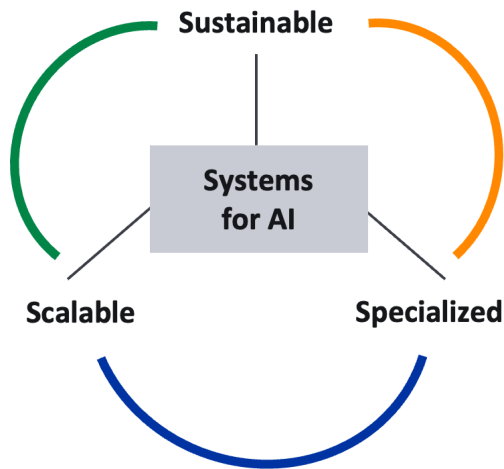


Quantifying the carbon footprint of AI and computing: Past, Present, and Future



Udit Gupta

Assistant Professor, ECE



1.2-2.2 Billion Metric-Tons CO₂

2.1 - 3.9% of worldwide emissions (*Freitag'21*)



On par with the aviation industry's footprint

Computing's emissions are rising given its growing demand!

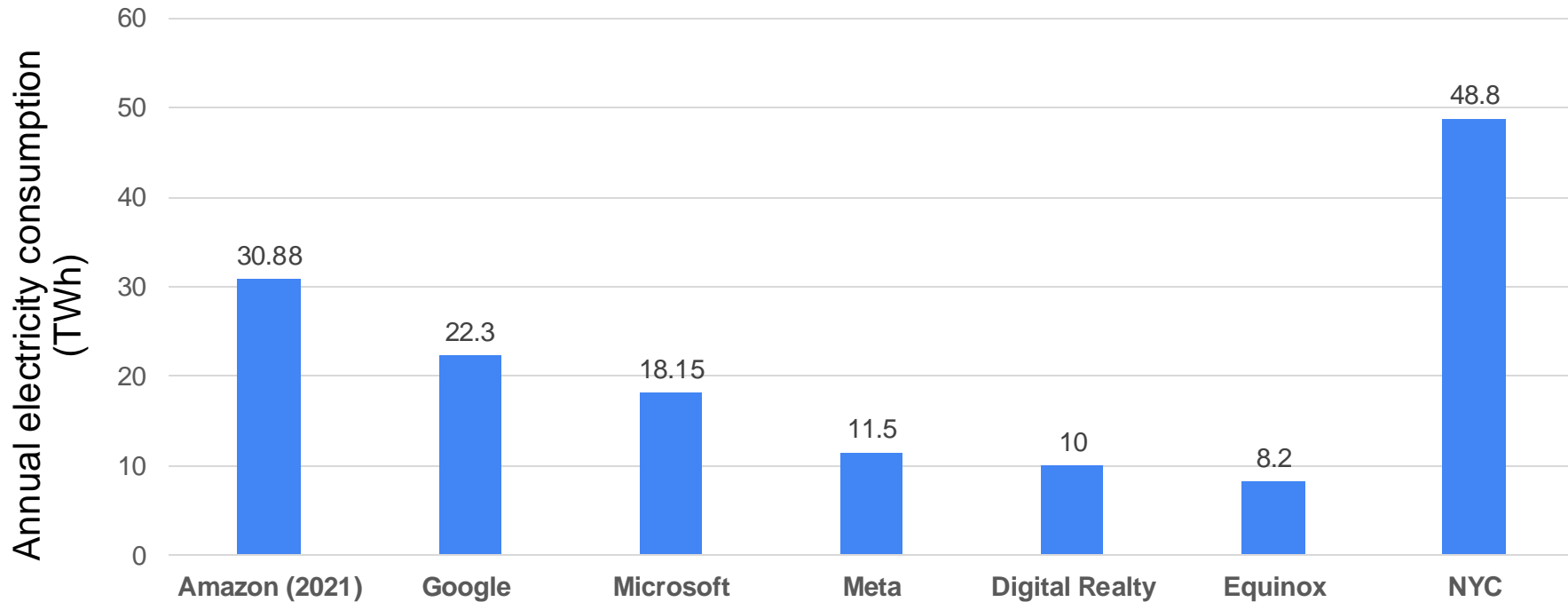
1.2-2.2 Billion Metric-Tons CO₂

2.1 - 3.9% of worldwide emissions (Freitag'21)

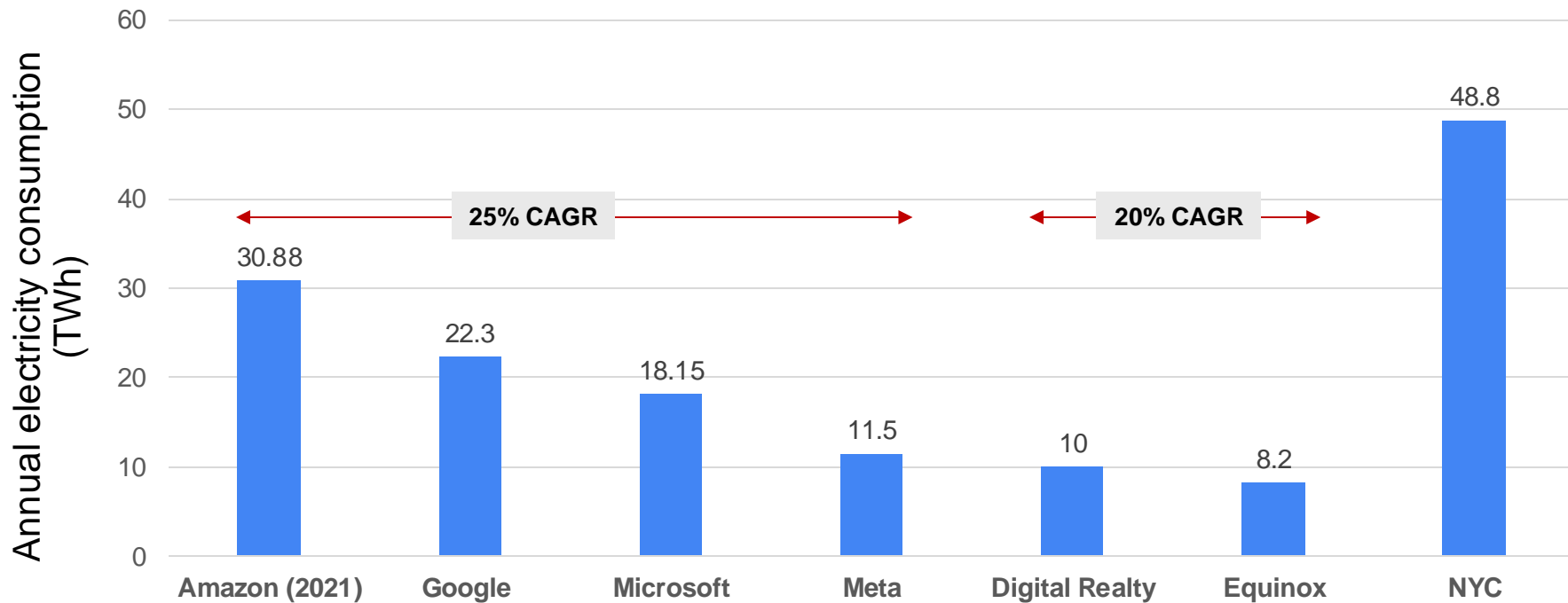


On par with the aviation industry's footprint

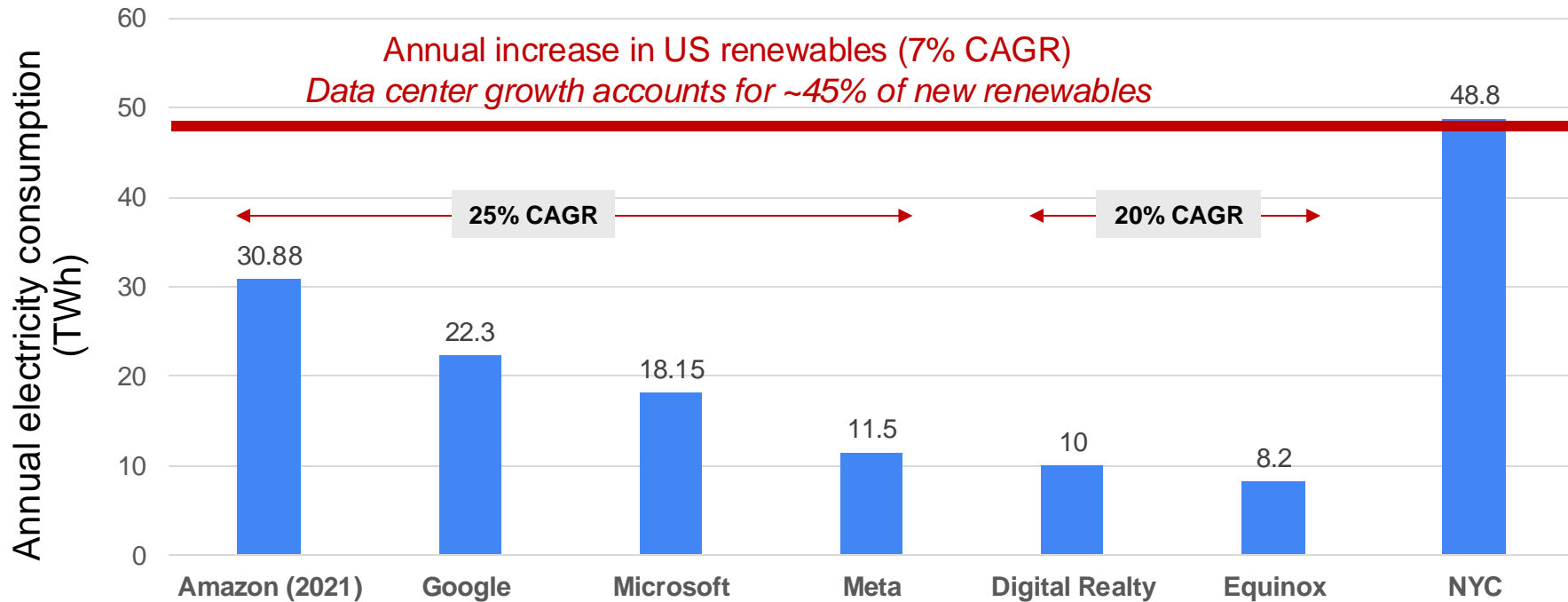
Growing rate of data center energy consumption



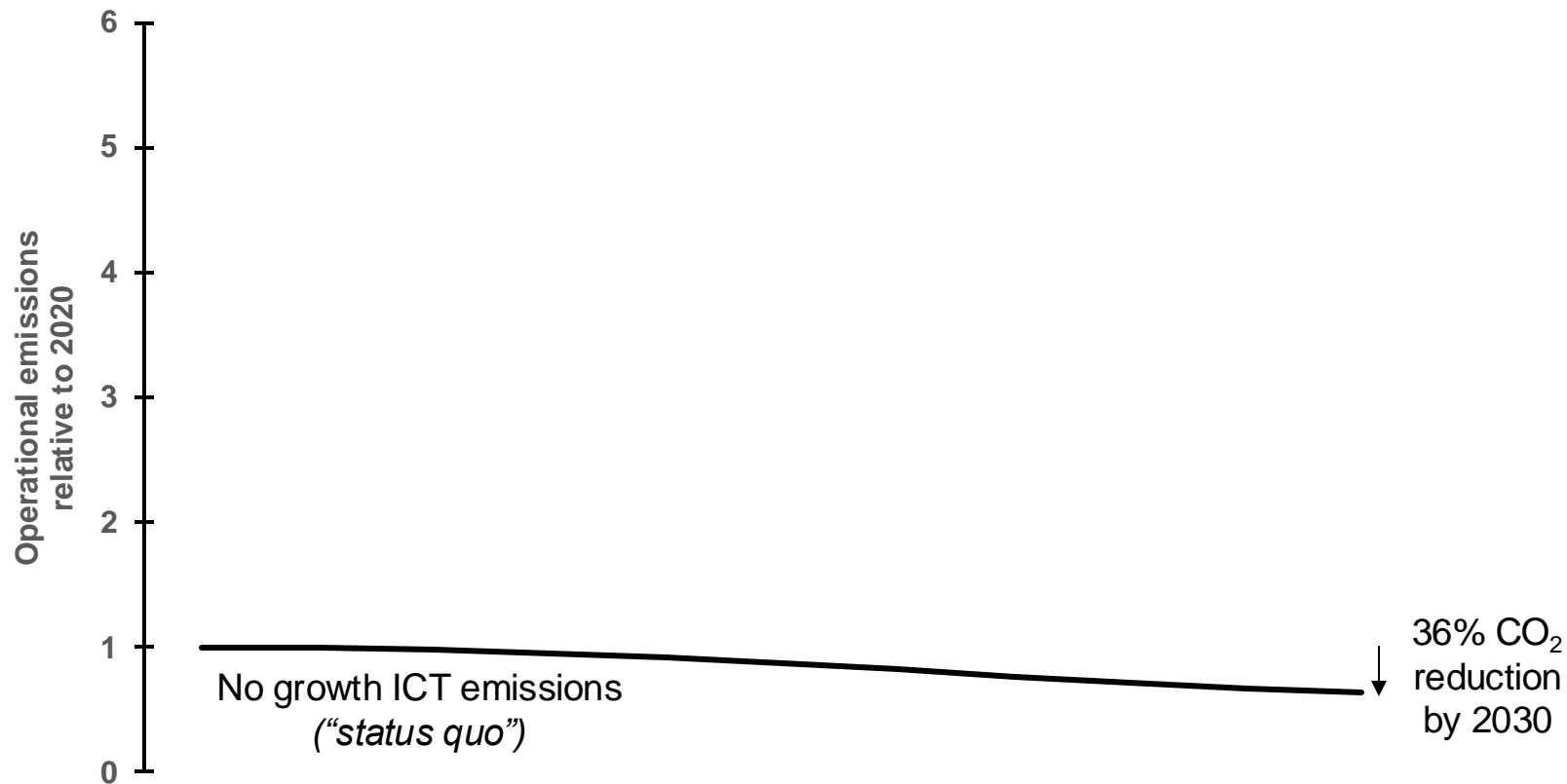
Growing rate of data center energy consumption



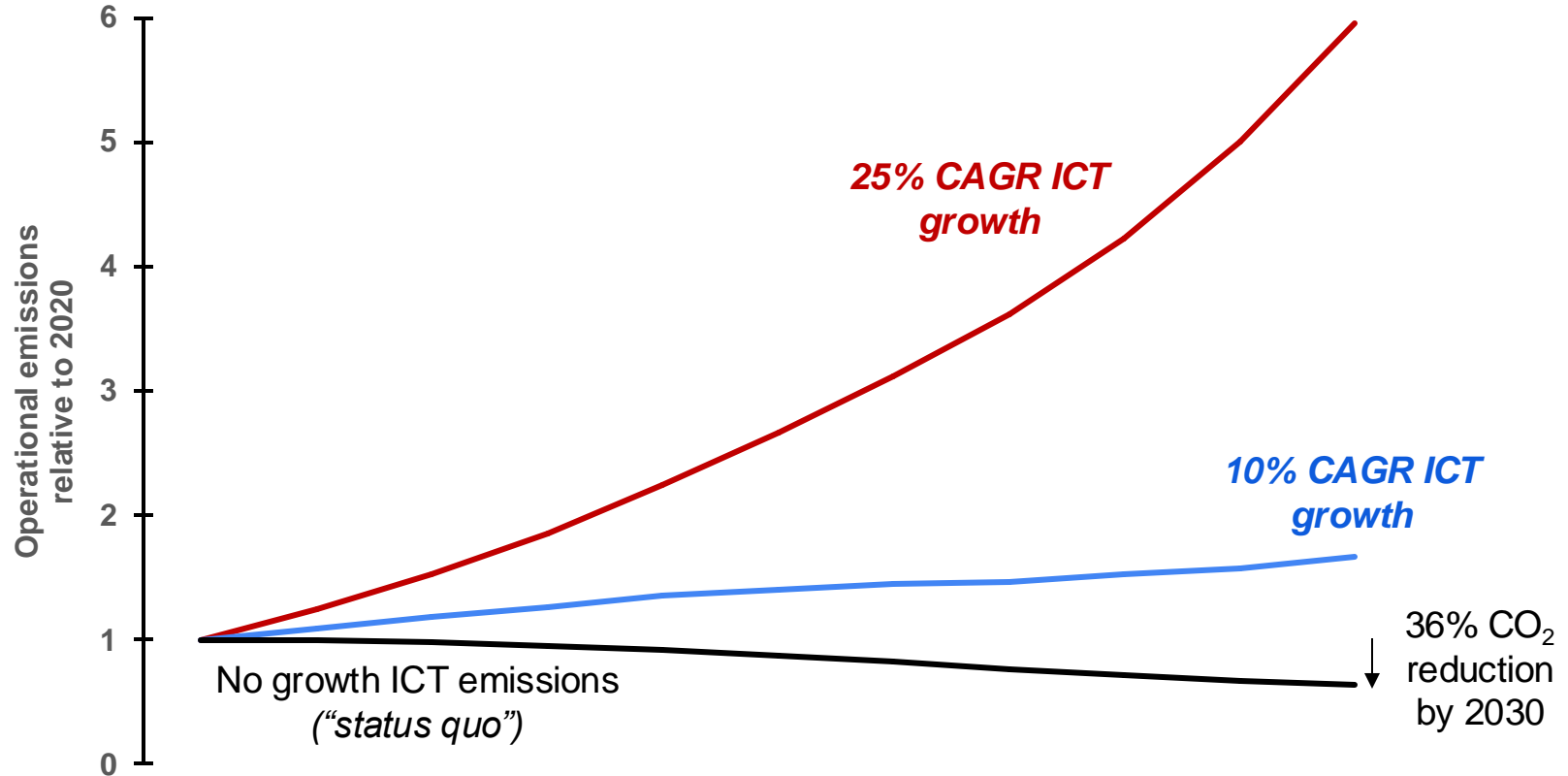
Growing rate of data center energy consumption



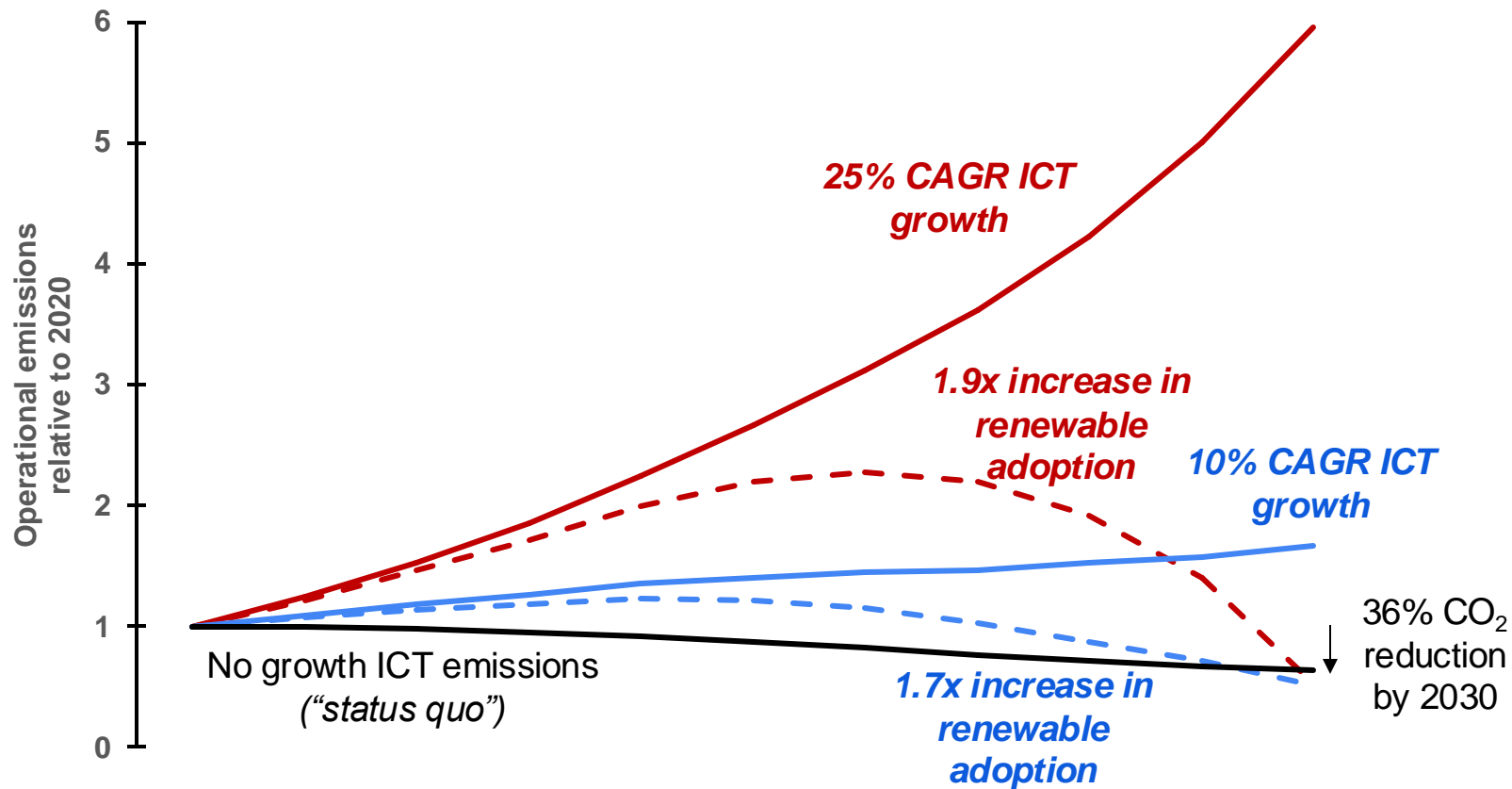
International Telecom. Union targets a 45% reduction in ICT emissions by 2030



International Telecom. Union targets a 45% reduction in ICT emissions by 2030



International Telecom. Union targets a 45% reduction in ICT emissions by 2030



Green
ESG & Investing

Google Is No Longer Claiming to Be Carbon Neutral

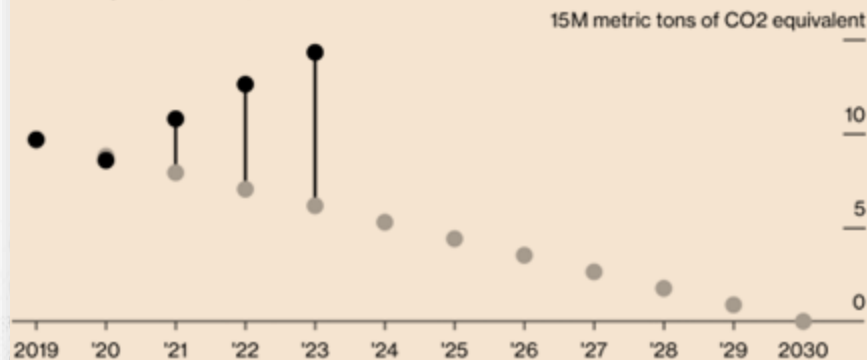
The tech giant, which has seen its planet-warming emissions rise because of artificial intelligence, has stopped buying cheap offsets behind the neutrality claim. The company now aims to reach net-zero carbon by 2030.



