

Historical Analysis of Energy Consumption in Large Scale Computer Systems

From the 1990s to the 2020s

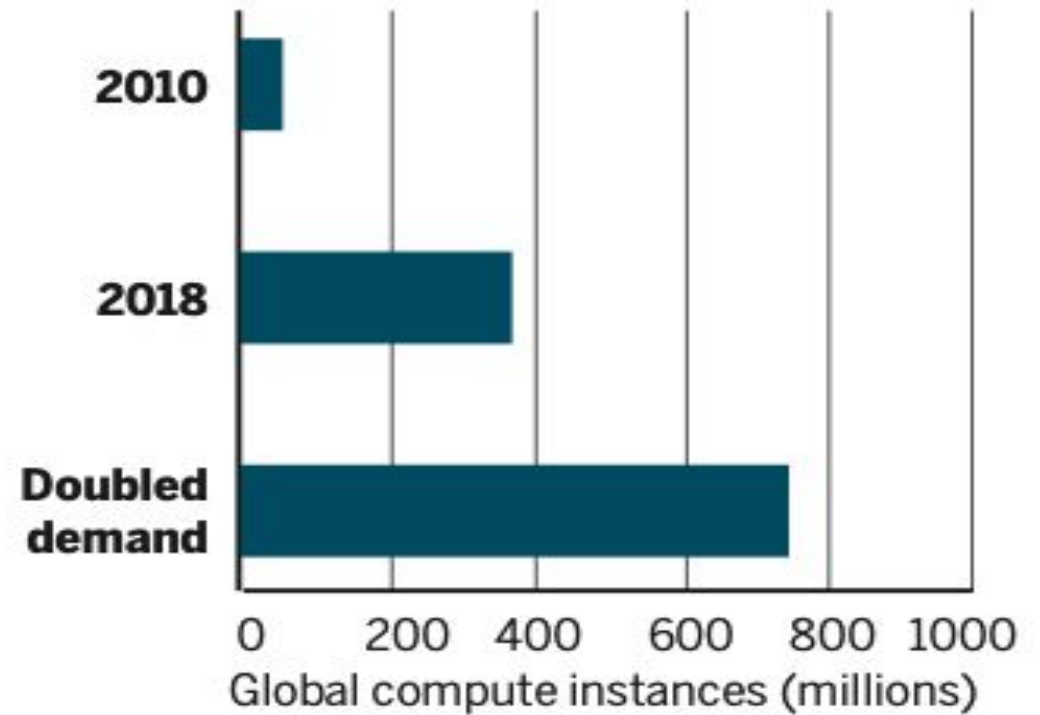
Dide Poyraz, Dante Niewenhuis, Alexandru Iosup



Data Centers are Vital to Digital Society

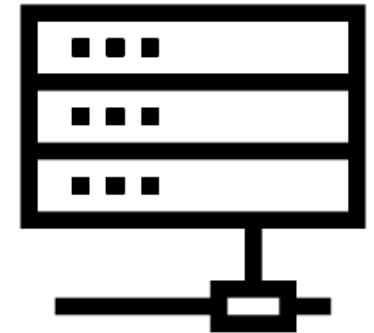
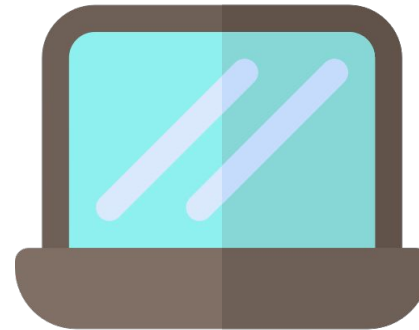


**Global data center
compute instances**



Analyzing Energy Consumption is Hard

- Lack of tools / guidelines
- Experimentation is **costly and time-consuming**
- Existing predictions based on the peak power usage

