	0.51	0.5	0.56	0.54	0.56	0.5	0.5	0.52	0.5	0.51	0.52	0.52	0.54	0.48
stl10	0.92	0.94	0.96	0.97	0.98	0.97	0.92	0.98	0.96	0.96	0.96	0.95	0.96	0.97
sun397	0.59	0.66	0.68	0.69	0.71	0.69	0.56	0.71	0.67	0.67	0.68	0.65	0.67	0.68
voc2007	0.7	0.73	0.75	0.77	0.79	0.77	0.62	0.82	0.77	0.77	0.77	0.77	0.78	0.81
voc2007 multilabel	0.75	0.79	0.81	0.81	0.83	0.82	0.74	0.81	0.79	0.8	0.8	0.76	0.79	0.79
vtab/caltech101	0.8	0.83	0.84	0.84	0.85	0.82	0.72	0.85	0.83	0.83	0.85	0.82	0.82	0.84
vtab/cifar10	0.73	0.83	0.89	0.9	0.93	0.91	0.73	0.96	0.91	0.92	0.92	0.91	0.92	0.96
vtab/cifar100	0.47	0.55	0.65	0.67	0.75	0.68	0.42	0.82	0.69	0.7	0.73	0.69	0.72	0.8
vtab/clevr closest object distance	0.16	0.15	0.17	0.16	0.19	0.2	0.16	0.24	0.16	0.16	0.16	0.17	0.21	0.21
vtab/clevr count all	0.18	0.19	0.21	0.34	0.25	0.2	0.13	0.33	0.16	0.2	0.33	0.19	0.24	0.13
vtab/diabetic retinopathy	0.04	0.09	0.16	0.04	0.2	0.03	0.02	0.11	0.05	0.07	0.05	0.1	0.24	0.42
vtab/dmlab	0.14	0.15	0.15	0.15	0.14	0.13	0.18	0.19	0.21	0.2	0.14	0.21	0.15	0.16
vtab/dsprites label orientation	0.02	0.02	0.02	0.02	0.02	0.02	0.03	0.02	0.02	0.04	0.03	0.02	0.02	0.03
vtab/dsprites label x position	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.02	0.03	0.03	0.03	0.03
vtab/dsprites label y position	0.03	0.03	0.03	0.03	0.03	0.03	0.04	0.03	0.03	0.03	0.03	0.03	0.03	0.03
vtab/dtd	0.41	0.51	0.54	0.51	0.58	0.47	0.29	0.58	0.48	0.51	0.53	0.51	0.52	0.57
vtab/eurosat	0.39	0.48	0.54	0.59	0.58	0.53	0.23	0.59	0.5	0.5	0.56	0.48	0.45	0.57
vtab/flowers	0.55	0.66	0.64	0.66	0.71	0.7	0.58	0.76	0.63	0.66	0.68	0.62	0.64	0.73
vtab/kitti_closest_vehicle_distance	0.37	0.32	0.3	0.26	0.35	0.11	0.15	0.29	0.28	0.19	0.32	0.15	0.17	0.16
vtab/pcam	0.59	0.57	0.59	0.63	0.6	0.61	0.53	0.6	0.61	0.58	0.54	0.52	0.57	0.53
vtab/pets	0.76	0.85	0.87	0.89	0.91	0.81	0.46	0.93	0.87	0.88	0.89	0.85	0.88	0.9
vtab/resisc45	0.45	0.57	0.62	0.64	0.66	0.55	0.21	0.65	0.58	0.64	0.63	0.54	0.58	0.63
vtab/smallnorb label azimuth	0.05	0.06	0.06	0.05	0.05	0.06	0.06	0.05	0.05	0.05	0.05	0.05	0.05	0.05
vtab/smallnorb label elevation	0.11	0.11	0.12	0.11	0.11	0.1	0.12	0.11	0.11	0.1	0.11	0.11	0.11	0.1
vtab/svhn	0.15	0.3	0.29	0.29	0.38	0.16	0.14	0.61	0.33	0.45	0.42	0.35	0.39	0.61

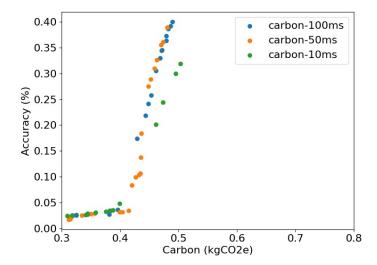


Figure 12 Latency and carbon footprint trade-offs for the CLIP-ViT-B/16 architecture across the Pareto frontier.

## G Case Study: Varying Latency Constraints

We evaluate the impact of latency constraints on model and hardware configurations, as well as their carbon footprints. We categorize use cases into three categories based on latency requirements: critical real-time (<10ms), interactive (<50ms), and non-critical (<100ms) Žádník et al. (2022). Figure 12 shows the Pareto frontiers for each category.

Lower latency constraints typically result in less carbon-efficient designs. For example, a 10ms latency constraint achieves only an 17% carbon reduction compared to latency-optimized models, although with comparable latency values. However, increasing the constraint to 100ms does not consistently improve efficiency, as many optimal designs already meet the 50ms threshold.

## **H** Carbon Footprint of the Framework

We quantify the carbon cost of running CATransformers itself via CodeCarbon Courty et al. (2024): On average, 100 optimization trials emit 57 kgCO2e, and final model training requires 454 kgCO2e per model.