Google's Emissions

Artificial intelligence is putting the tech giant's climate goals in peril

● Climate plan (simulated) ● Actual



Source: Google (Scope 1, 2 and 3 data)

Note: Green dots represent linear decline to net-zero emissions goal.

company now aims to reach net-zero carbon emissions by 2030.

The Alphabet Inc. unit has claimed that it's been carbon neutral in its operations since 2007. The status was based on purchasing carbon offsets to match the volume of emissions that were

Interface Inc
14.67 ▲+1.24%

Follow

The AI Race: Why It's So Expensive [Chip Arms Race](#) [Global Energy Strain](#) [DOJ Scrutiny](#) [How Chatbots Work](#)

 Green
Cleaner Tech

Microsoft's AI Push Imperils Climate Goal as Carbon Emissions Jump 30%

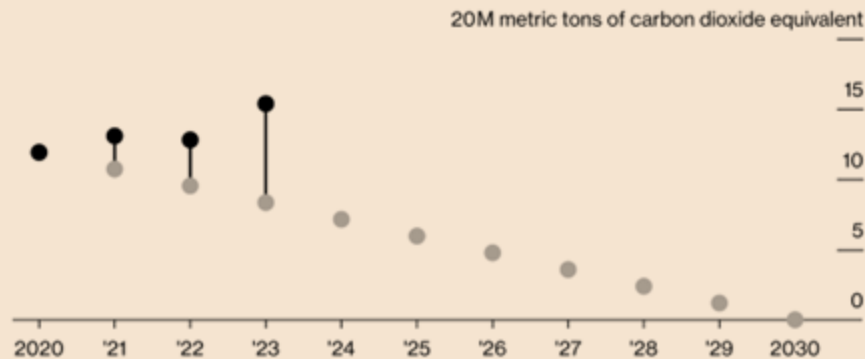
The company's goal to be carbon negative by 2030 is harder to reach, but President Brad Smith says the good AI can do for the world will outweigh its environmental impact.



Microsoft's Emissions

Artificial intelligence is putting the tech giant's climate goals in peril

● Climate plan (simulated) ● Actual



Source: Microsoft (Scope 1, 2 and 3 "management criteria" data)

Note: Green dots represent linear decline to carbon negative goal.

Follow

Amazon.com Inc
199.29 ▼-0.35%

Follow

ambitious and comprehensive plans to tackle climate change. Now the software giant's relentless push to be the global leader in artificial intelligence is putting that goal in peril.

The Seattle-based company's total planet-warming impact is about

Goal:

Mitigate ICT carbon emissions by co-designing solutions across the stack

Goal:

Mitigate ICT carbon emissions by co-designing solutions across the stack

Economics and policy



Goal:

Mitigate ICT carbon emissions by co-designing solutions across the stack

Economics and policy



Education and workforce development



ECE 6960: Sustainable Computing (Spring 2024)

ECE 6960: Sustainable Computing (Spring 2024)

Description

This graduate level course provides an overview of the holistic environmental impact of computing platforms over the course of their lifetime. Topics include life cycle analyses of computing devices, carbon footprint of computing, computer architecture and systems, renewable energy driven data centers, intermittent computing, sustainable applications (e.g., AI), and emerging technologies. We will understand how to evaluate and consider the holistic environmental impact of computing platforms including carbon, water, e-waste, and materials used. Through reading, analyzing, and discussing papers, and an open-ended project students will develop a holistic understanding of the environmental impact of computing and designing sustainable platforms.

Logistics

- Room: Bloomberg Center 91 (Cornell Tech)
- Time: Tuesdays and Thursdays at 1:25pm - 2:40pm
- Please read the [syllabus](#)
- [Gradescope link](#)
- [Zoom link](#)

Course Staff



20 students

- 6 PhD, 11 Master's, 3 Undergraduate

Surveyed a range of topics:

- Metrics, materials, tools, embedded devices, data center power and renewable energy integration, and AI

Roughly 10 final projects:

- 2 final projects exploring integration into **startups**
- At least 3 final projects looking to extend into follow-on **research**

Goal:

Mitigate ICT carbon emissions by co-designing solutions across the stack

Economics and policy



**Education and
workforce development**



**Carbon accounting and
reporting**



Today: Quantifying the carbon footprint of computing

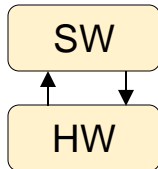
with insights, opening new research and sustainable development opportunities



Understanding the source of computing's emissions

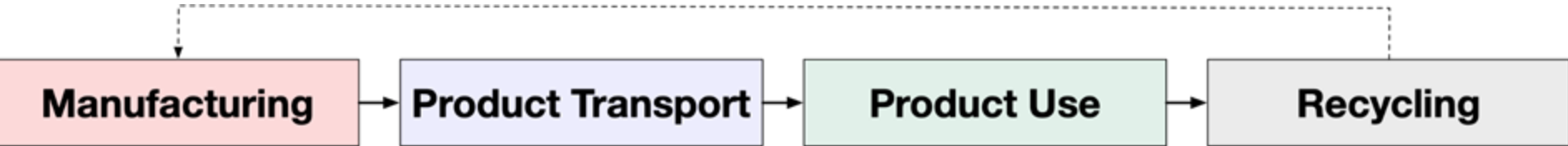


Deep dive: developing computer architectural models to estimate CO₂ emissions

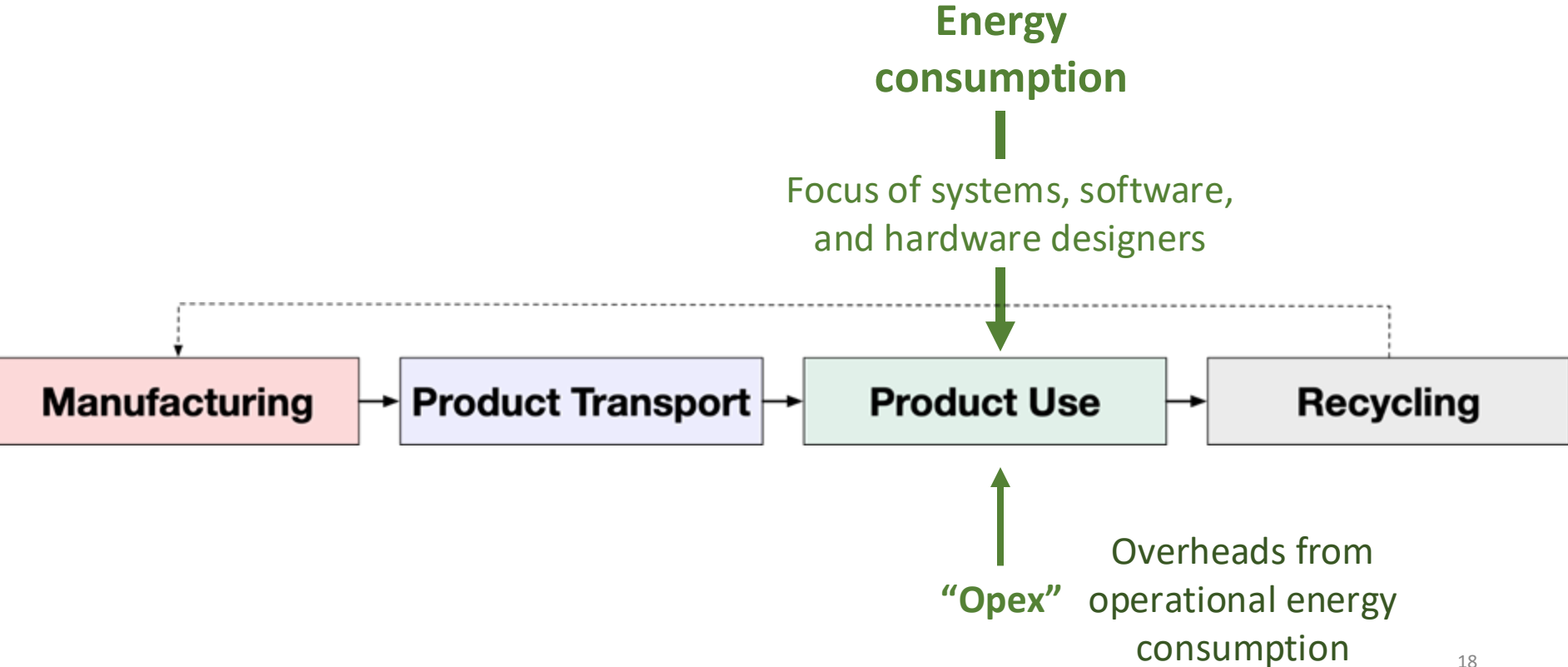


Cross stack: Developing modeling methods across the computing stack

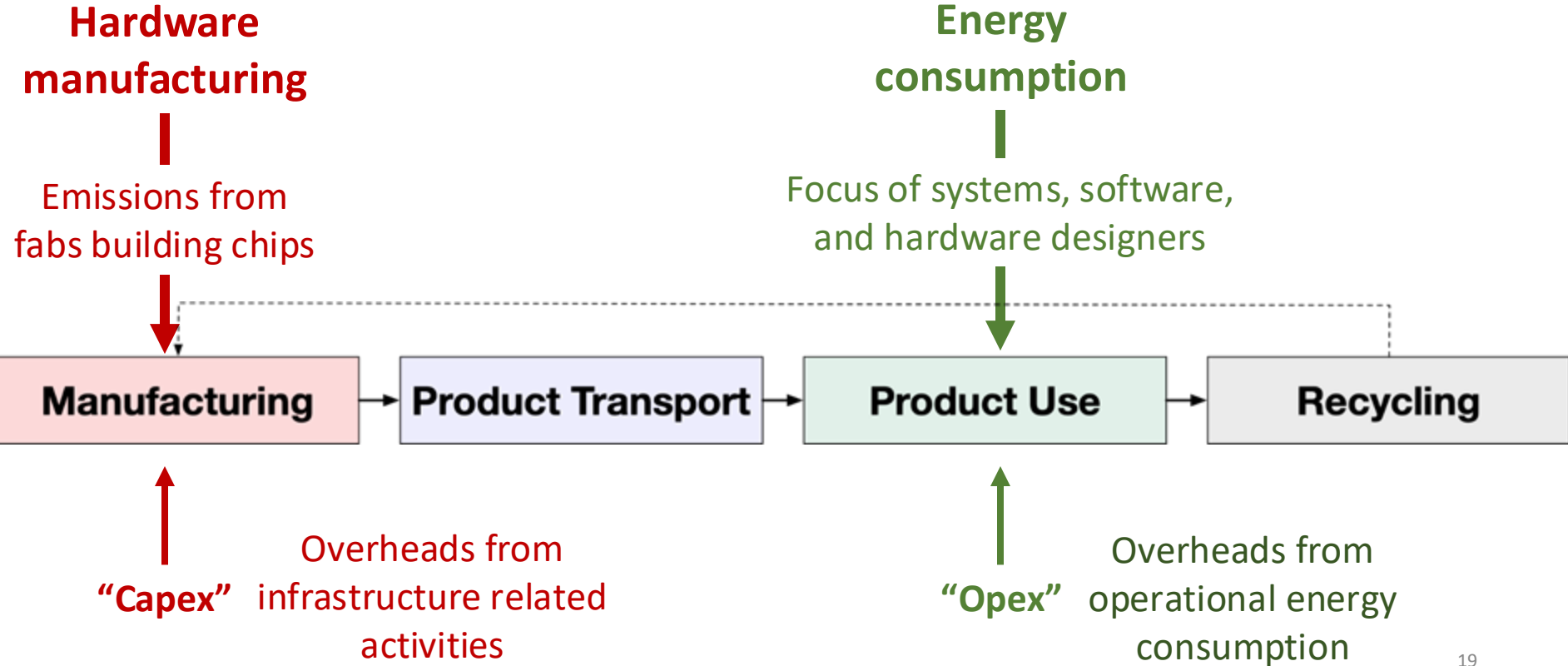
Life Cycle Analysis: key to understanding carbon emissions



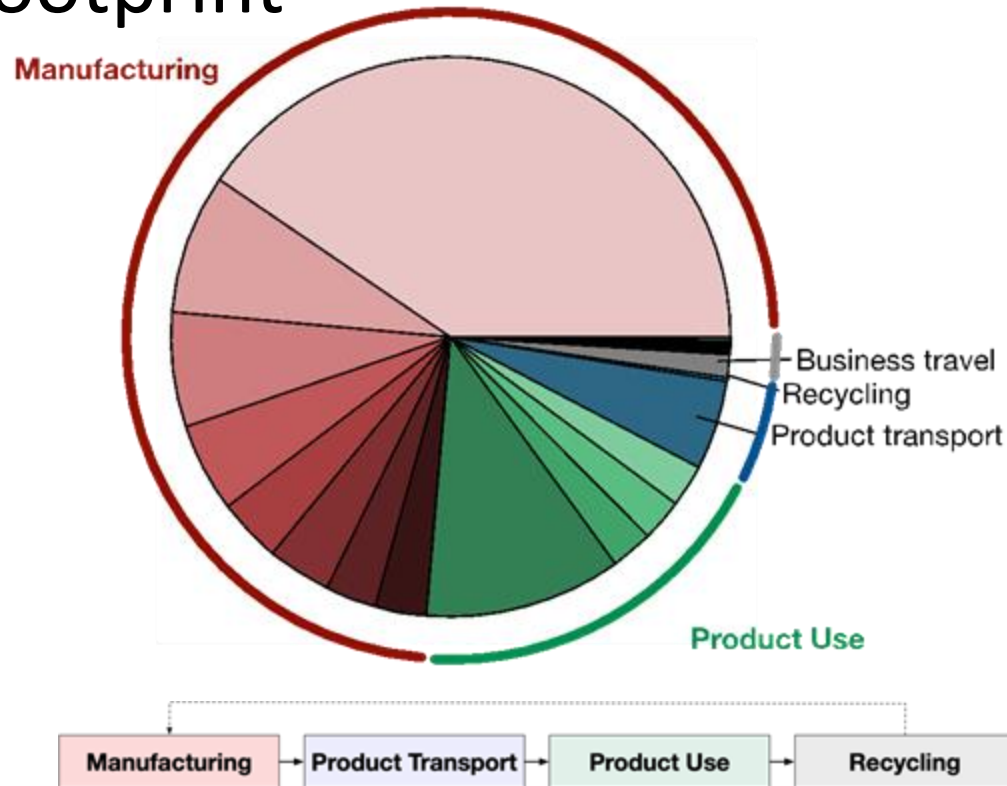
Life Cycle Analysis: key to understanding carbon emissions



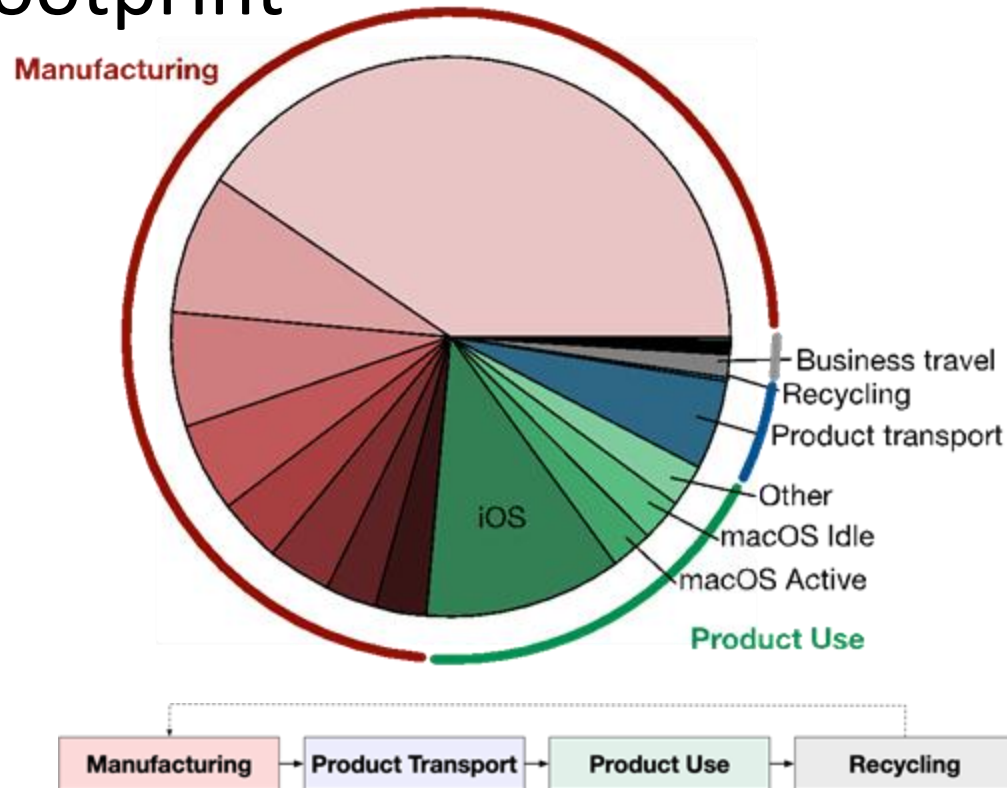
Life Cycle Analysis: key to understanding carbon emissions



Manufacturing dominates Apple's overall carbon footprint

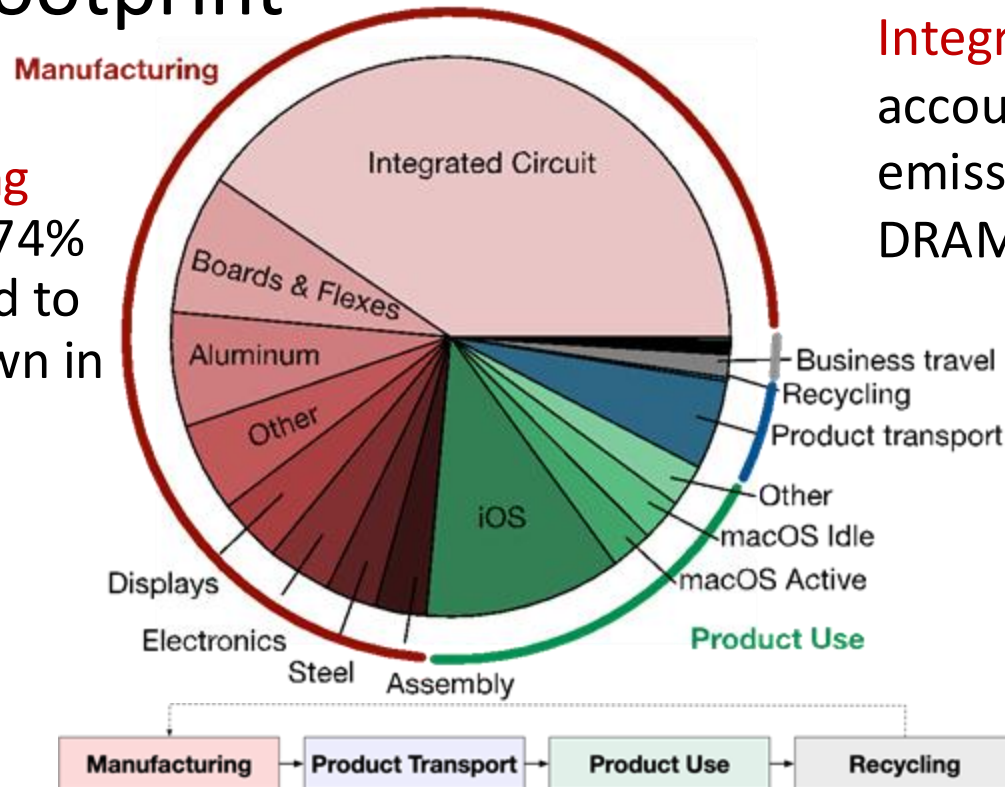


Manufacturing dominates Apple's overall carbon footprint



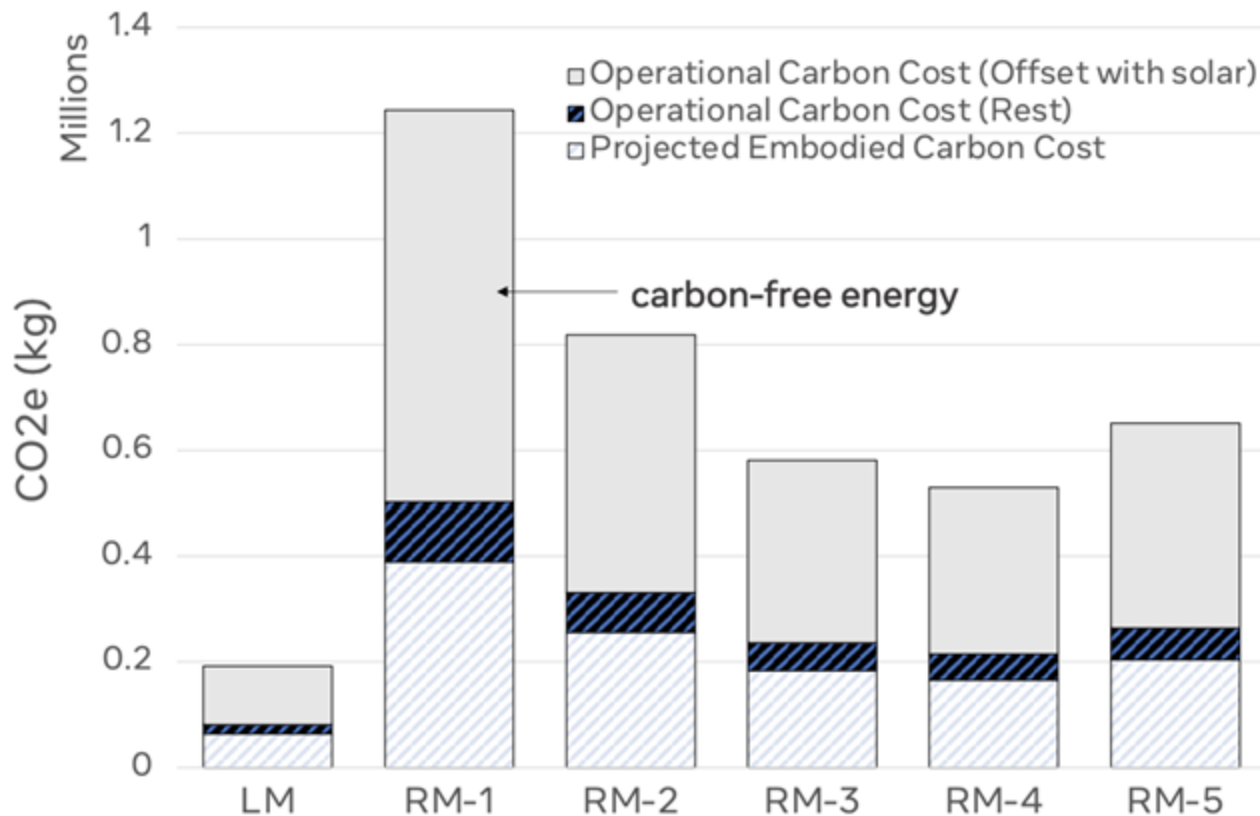
Manufacturing dominates Apple's overall carbon footprint

Manufacturing accounts for 74% of Apple's end to end breakdown in 2019



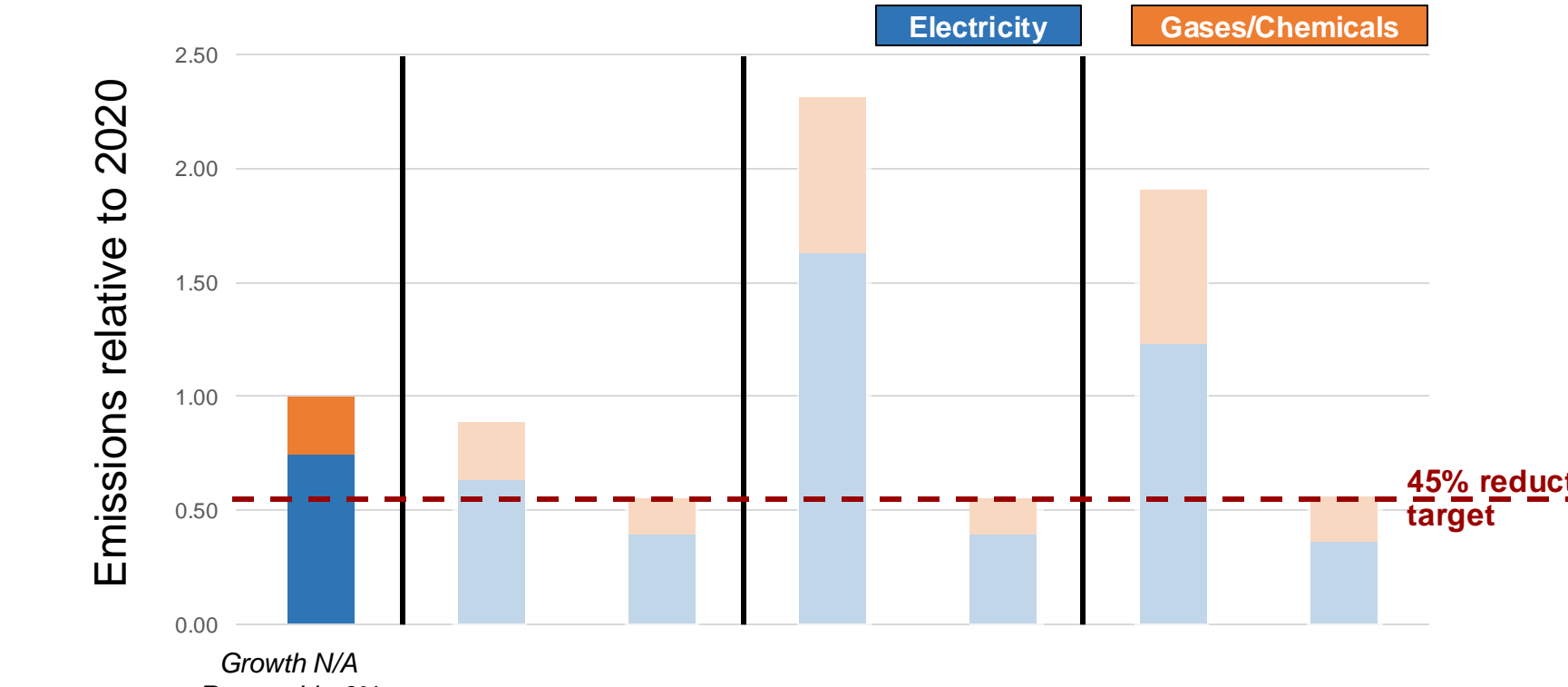
Integrated circuits account for 33% of emissions (SoCs, DRAMs, NAND Flash)

Crucial to look at emissions across HW cycle



The chart displays emissions relative to 2020 for two categories: Electricity (blue) and Gases/Chemicals (orange). The y-axis ranges from 0.00 to 2.50. A dashed red line indicates a 45% reduction target at approximately 0.55. Vertical black lines mark the 2030 and 2050 targets. The scenarios are: Growth N/A, 2030 target, 2050 target, 2030 target, and 2050 target.

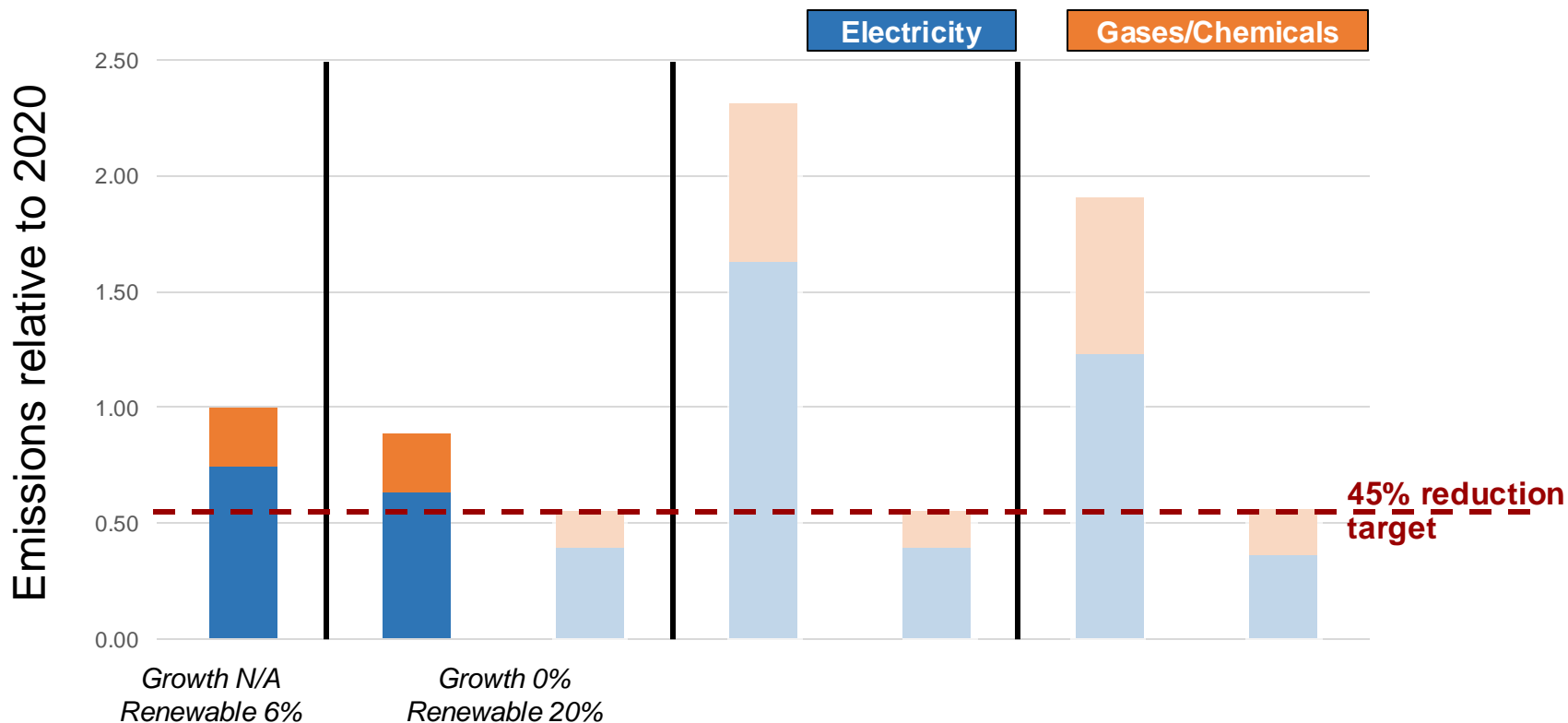
Scenario	Electricity (Blue)	Gases/Chemicals (Orange)	Total Emissions
Growth N/A	0.75	0.25	1.00
2030 target	0.60	0.25	0.85
2050 target	0.40	0.15	0.55
2030 target	1.65	0.65	2.30
2050 target	0.40	0.15	0.55
2030 target	1.25	0.65	1.90
2050 target	0.35	0.20	0.55



“Carbon Connect: An Ecosystem for Sustainable Computing” Lee et. Al. (arxiv 2024)

“Carbon Connect: An Ecosystem for Sustainable Computing” Lee et. Al. (arxiv 2024)

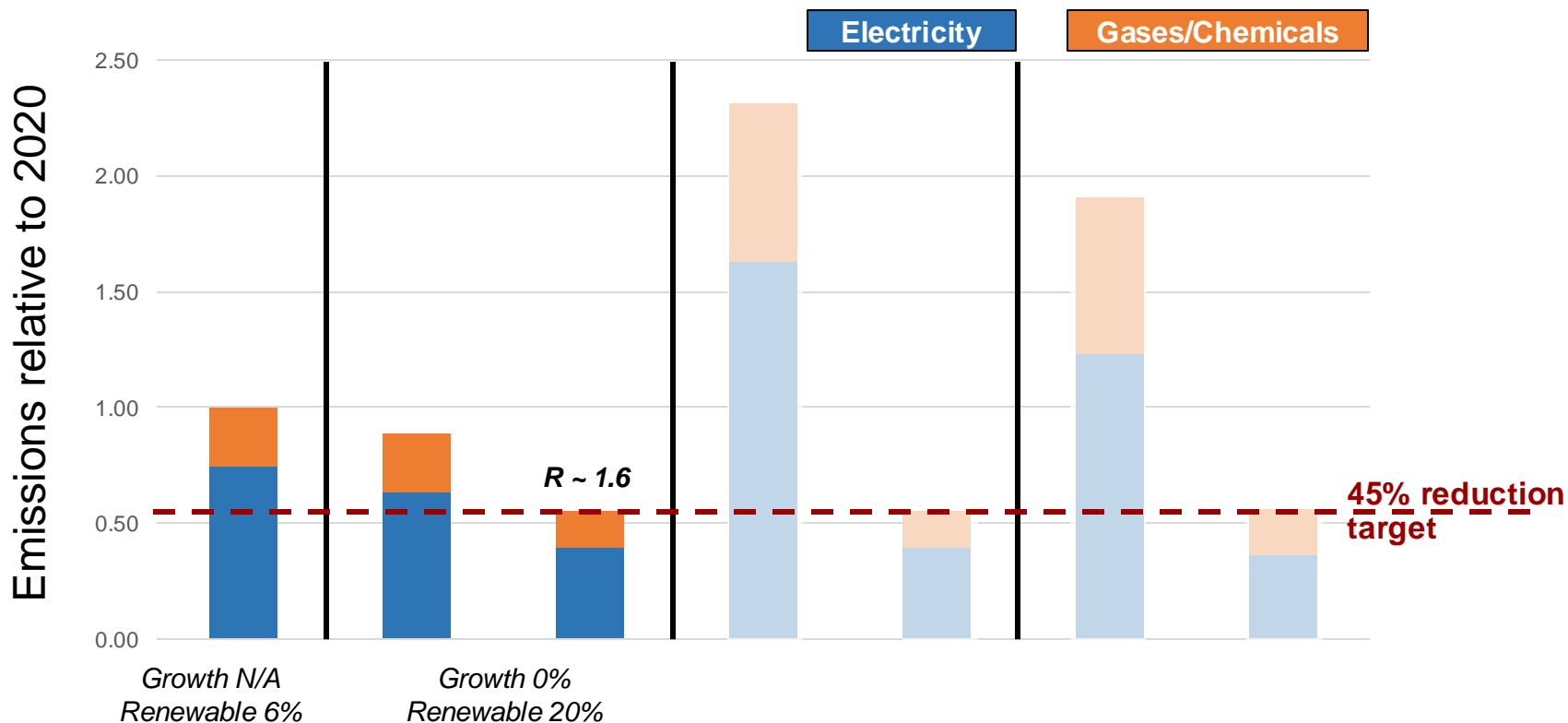
Chip manufacturing requires cross-stack optimization



Source: 2021 corporate sustainability reports

"Carbon Connect: An Ecosystem for Sustainable Computing" Lee et. Al. (arxiv 2024)

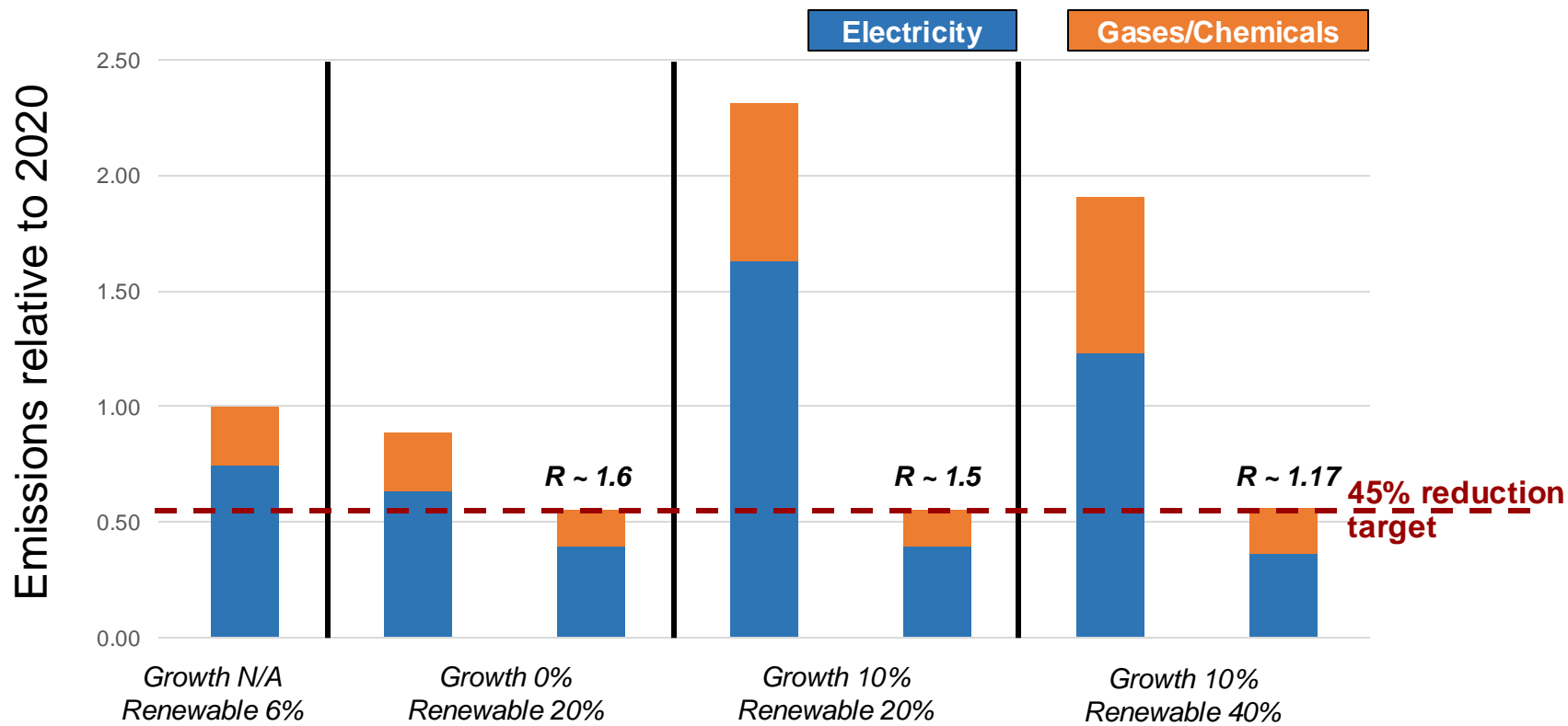
Chip manufacturing requires cross-stack optimization



Source: 2021 corporate sustainability reports

"Carbon Connect: An Ecosystem for Sustainable Computing" Lee et. Al. (arxiv 2024)

Chip manufacturing requires cross-stack optimization



Source: 2021 corporate sustainability reports

"Carbon Connect: An Ecosystem for Sustainable Computing" Lee et. Al. (arxiv 2024)

Today: Quantifying the carbon footprint of computing

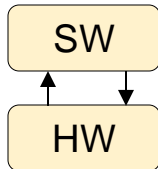
with insights, opening new research and sustainable development opportunities



Understanding the source of computing's emissions



Deep dive: developing computer architectural models to estimate CO₂ emissions



Cross stack: Developing modeling methods across the computing stack

Current carbon accounting methodologies

Economic Input/Output (EIO)



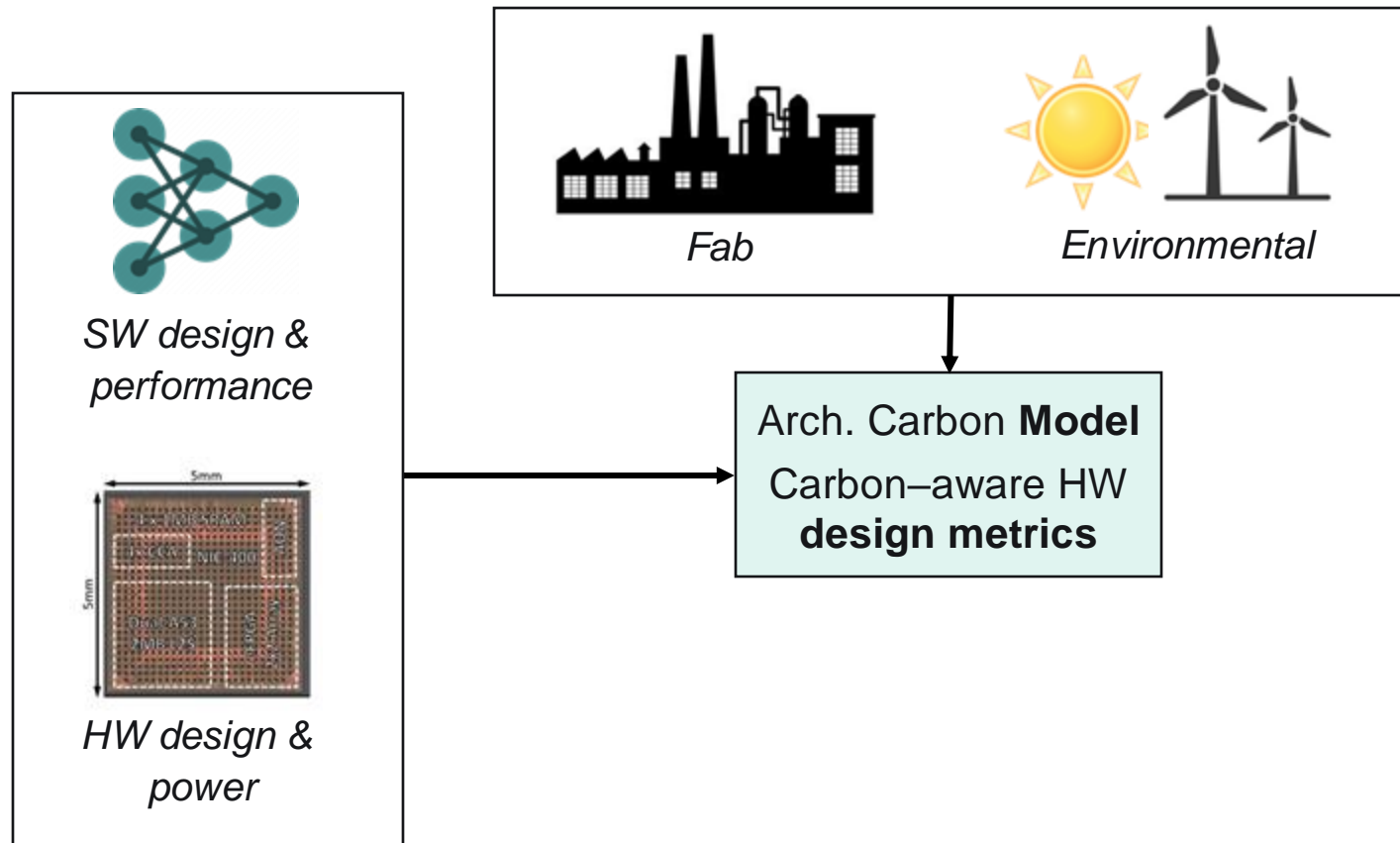
Carbon is tied directly to economic cost which is susceptible to market effects.

Life cycle analysis (LCA)

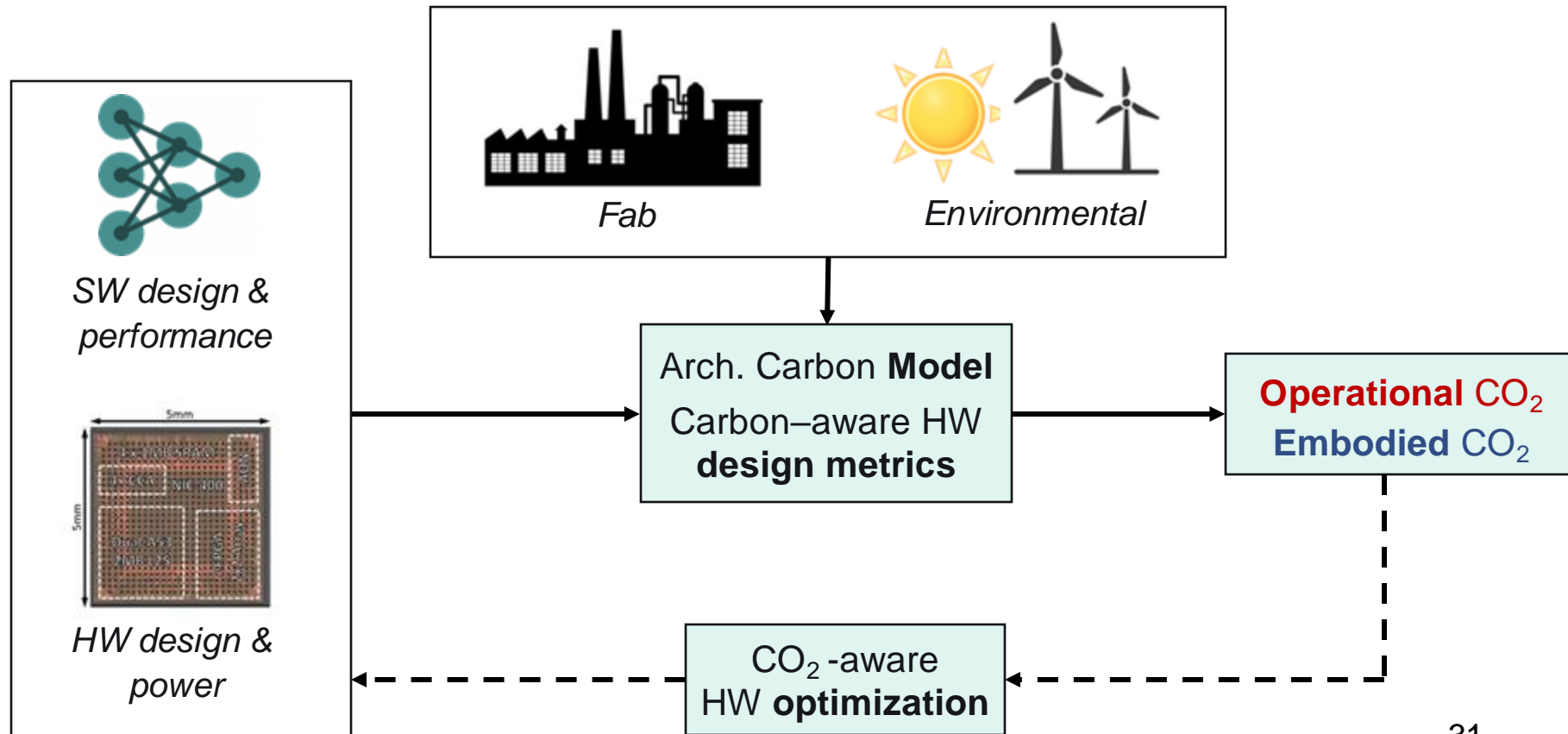


Current databases are out-of-date (45nm or older nodes).
LCA's take high \$\$ and time to conduct.

Architectural Carbon Modeling Tools (ACT)



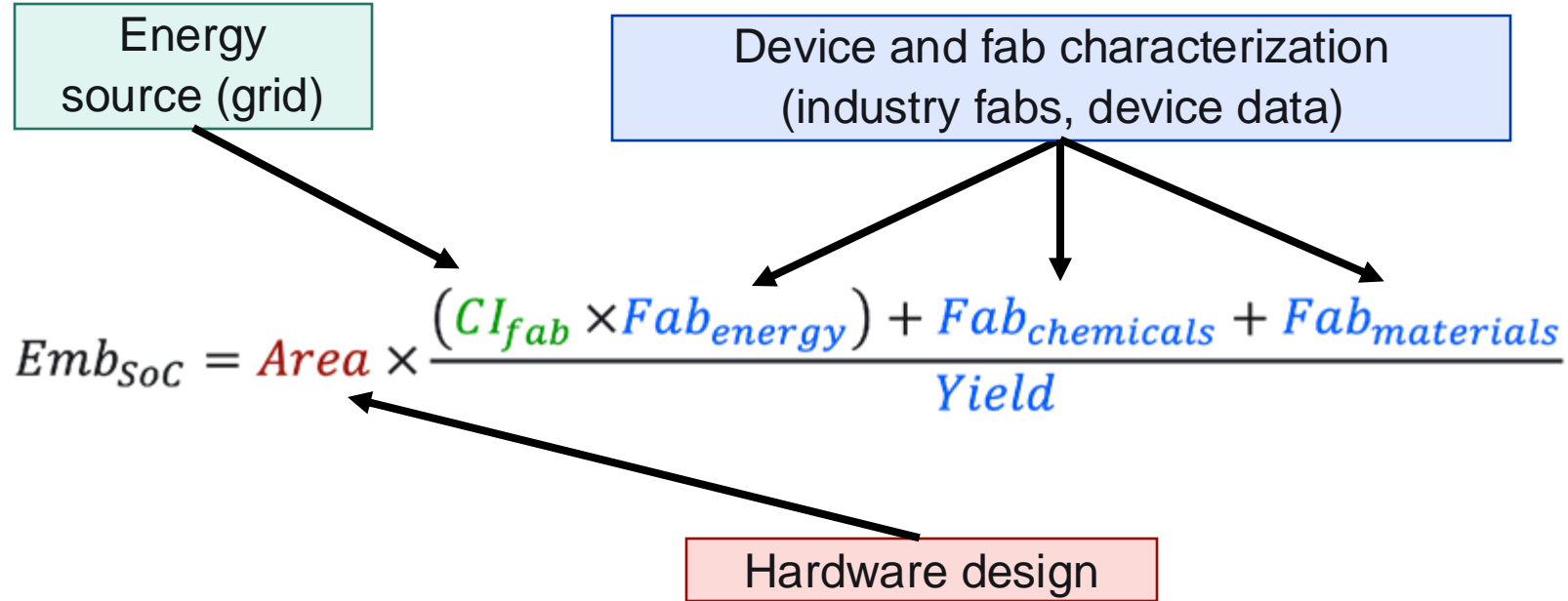
Architectural Carbon Modeling Tools (ACT)



Architectural Carbon Model

Model	Hardware/software input
$Carbon = OP_{CF} + \frac{Runtime}{Lifetime} Emb_{CF}$	Performance/power/energy and lifetime of hardware
$OP_{CF} = CI_{use} \times Energy$	Energy efficiency and environment (carbon intensity)
$Emb_{CF} = Packaging + \sum_r^{SoC, Memory, Storage} Emb_r$	Overhead of hardware manufacturing

Embodied carbon of application processors (SoC's)

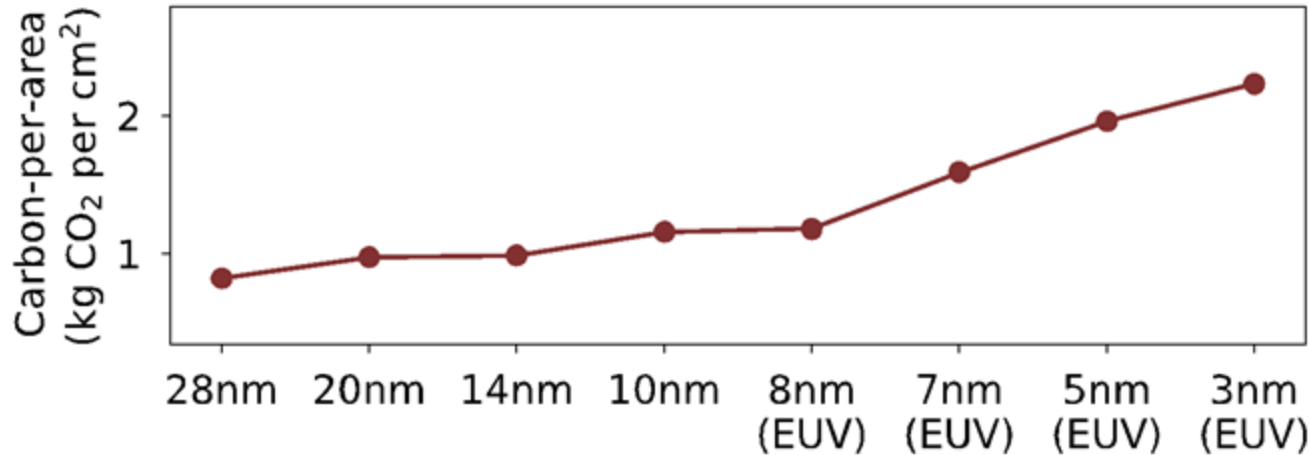


Embodied carbon of application processors (SoC's)

$$Emb_{SoC} = Area \times \textcolor{red}{CPA}$$

Embodied carbon of application processors (SoC's)

$$Emb_{SoC} = Area \times \textcolor{red}{CPA}$$

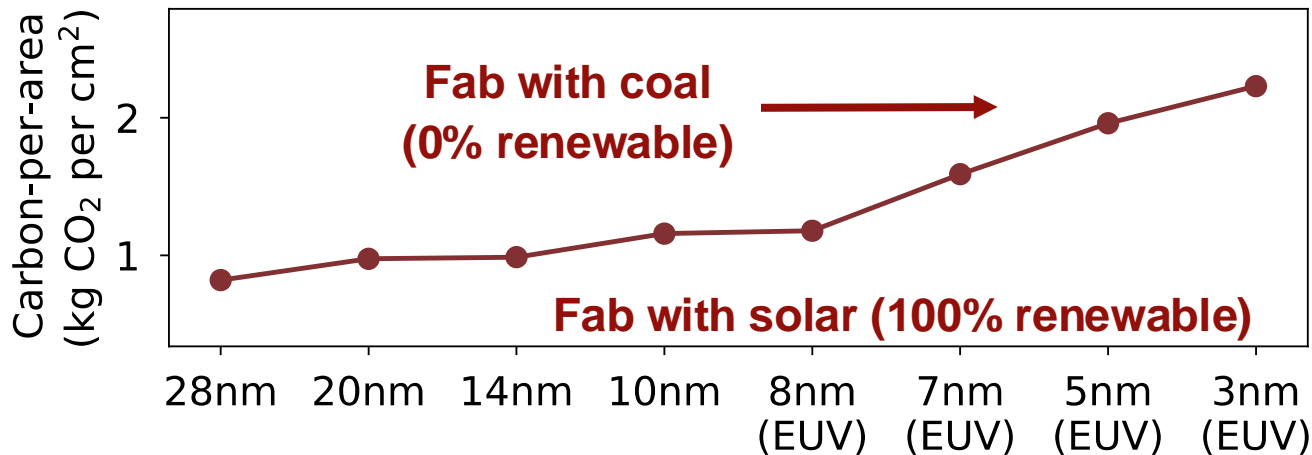


Data sources:

- [IMEC] DTCO including Sustainability: Power-Performance-Area-Cost-Environmental score (PPACE) Analysis for Logic Technologies. Bardon et. al (IEDM 2020)
- [TSMC] TSMC Sustainability Reports 2018-2020

Embodied carbon of application processors (SoC's)

$$Emb_{SoC} = Area \times CPA$$

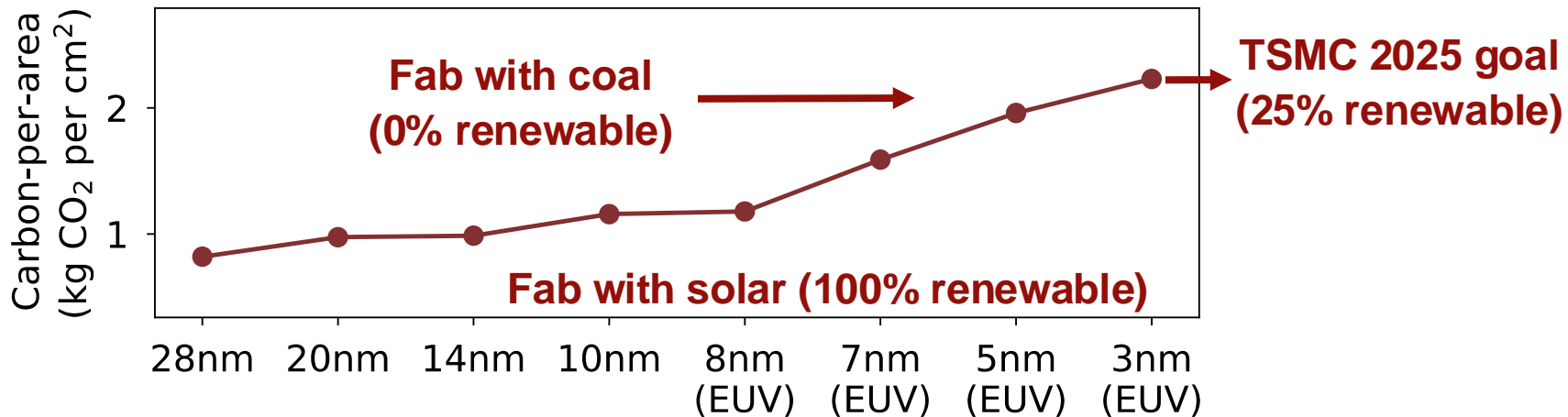


Data sources:

- [IMEC] DTCO including Sustainability: Power-Performance-Area-Cost-Environmental score (PPACE) Analysis for Logic Technologies. Bardon et. al (IEDM 2020)
- [TSMC] TSMC Sustainability Reports 2018-2020

Embodied carbon of application processors (SoC's)

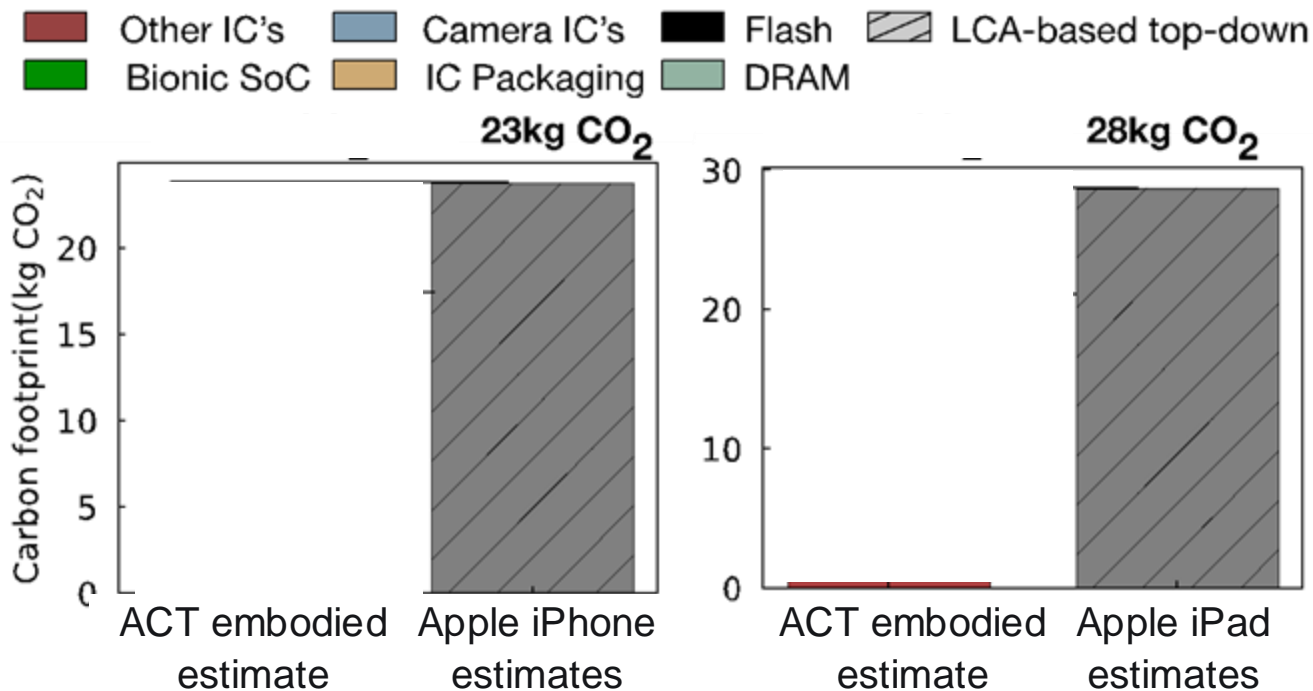
$$Emb_{SoC} = Area \times CPA$$



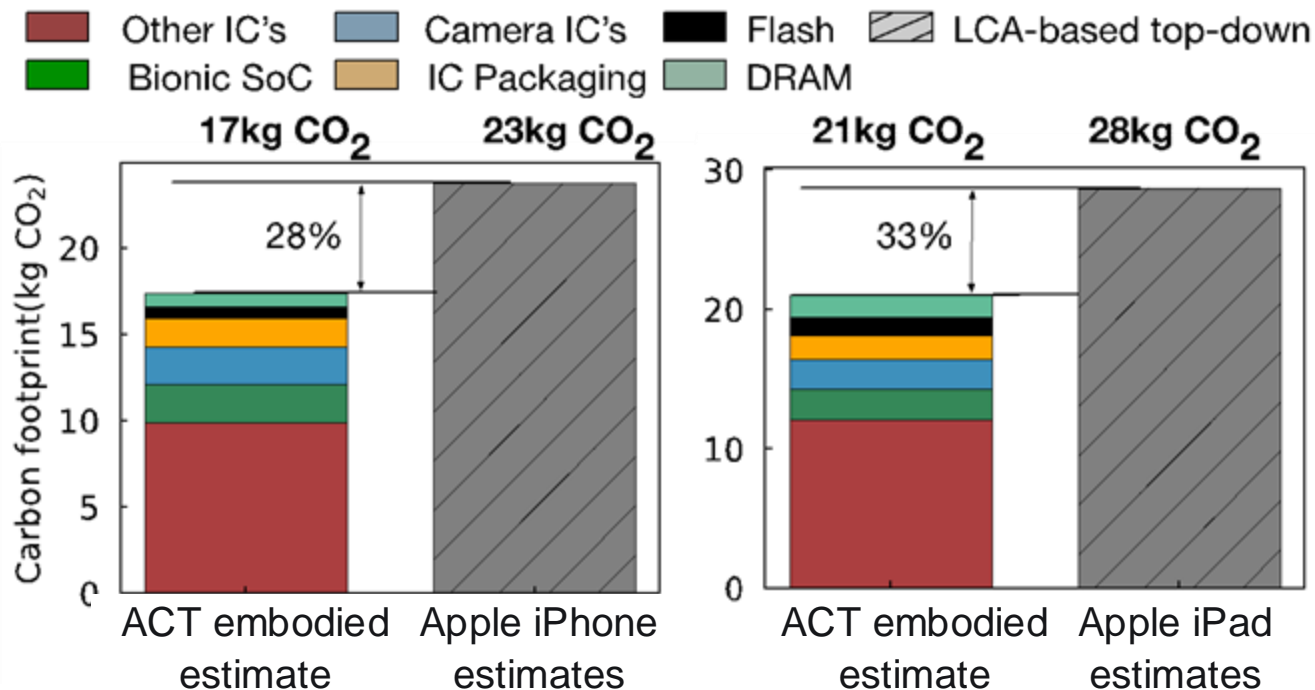
Data sources:

- [IMEC] DTCO including Sustainability: Power-Performance-Area-Cost-Environmental score (PPACE) Analysis for Logic Technologies. Bardon et. al (IEDM 2020)
- [TSMC] TSMC Sustainability Reports 2018-2020

Comparing ACT with Apple's product environmental reports



Comparing ACT with Apple's product environmental reports



Setting the standard for data center sustainability



Understanding the life cycle impact of data center components

We cannot reduce what we do not measure. In 2022, we conducted Life Cycle Assessments (LCAs) on several data center hardware products and developed internal visualization tools to identify the highest carbon emitting components of each product.

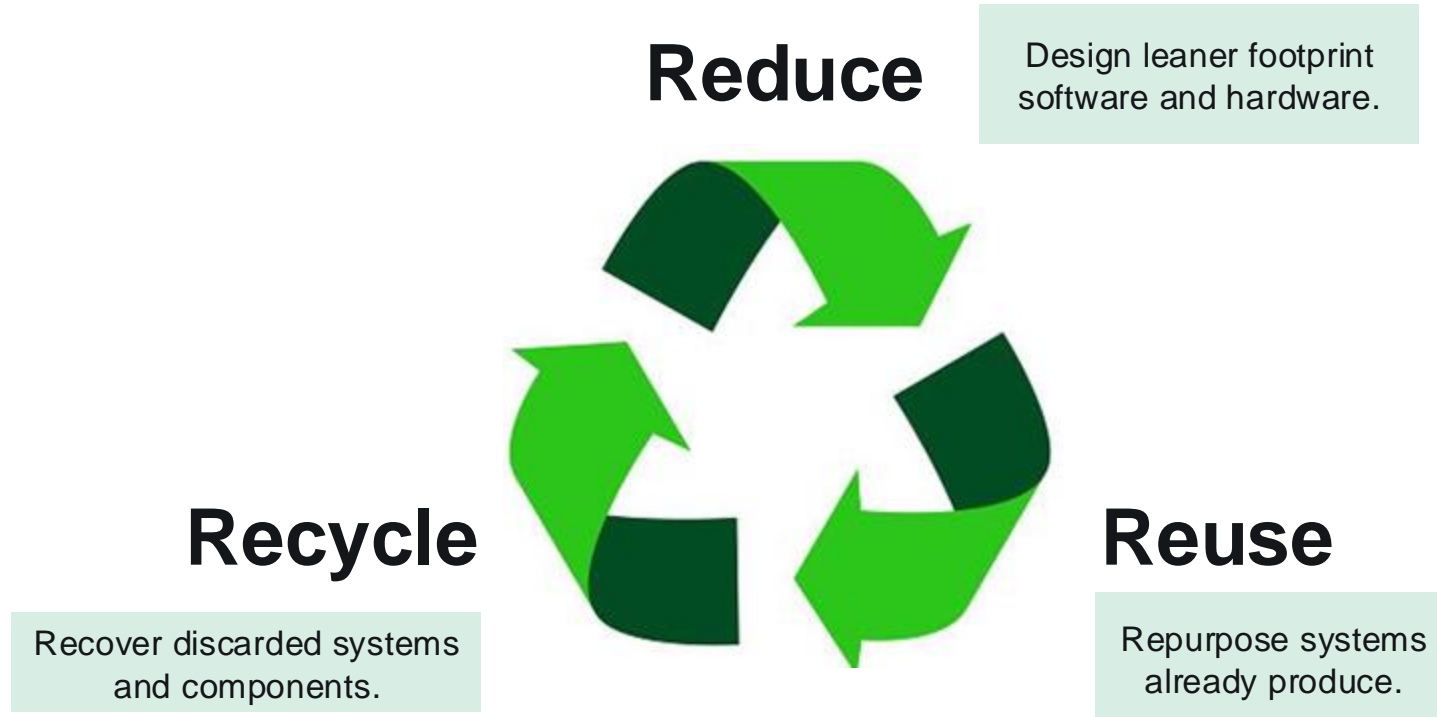
At the data center fleet level, the Sustainability, Physical Modeling, and Meta AI Systems and Machine Learning teams have partnered on a large-scale project to develop and scale a dataset containing the best available

embodied carbon estimates at the scale of the hundreds of millions of components in our data center hardware.

In 2022, the teams reached more than 90% coverage, meaning there is primary data, an LCA, or a [modeled](#) value assigned to each asset. This dataset lays

carbon reductions by helping us use less or choose low-carbon options, engage suppliers, and drive value chain and system-level interventions in line with Meta's net zero strategy.

Tenets of Environmental Design



Tenets of Environmental Design

Reduce

Design leaner footprint software and hardware.



Recycle

Recover discarded systems and components.

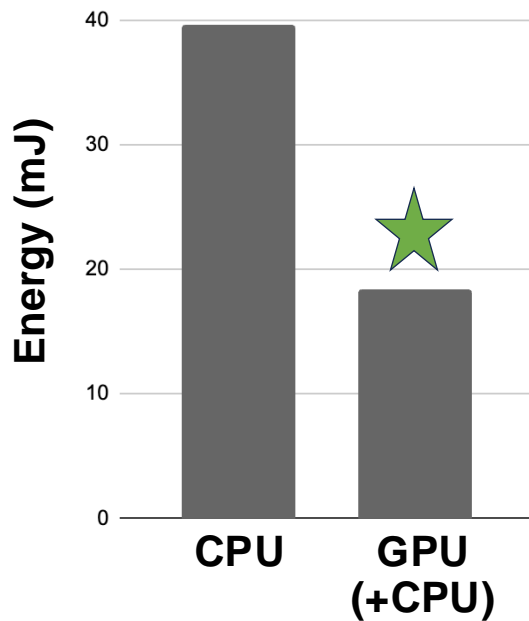
Reuse

Repurpose systems already produce.

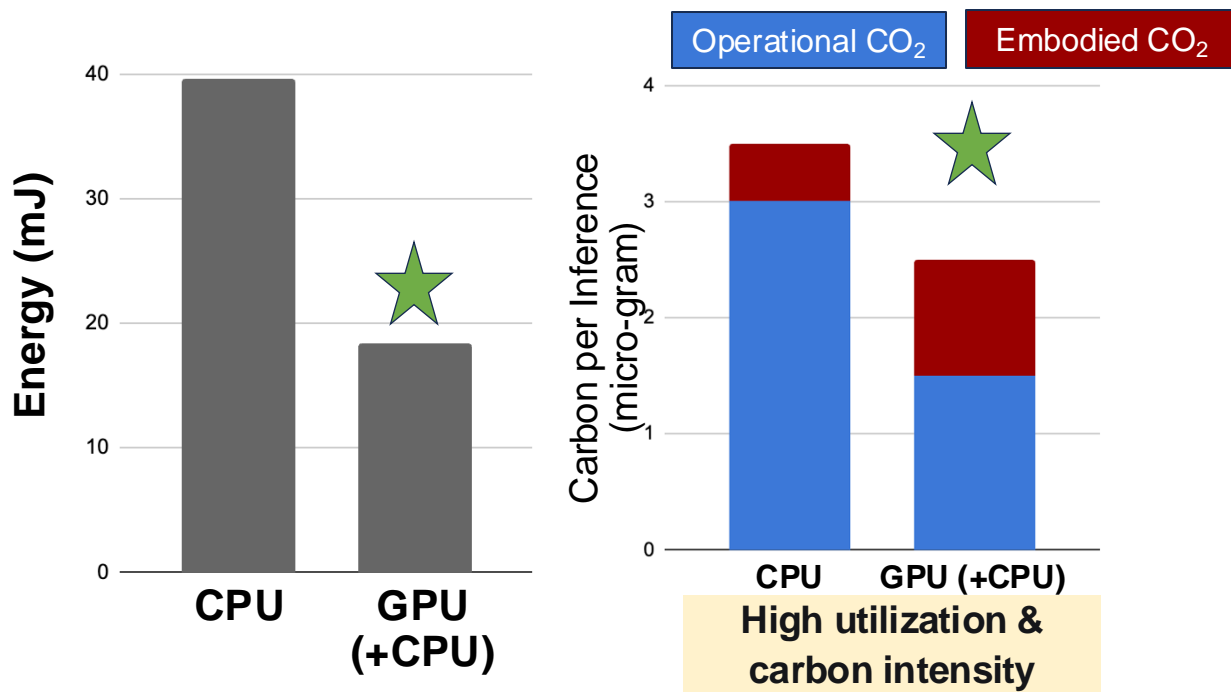
Reuse: General purpose versus custom mobile HW



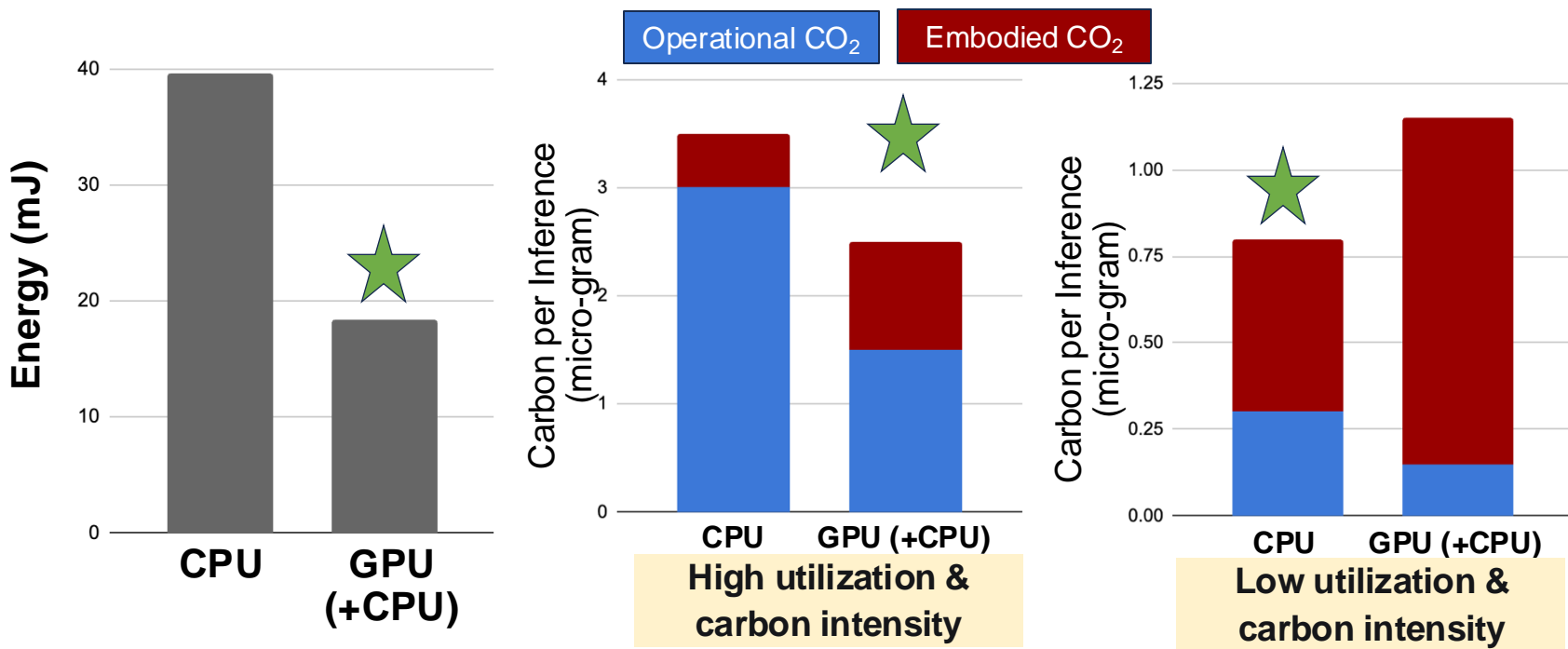
Reuse: General purpose versus custom mobile HW



Reuse: General purpose versus custom mobile HW



Reuse: General purpose versus custom mobile HW



ACT: Architectural Carbon Modeling Tool

ACT An Architectural Carbon Modeling Tool for Designing Sustainable Computer Systems

[View on GitHub](#)

ACT: Architectural Carbon Modeling Tool

Abstract

Motivation: Over the past two decades, the world has witnessed a dramatic rise in computing across data centers, mobile, and communication technologies. As of 2015, information computing technology (ICT) accounts for up to 3% of global carbon emissions. Unfortunately, as the demand for computing grows with new applications and platforms so does energy demand. Many technology companies, including Microsoft, Google, Facebook, and Amazon, have pledged to reduce their carbon footprints over the next decade. Meeting these pledges and enabling sustainable computing requires immediate action from the systems and hardware community.

Background: In addition to its convenience, enabling environmentally sustainable computing introduces unique challenges for system hardware designers. First, carbon emissions are shifting from being dominated by open energy consumption to hardware manufacturing. Traditionally the majority of emissions come from operational hardware use (i.e., energy consumption). However, given the energy efficiency optimizations and increasing fabrication complexity, the majority of carbon emissions have shifted to hardware manufacturing. Given these new challenges, enabling environmentally sustainable computing demands distinct solutions across the computing stack, hardware life cycles, and end-to-end systems.

Come join us at MICRO
2024 for the next ACT
tutorial!

Our giganon aspiration

We believe that Google has a unique opportunity that extends beyond reducing the environmental impacts of our own operations and value chain.²¹

In 2020, Google shared our aspiration to help others reduce 1-gigaton (GT) of their carbon equivalent emissions annually by 2030. This is an aspiration that we believe we can help achieve by contributing meaningfully beyond our own operations and value chain.

We initially focused on helping cities and local governments reduce their carbon footprint. Our effort is the Environmental Insights Explorer, which provides actionable climate and sustainability data to government officials. It has been used in multiple ways across the globe, including by city leaders in Dubai to inform smart transportation policies, and by the city of Austin to prioritize planting trees in areas with the highest need.

To better reflect the broader group of partners we aim to help, we're updating our shared ambition:



For context, 1 GT is comparable to the entire annual emissions of Japan.²² Helping others to reduce 1 GT of carbon equivalent emissions per year, starting in 2030, is a bold aspiration focused on where we can have the most impact—enabling others to reduce emissions in key areas like energy and transportation. Our ultimate measure of success will be how much we've helped individuals, cities, and other partners to achieve their own greenhouse gas

Many of the solutions to achieve a reduction of carbon emissions are complex and require a lot of resources. This pushes us to iterate and be meticulous in our approach. We'll share progress and learnings along the way.

Estimating impact

Estimating the carbon impact of actions taken by many millions of people, communities, and organizations will be inherently difficult, imprecise, and fundamentally different from measuring a corporate carbon footprint. However, it's also useful to enable us to prioritize the most helpful solutions for others.

After reviewing emerging best practices and applying our expertise to advance our mission, we've developed a set of approaches and knowledge that will inform how we estimate

Approaches

Consider carbon accounting principles

Established carbon accounting principles (such as well-defined baselines and true and fair representation of data) provide helpful insights as we develop estimation methodologies for enabled emissions reductions.

Quantify and evaluate real world action

The data available to us from our technology, products, or services may be several steps removed from actual real-world actions and impact that resulted in reduced emissions. We'll use our best judgment to evaluate the effect of those actions.

Challenges

Enabling emissions reductions is a complex and multifaceted challenge. Emissions reduction efforts don't happen in isolation. For example, a city's transportation emissions may be reduced by a city's investment in public transit, but this may not be considered when estimating enabled emissions reductions.

Inherent uncertainty

Uncertainty is inherent to most GHG accounting methodologies and results, and it increases when considering enabled emissions reductions due to a lack of primary data and precise information about real-world actions and their effects. However, understanding the sources, types, and magnitude of uncertainty is crucial to deploy conservative estimates, inform improved data inputs, and properly interpret results.

While carbon accounting principles and these concepts are a good start towards estimating enabled emissions reductions, they are not perfect. Emissions will rapidly evolve, and we welcome the opportunity to collaborate with others to improve our practices.

Supporting individuals and working together to reduce emissions are examples of how we're already supporting and enabling individuals and partners to

“Inherent Uncertainty: Uncertainty is inherent to most GHG accounting methodologies and results, and it increases when consideration enabled emissions reductions due to a lack of primary data and precise information about real-world actions and their effects. However, understanding the sources, types, and magnitude of uncertainty is crucial to deploy conservative estimates, inform improved data inputs, and properly interpret results.”



The Environmental Insights Explorer makes actionable climate data available to more than 40,000 cities and regions worldwide.

Count

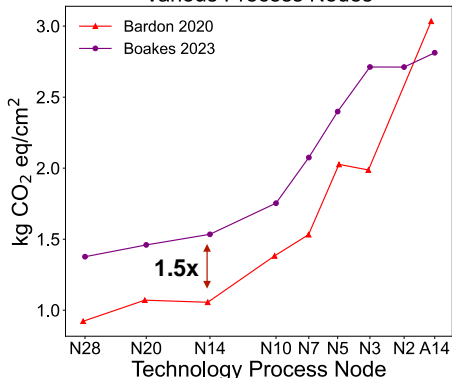
Storage Embodied Factor (SEF)

SEF = $\frac{\text{kgCO}_2}{\$B}$

> 4x

Uncertainty is inherent in carbon accounting

Total Carbon Emissions for Various Process Nodes



Open research questions:

- What magnitude uncertainty exists across all IC components?
- What degree of uncertainty exists in embodied versus operational carbon?
- How do we consider uncertainty in carbon-aware hardware design to enable robust sustainable computing decisions?

The Dirty Secret of SSDs: Embodied Carbon

SWAMIT TANNU, University of Wisconsin, Madison, USA
PRASHANT J. NAIR, University of British Columbia, Canada

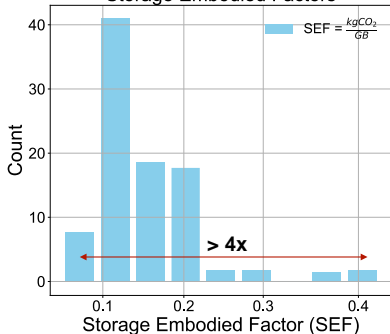
Scalable Solid State Drives (SSDs) have ushered in a transformative era in data storage and accessibility, powering both consumer and enterprise devices. However, the strides made in scaling this technology can have significant environmental consequences. In the global race to develop sustainable semiconductor manufacturing relies on electricity derived from coal and natural gas. A striking example of this is the manufacturing process for a single Gigabyte of Flash memory, which currently amounts to 1.5 kg of CO₂ – a considerable fraction of the total carbon emissions attributed to the system. Remarkably, the manufacturing of storage devices alone contributed to an estimated 10 million metric tonnes of CO₂ emissions in the year 2021. In light of these environmental concerns, this paper delves into an analysis of the sustainability trade-offs inherent in Solid State Drives (SSDs) when compared to traditional Hard Disk Drives (HDDs). Moreover, this study proposes sustainability to gauge the embodied carbon associated with storage systems effectively. The research encompasses five key strategies to enhance the sustainability of storage systems.

Firstly, the paper offers insightful guidance for selecting the most suitable storage medium, be it HDD or SSD, considering the broader societal impact. Secondly, the paper advocates for implementing techniques that reduce the lifespan of SSDs, thereby mitigating premature device replacement and their attendant environmental toll. Thirdly, the paper emphasizes the need for efficient recycling and reuse of high-density media and hard drives, underscoring the significance of minimizing electronic waste.

Lastly, for bandwidth devices, the research examines the potential of harnessing the electricity offered by direct storage solutions as a means to curtail the ecological repercussions of localized data storage. In summation, this study critically addresses the embodied carbon issues associated with SSDs, comparing them with HDDs, and proposes a comprehensive framework of strategies to enhance the sustainability of storage systems.

CC0 Creative Commons - Social and professional topics - Sustainability - Hardware - External storage - Applied computing - Data centers, datacenters
Additional Key Words and Phrases: Embodied Carbon, Solid State Drives, Hard Disk Drives, Sustainability
ACM Reference Format
Swamit Tannu and Prashant J. Nair. 2023. The Dirty Secret of SSDs: Embodied Carbon. In 1 (October 2023), 4 pages. <https://doi.org/XXXXXX.XXXXXX>

Histogram of SSD Storage Embodied Factors



Today: Quantifying the carbon footprint of computing

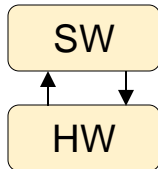
with insights, opening new research and sustainable development opportunities



Understanding the source of computing's emissions

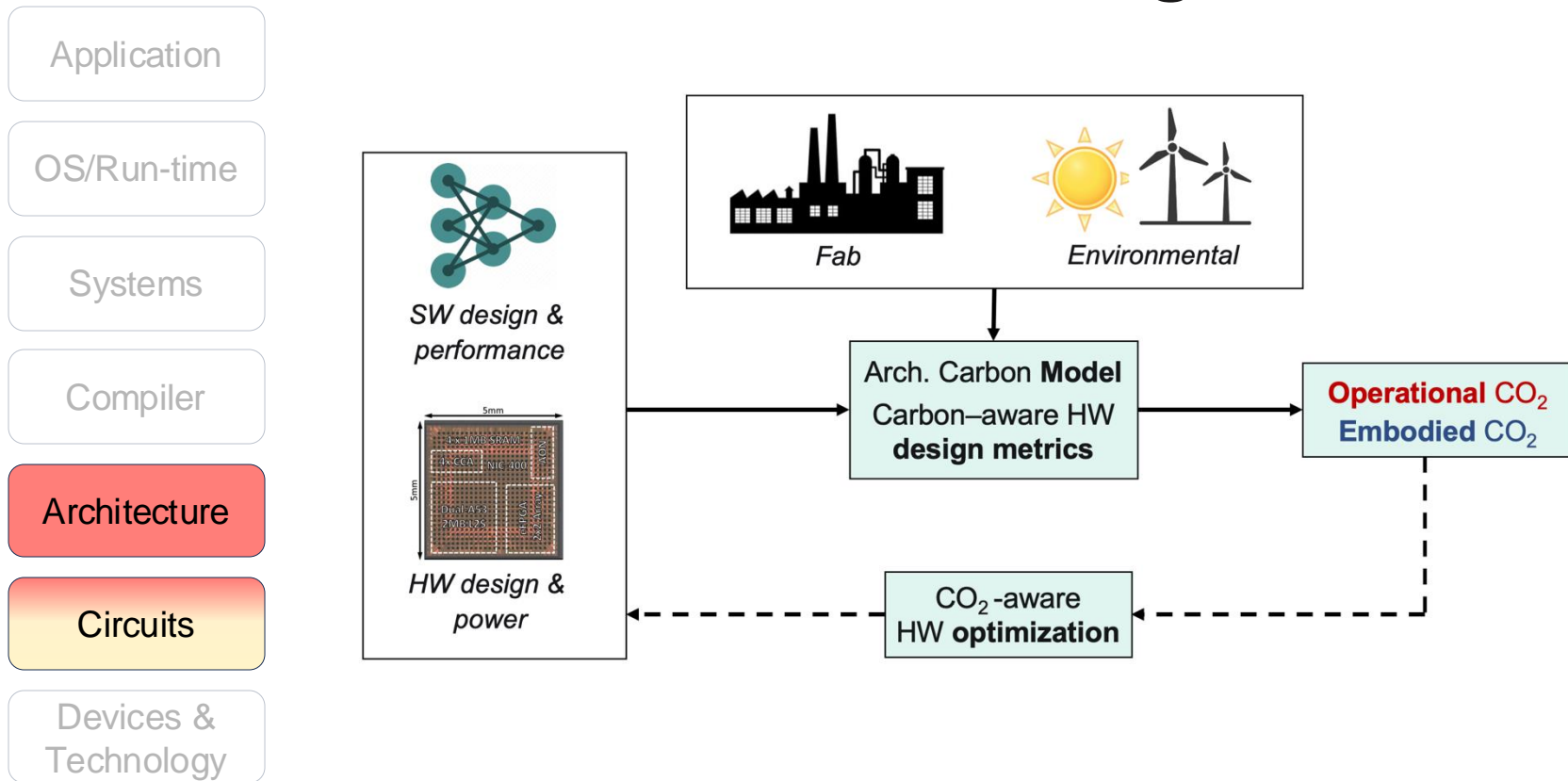


Deep dive: developing computer architectural models to estimate CO₂ emissions



Cross stack: Developing modeling methods across the computing stack

Need to go beyond architecture centric-view for cross-stack carbon accounting



Application

OS/Run-time

Systems

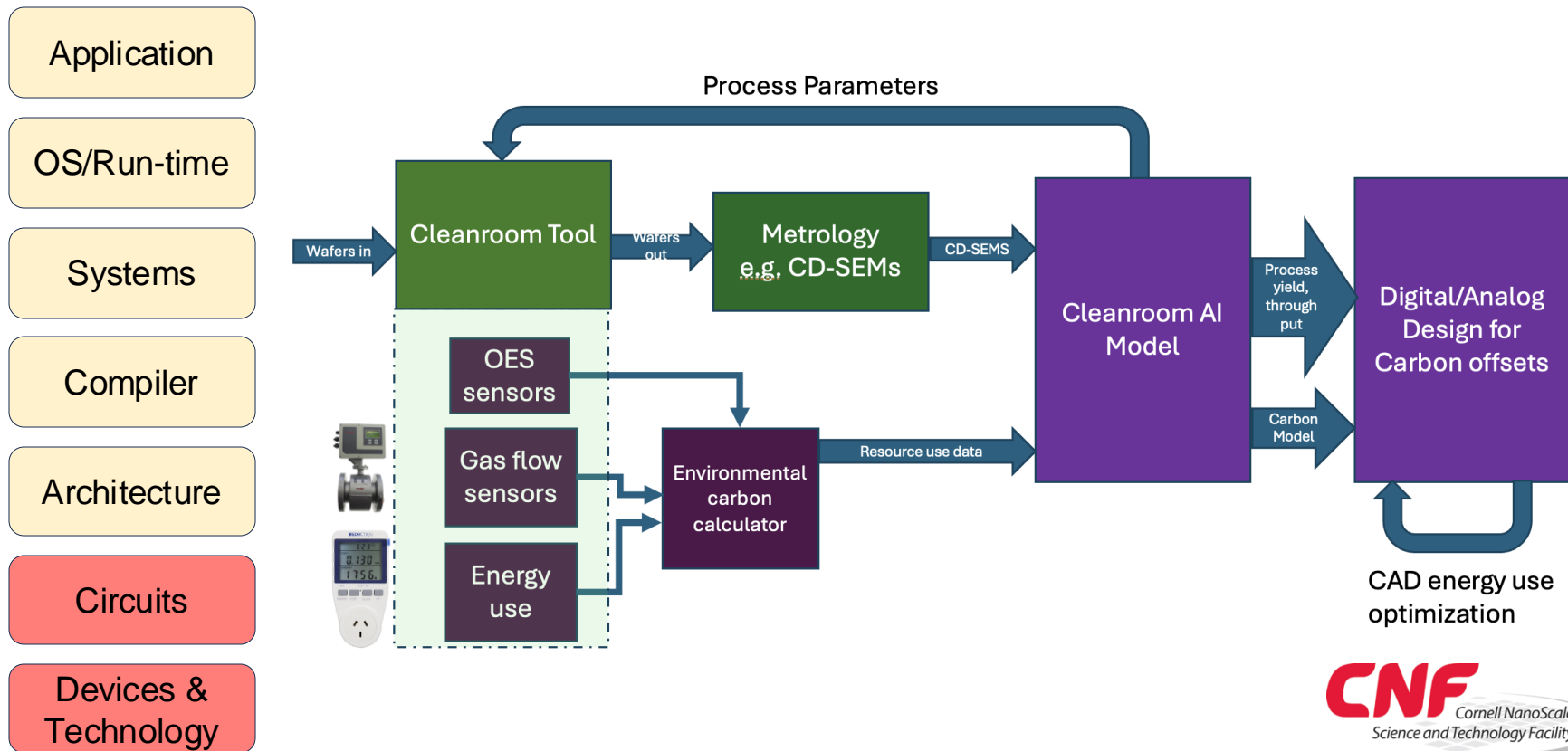
Compiler

Architecture

Circuits

Devices &
Technology

Instrumenting Cornell's NanoScale Facility



Rising emissions from AI Footprint

Application

OS/Run-time

Systems

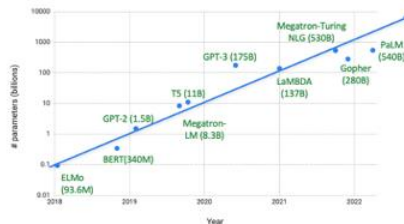
Compiler

Architecture

Circuits

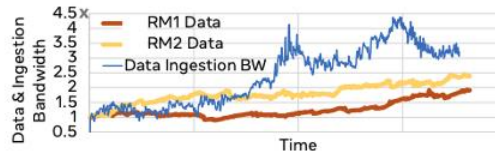
Devices &
Technology

Model Growth



- State-of-the-art DL models growing in capacity by **10x/year**

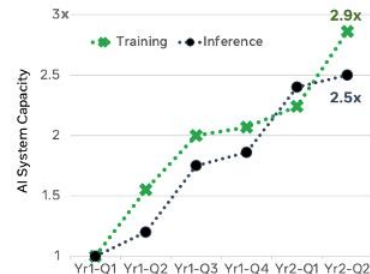
Data Growth



For recommendation models, in two years:

- Data stored and used **doubled**
- Storage bandwidth grew by **3.2x**

Systems Growth



In 1.5 years:

- Training capacity grew by **2.9x**
- Inference capacity grew by **2.5x**

Crucial to look at emissions across ML cycle

Application

OS/Run-time

Systems

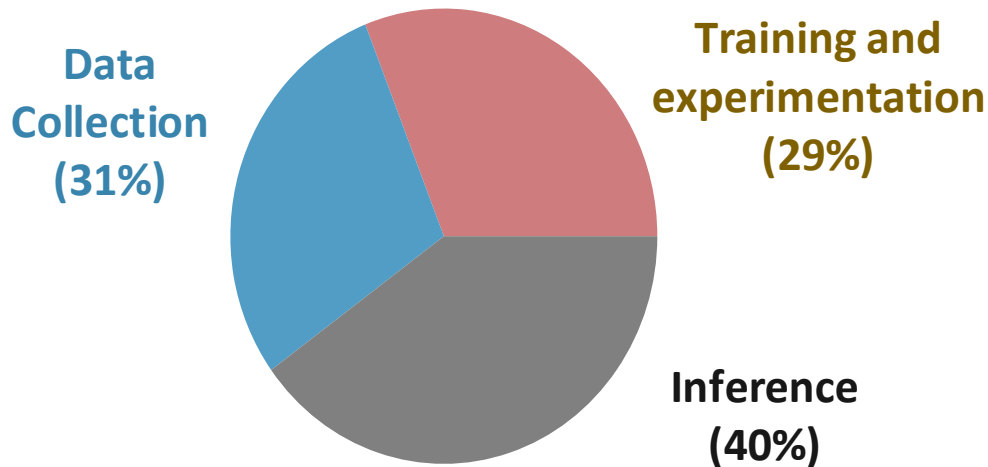
Compiler

Architecture

Circuits

Devices &
Technology

Machine learning life cycle



Crucial to look at emissions across ML cycle

Application

OS/Run-time

Systems

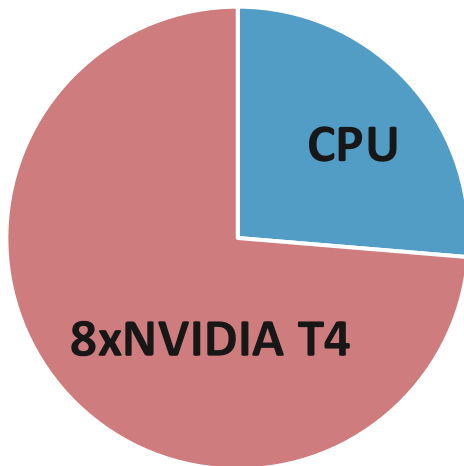
Compiler

Architecture

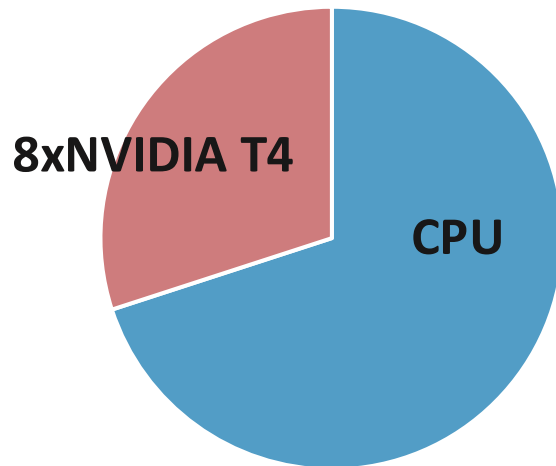
Circuits

Devices &
Technology

Inference Power



**Inference Server
Embodied Carbon**



Crucial to look at emissions across ML cycle

Application

OS/Run-time

Systems

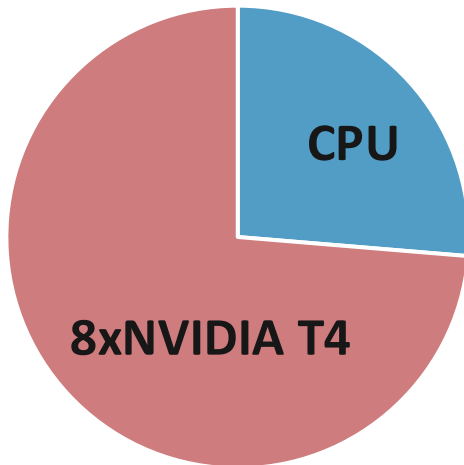
Compiler

Architecture

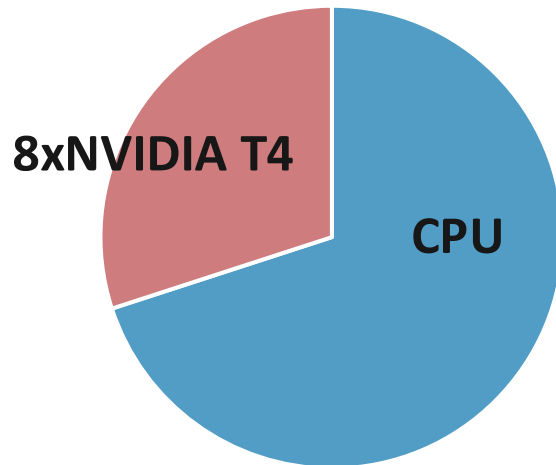
Circuits

Devices &
Technology

Inference Power



**Inference Server
Embodied Carbon**



Research questions:

- How do the breakdowns scale for different AI applications and hardware?
- How do we co-optimize AI systems for end-to-end carbon?

“Towards Carbon-efficient LLM Life Cycle” Yueying Li, Omer Graif, Udit Gupta (HotCarbon 2024)

Carbon Explorer: Carbon-Aware Datacenter Design

Application

OS/Run-time

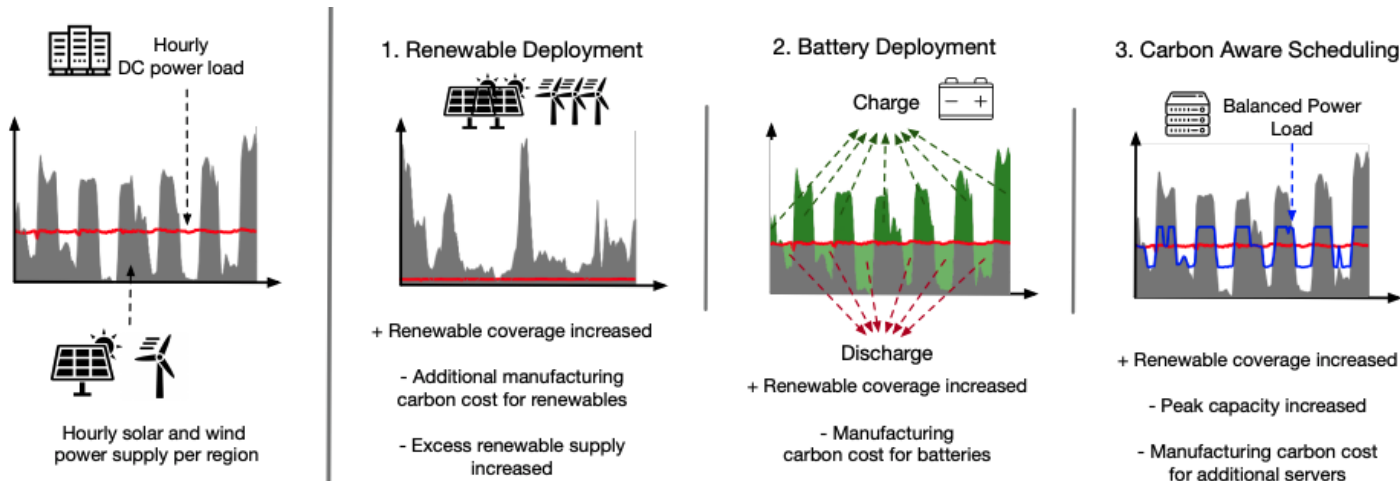
Systems

Compiler

Architecture

Circuits

Devices &
Technology



Optimizing both *operational* and *embodied* carbon requires balancing
(1) **renewable deployment**, (2) **battery deployment**, (3) **carbon-aware scheduling**

Attributing carbon footprint of cloud usage

Application

OS/Run-time

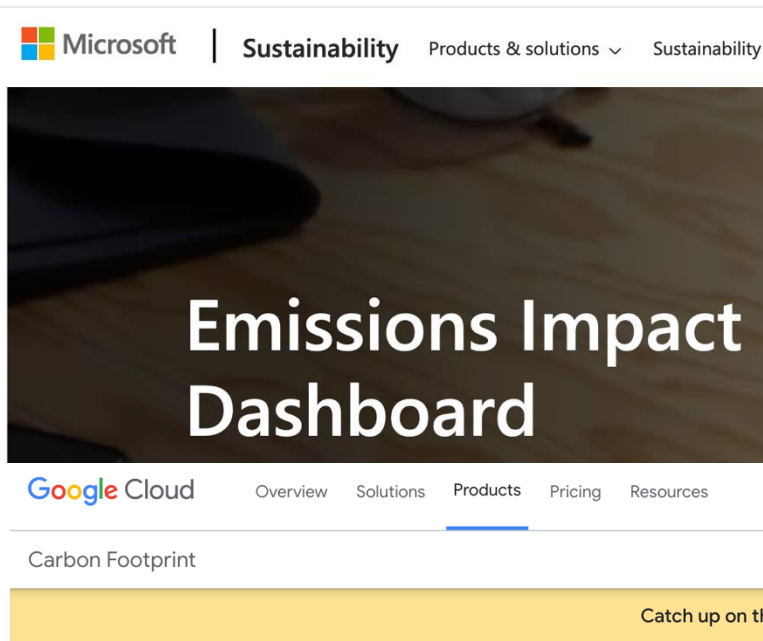
Systems

Compiler

Architecture

Circuits

Devices &
Technology



Carbon Footprint

Benefits



Carbon Footprint

Measure, report, and reduce your cloud carbon emissions.

Fairly attributing carbon

Application

OS/Run-time

Systems

Compiler

Architecture

Circuits

Devices &
Technology

Carbon accounting in the Cloud:
a methodology for allocating emissions across
data center users

Ian Schneider*, Taylor Mattia*†

June 2024

1 Introduction

Google has undertaken considerable efforts to reduce electricity consumption and the associated greenhouse gas (GHG) emissions from its electricity use. By 2022, Google delivered approximately three times as much computing power with the same amount of electrical power as it did five years prior [1].¹ Google uses 5.5 times less overhead energy for every unit of information-technology (IT) equipment energy, compared to the industry average [1]. Even with these dramatic improvements in efficiency, Google consumed 22 TWh of electricity in 2022, with the majority of its electricity consumption coming from data center operations [1].

Fairly attributing carbon

Application

OS/Run-time

Systems

Compiler

Architecture

Circuits

Devices &
Technology

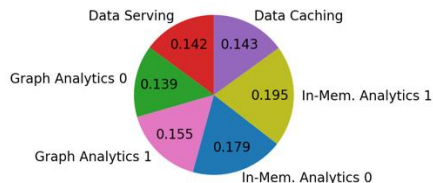
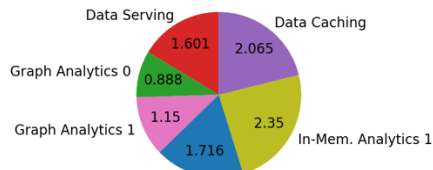
Carbon accounting in the Cloud:
a methodology for allocating emissions across
data center users

Ian Schneider*, Taylor Mattia*†

June 2024

1 Introduction

Google has undertaken considerable efforts to reduce electricity consumption and the associated greenhouse gas (GHG) emissions from its electricity use. By 2022, Google delivered approximately three times as much computing power with the same amount of electrical power as it did five years prior [1].¹ Google uses 5.5 times less overhead energy for every unit of information-technology (IT) equipment energy, compared to the industry average [1]. Even with these dramatic improvements in efficiency, Google consumed 22 TWh of electricity in 2022, with the majority of its electricity consumption coming from data center operations [1].



Open research questions:

- How do we fairly attribute operational and embodied carbon to individual cloud services?
- How do we consider varying demand in data centers in attributing carbon responsibility?
- How do we scale attribution mechanisms to cloud-scale?

Key takeaways

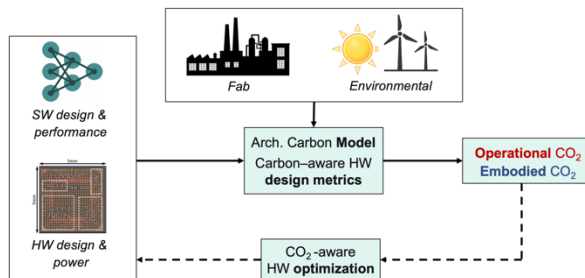
Quantifying the carbon footprint of AI and computing: Past, Present, and Future

Pillars of Advancing Sustainable Computing



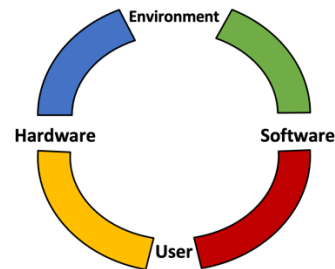
Economics & Policy
Education
Accounting and Reporting

Must analyze emissions across life cycles



Across hardware
manufacturing to
operational use

Dire need for cross-stack carbon optimization

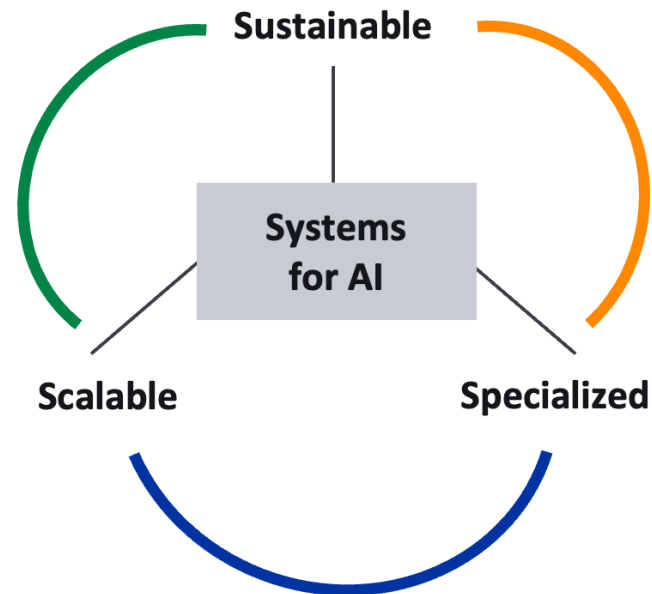


From applications to devices
Across eco-HW/SW design
loop

S⁴AI: Specialized, Scalable, Sustainable Systems for AI



Collaborators



Key takeaways

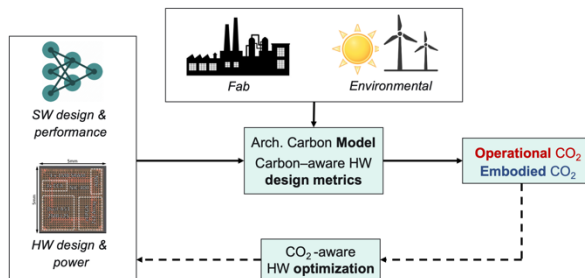
Quantifying the carbon footprint of AI and computing: Past, Present, and Future

Pillars of Advancing Sustainable Computing



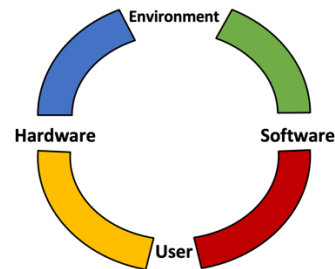
Economics & Policy
Education
Accounting and Reporting

Must analyze emissions across life cycles



Across hardware
manufacturing to
operational use

Dire need for cross-stack carbon optimization

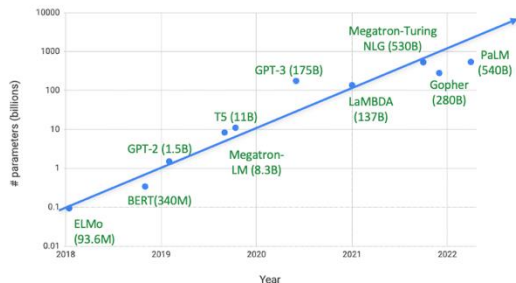


From applications to devices
Across eco-HW/SW design
loop

AI Growth Driving Datacenter Energy Consumption

AI Growth Driving Datacenter Energy Consumption

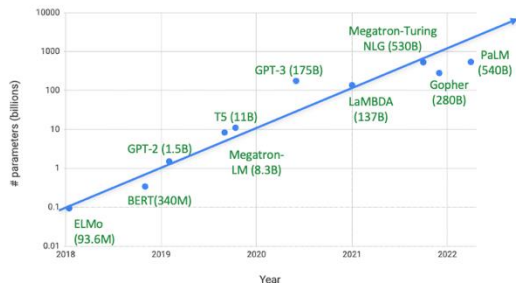
Model Growth



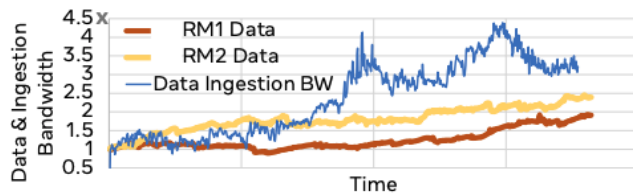
- State-of-the-art DL models growing in capacity by **10x/year**

AI Growth Driving Datacenter Energy Consumption

Model Growth



Data Growth



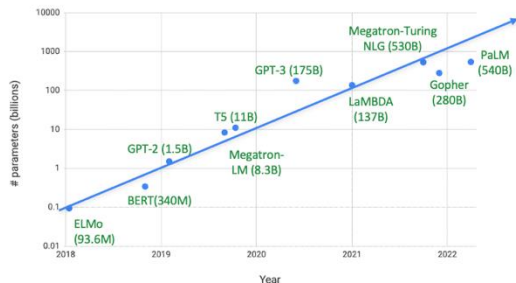
- State-of-the-art DL models growing in capacity by **10x/year**

For recommendation models, in two years:

- Data stored and used **doubled**
- Storage bandwidth grew by **3.2x**

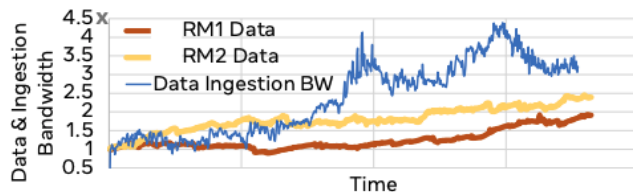
AI Growth Driving Datacenter Energy Consumption

Model Growth



- State-of-the-art DL models growing in capacity by **10x/year**

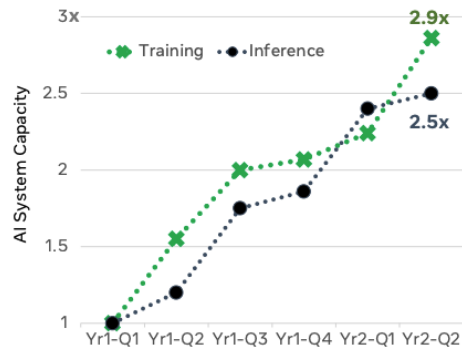
Data Growth



For recommendation models, in two years:

- Data stored and used **doubled**
- Storage bandwidth grew by **3.2x**

Systems Growth

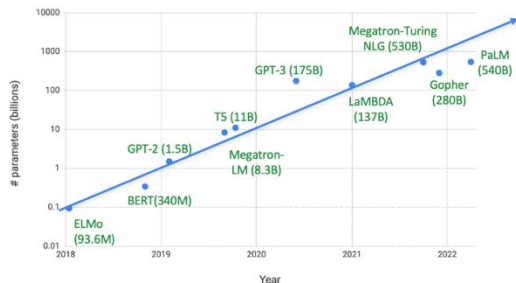


In 1.5 years:

- Training capacity grew by **2.9x**
- Inference capacity grew by **2.5x**

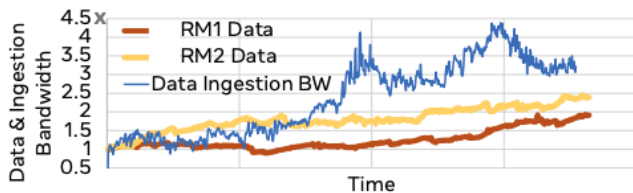
AI Growth Driving Datacenter Energy Consumption

Model Growth



- State-of-the-art DL models growing in capacity by **10x/year**

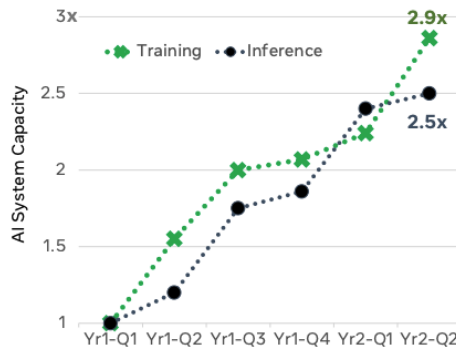
Data Growth



For recommendation models, in two years:

- Data stored and used **doubled**
- Storage bandwidth grew by **3.2x**

Systems Growth



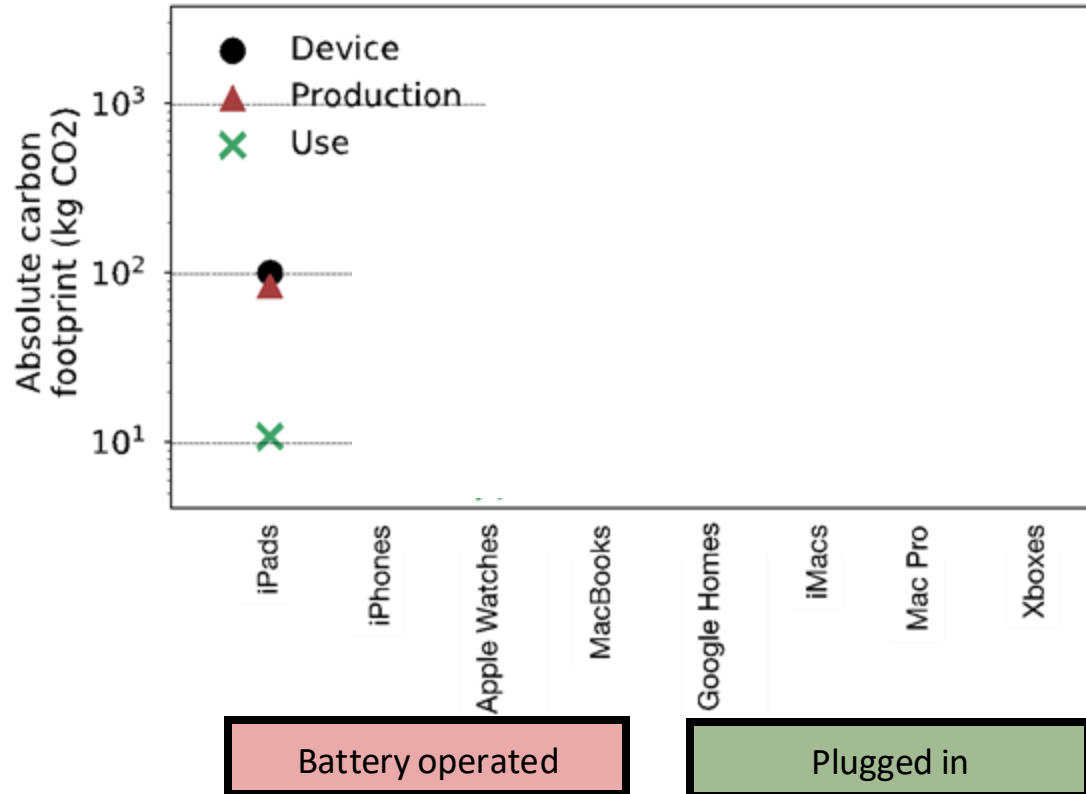
In 1.5 years:

- Training capacity grew by **2.9x**
- Inference capacity grew by **2.5x**

Recent advances in LLMs further exacerbating model, data, and systems trends!

Carbon footprint characteristics vary across devices

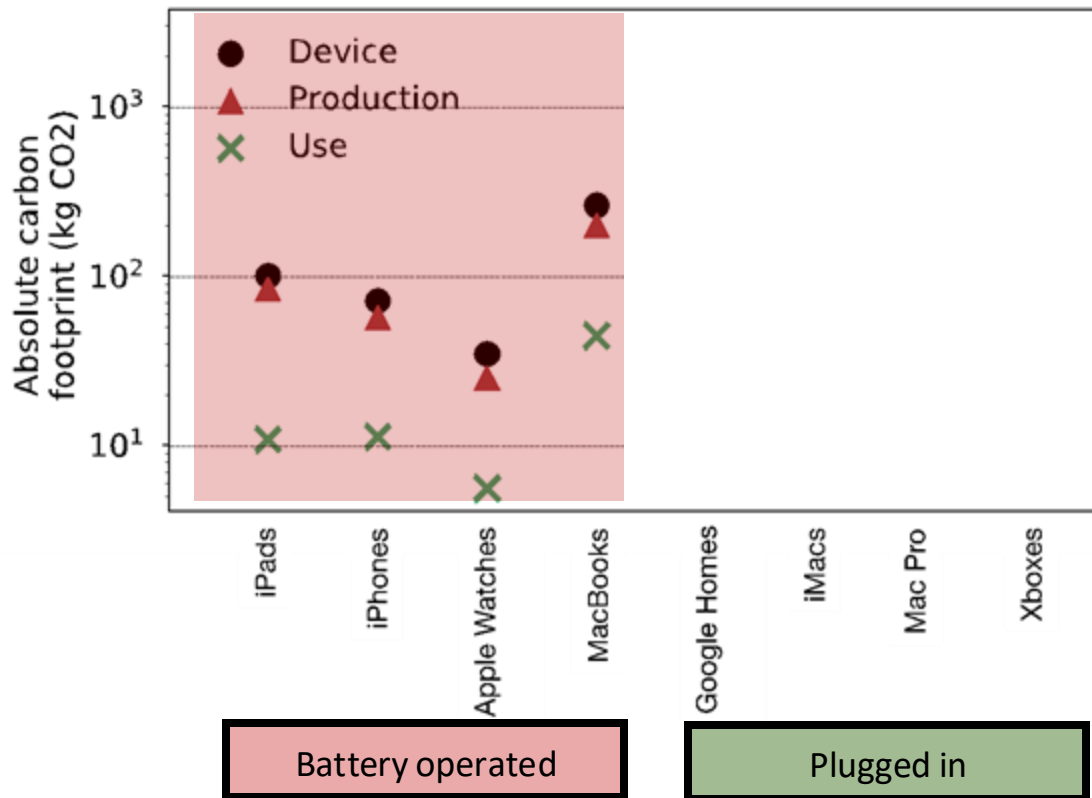
Data from public industry validated sustainability reports and life cycle analyses



Carbon footprint characteristics vary across devices

Data from public industry validated sustainability reports and life cycle analyses

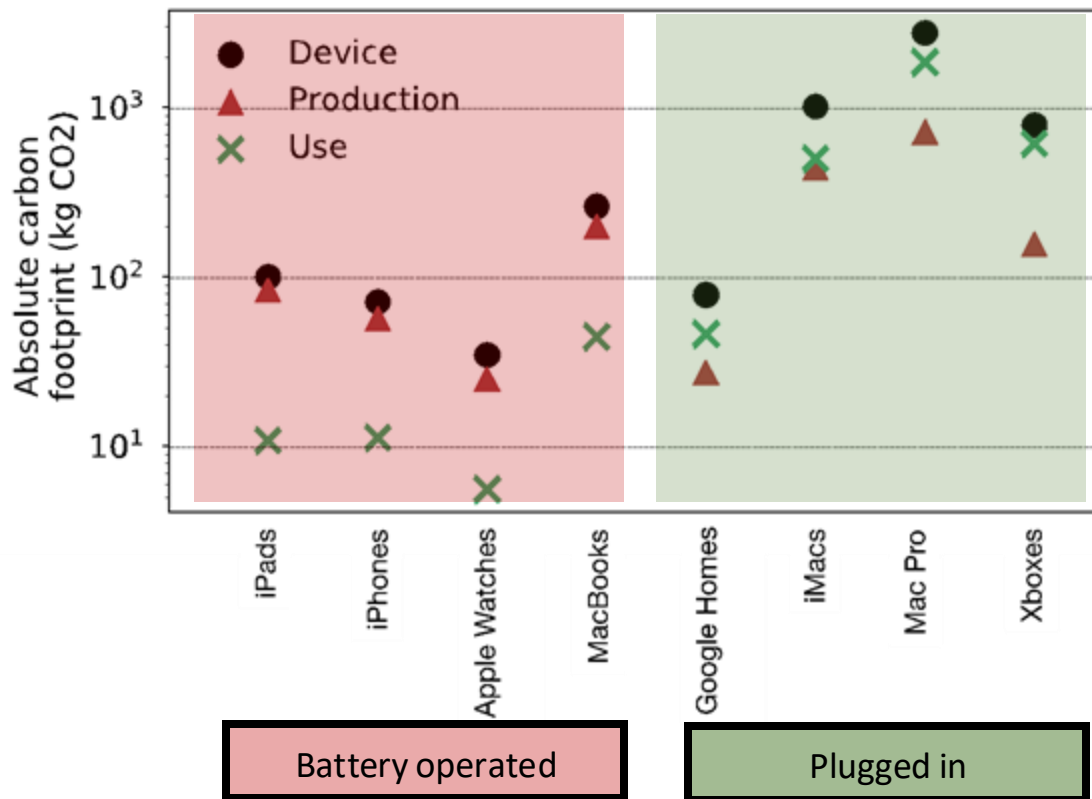
Roughly **75%** life cycle emissions for **battery operated devices** comes from hardware manufacturing.



Carbon footprint characteristics vary across devices

Data from public industry validated sustainability reports and life cycle analyses

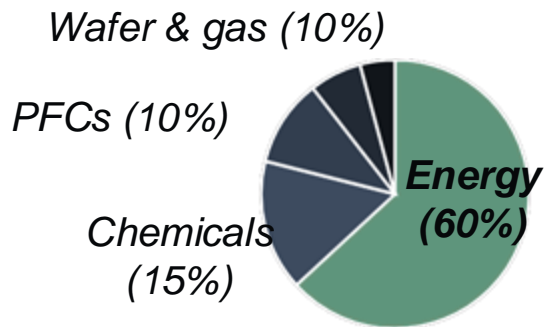
Roughly **75%** life cycle emissions for **battery operated devices** comes from hardware manufacturing.



Emissions for **always-connected devices** come mainly from **energy consumption**



*Semiconductor
fab*



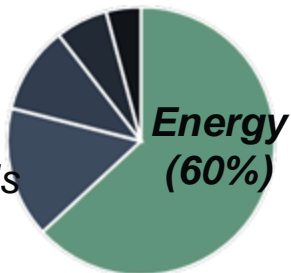


*Semiconductor
fab*

Wafer & gas (10%)

PFCs (10%)

*Chemicals
(15%)*



***100% Renewable** powered
semiconductor fab*

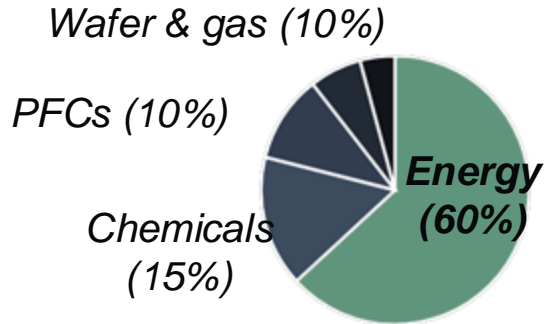


*Reduces manufacturing footprint by **2.5x***

“Green” powered fabs are not enough



Semiconductor
fab



TSMC plans for
25% renewable by 2025 and
100% renewable by 2050.



100% Renewable powered
semiconductor fab



Reduces manufacturing footprint by **2.5x**

Architectural Carbon Model

Model	Hardware/software input
-------	-------------------------

Architectural Carbon Model

Model	Hardware/software input
-------	-------------------------

$$\text{Carbon} = OP_{CF} + \frac{\text{Runtime}}{\text{Lifetime}} Emb_{CF}$$

Performance/power/energy and
lifetime of hardware

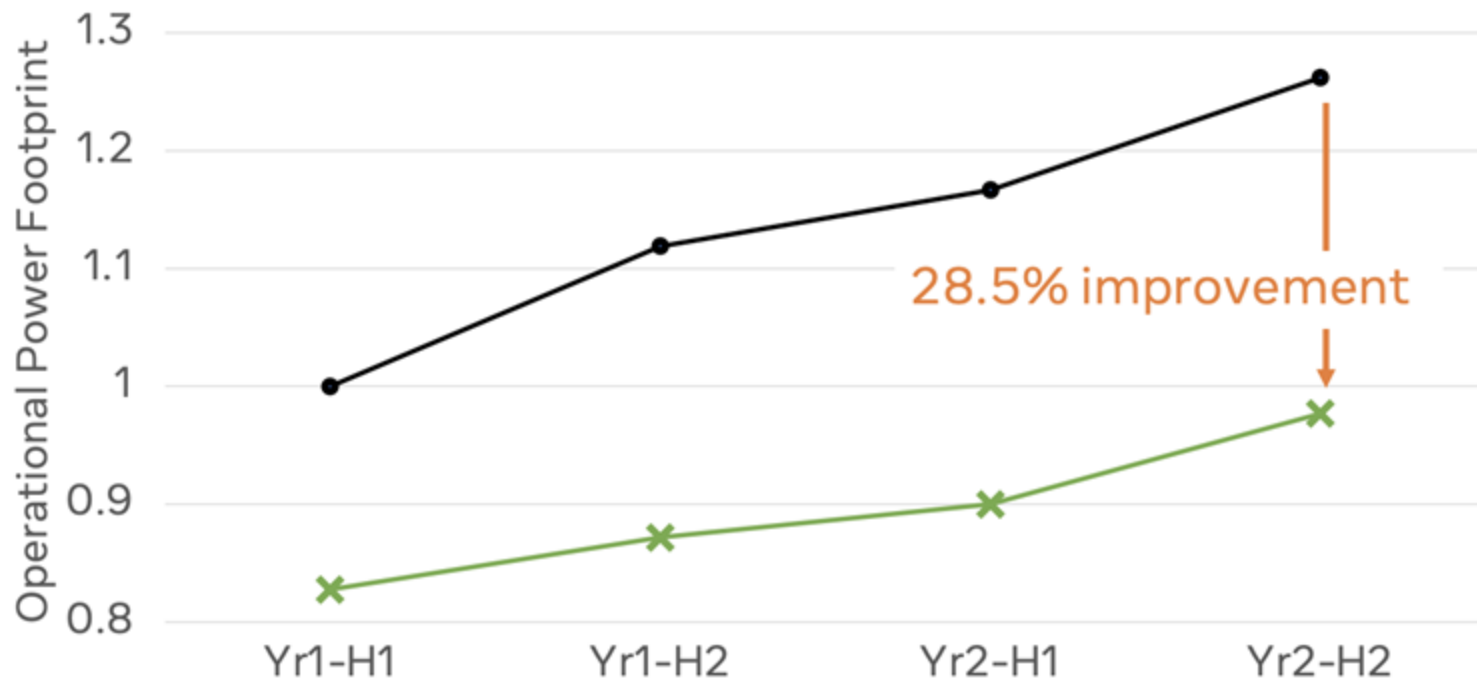
Architectural Carbon Model

Model	Hardware/software input
$Carbon = OP_{CF} + \frac{Runtime}{Lifetime} Emb_{CF}$	Performance/power/energy and lifetime of hardware
$OP_{CF} = CI_{use} \times Energy$	Energy efficiency and environment (carbon intensity)

Embodied carbon of application processors (SoC's)

$$Emb_{soc} = Area \times \frac{(CI_{fab} \times Fab_{energy}) + Fab_{chemicals} + Fab_{materials}}{Yield}$$

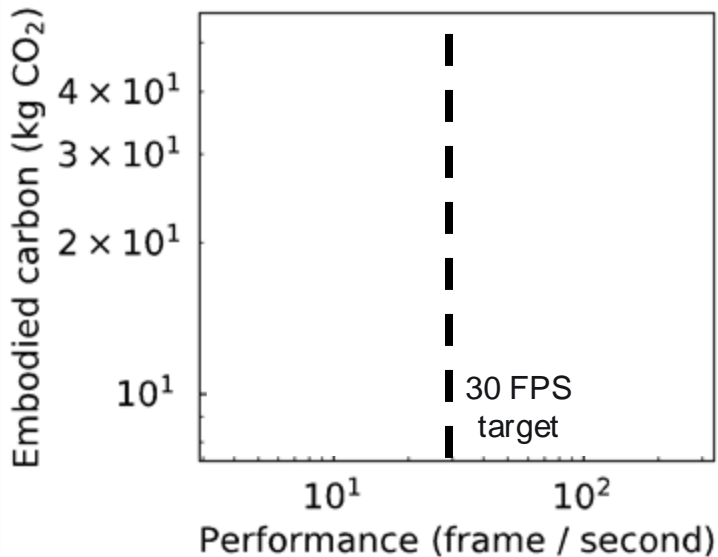
Jevon's paradox of AI at-scale



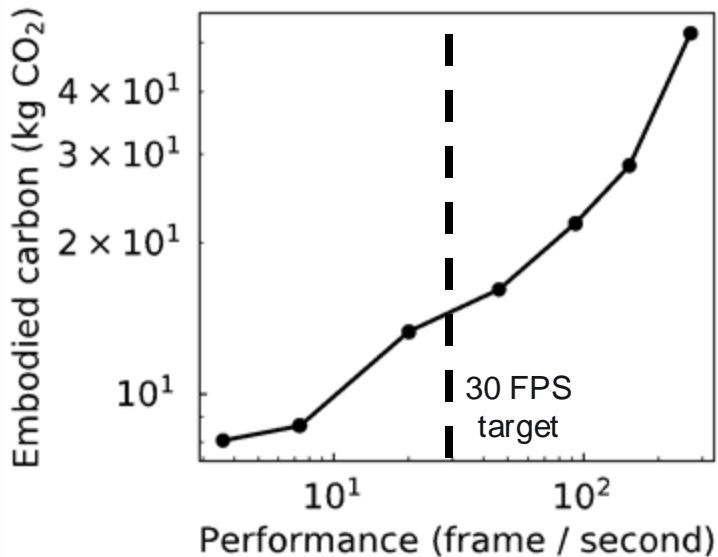
Reduce: Designing leaner hardware systems



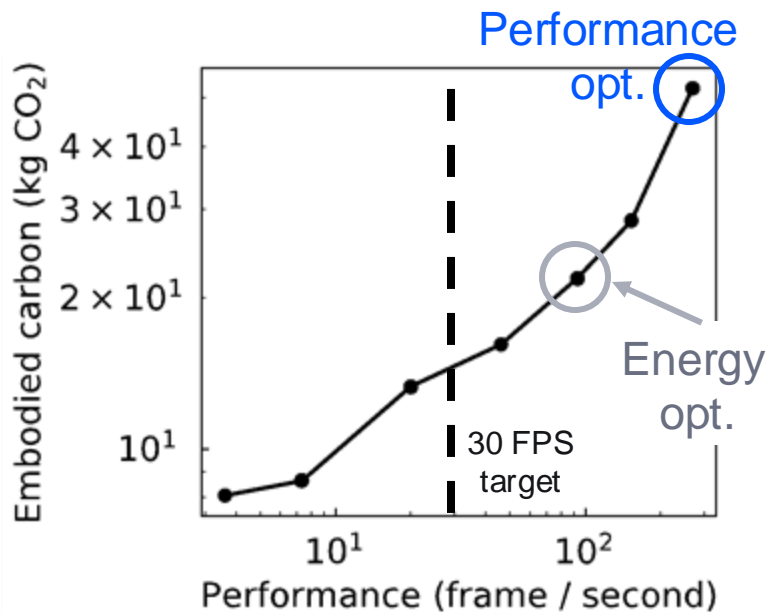
Reduce: Designing leaner hardware systems



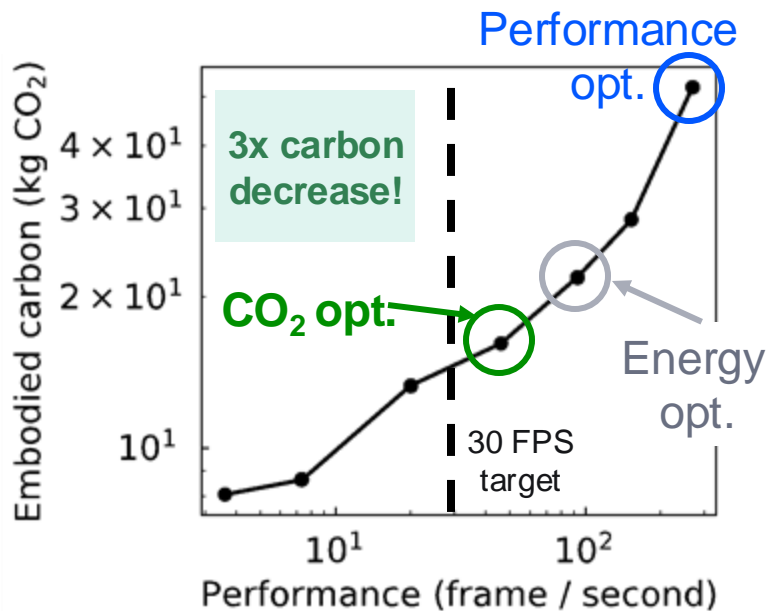
Reduce: Designing leaner hardware systems



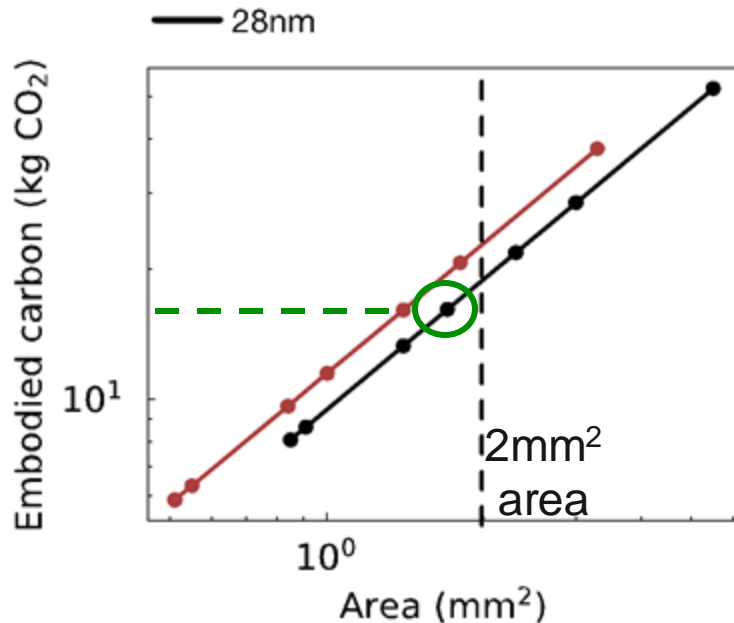
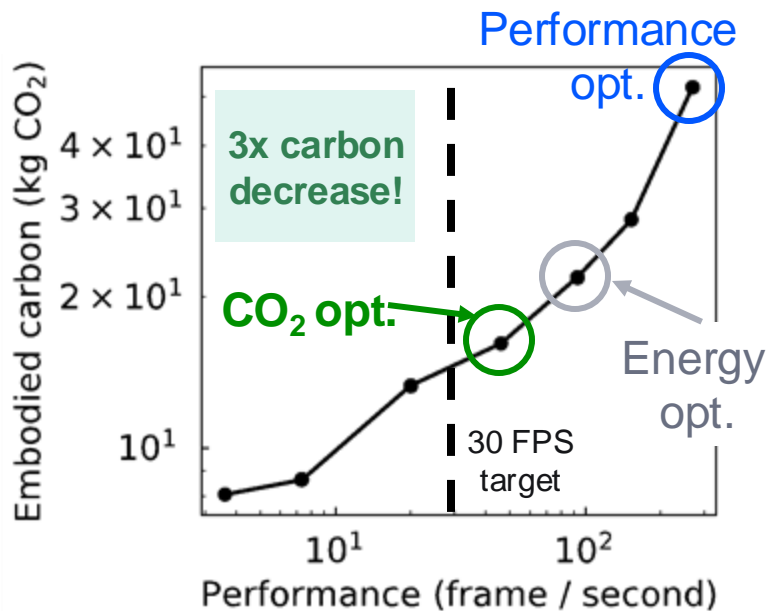
Reduce: Designing leaner hardware systems



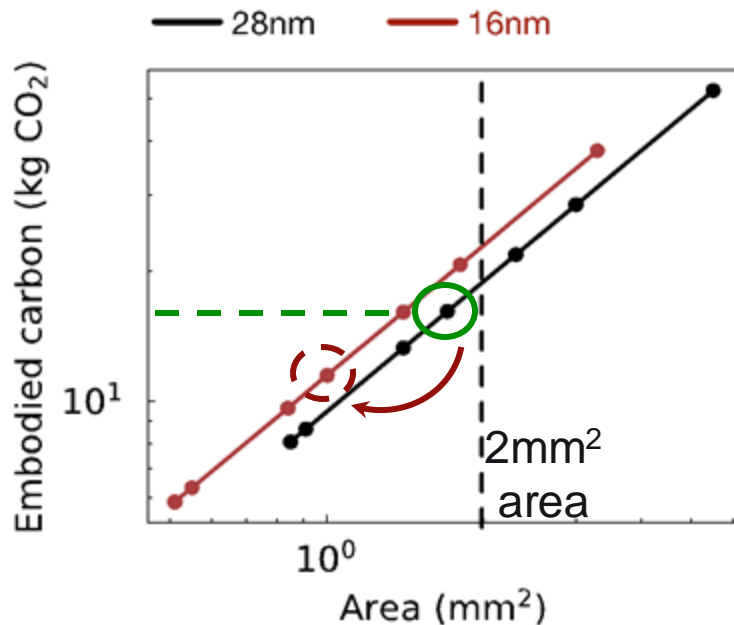
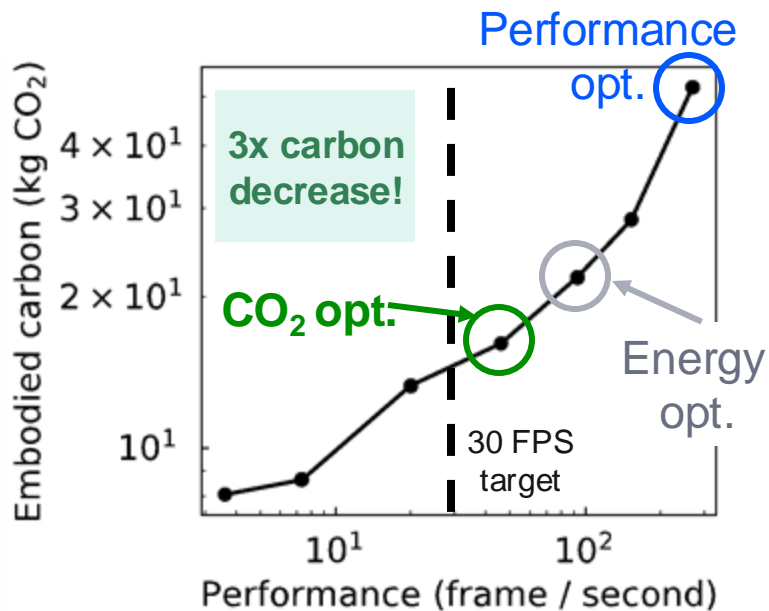
Reduce: Designing leaner hardware systems



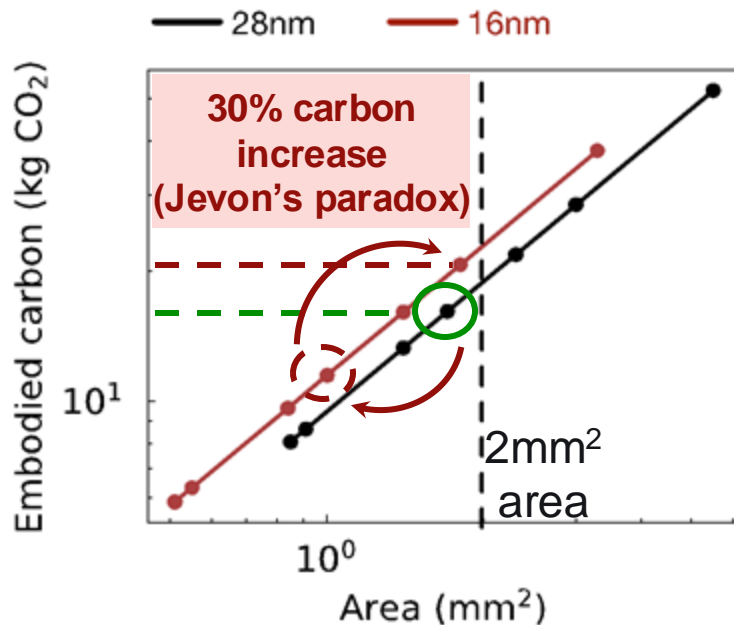
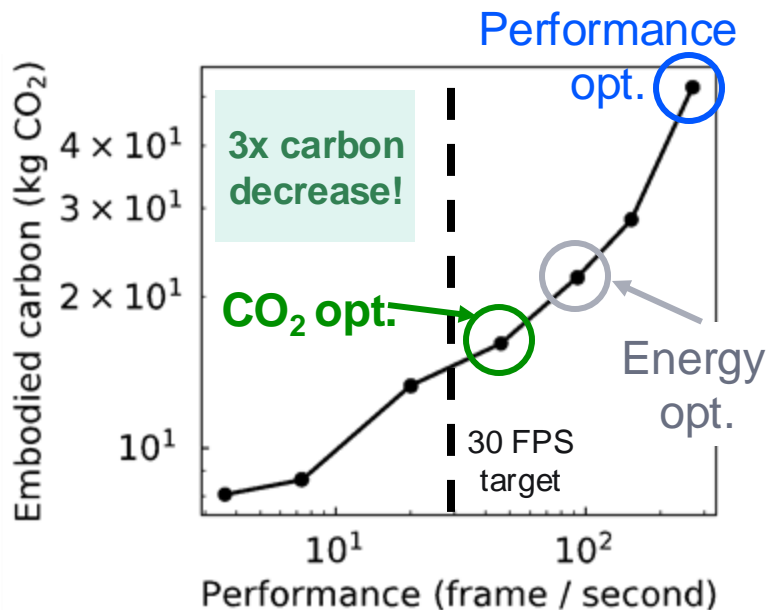
Reduce: Designing leaner hardware systems



Reduce: Designing leaner hardware systems



Reduce: Designing leaner hardware systems



DRAM Memory Embodied Carbon Emissions