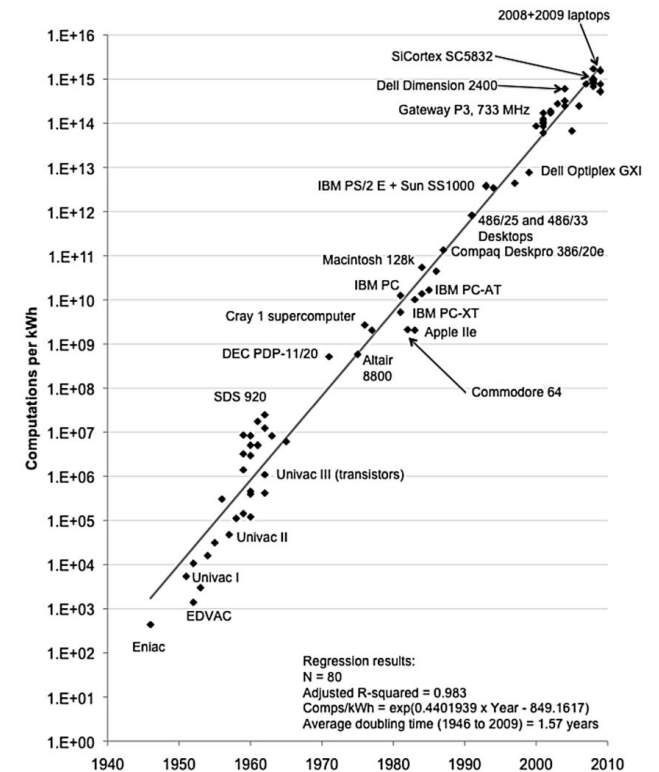


Koomey's Law

Energy efficiency of computers doubles around every 18 months

- **Simplified** Assumptions on models
- Static Conditions based on peak power

Num. of computations at maximum load
Energy Consumed per hour

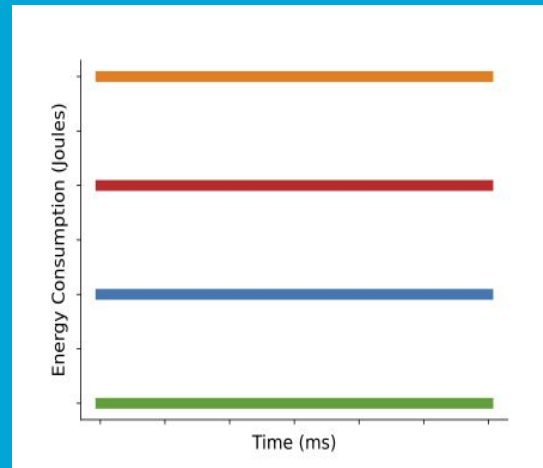


Dynamic Models - A New Approach



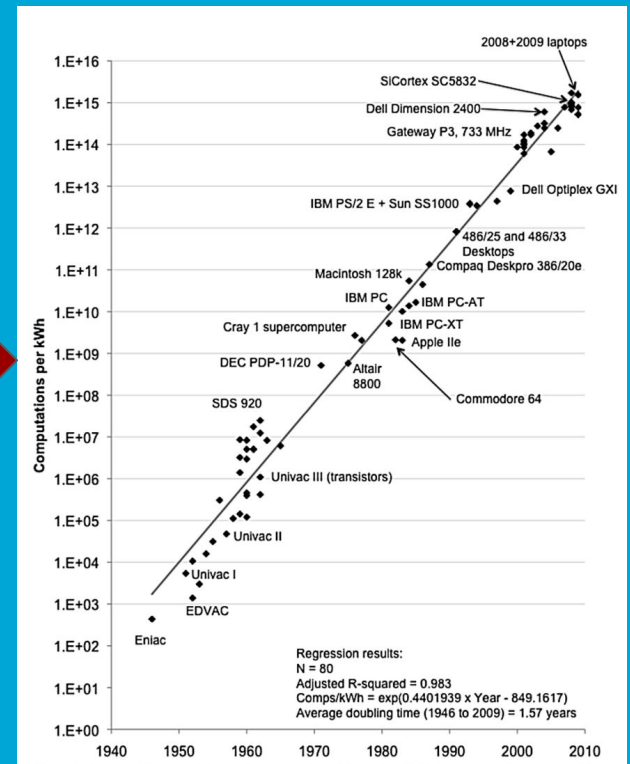
Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation



Static Workloads

Koomey's Law

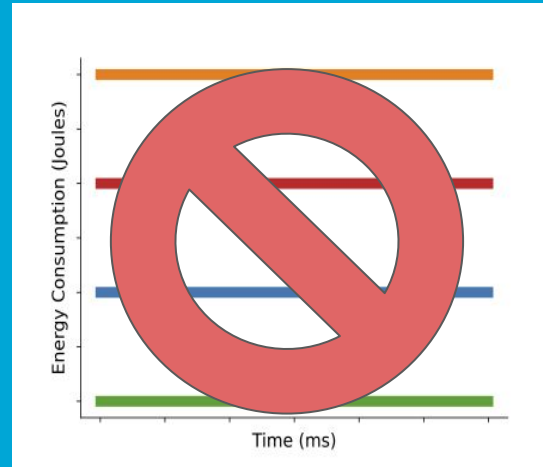


Dynamic Models - A New Approach



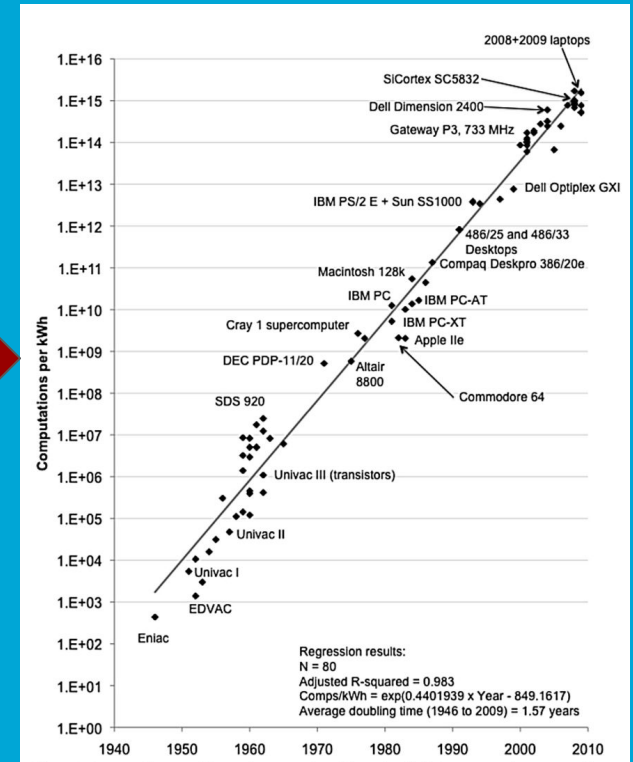
Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation



Static Workloads

Koomey's Law



Dynamic Models - A New Approach

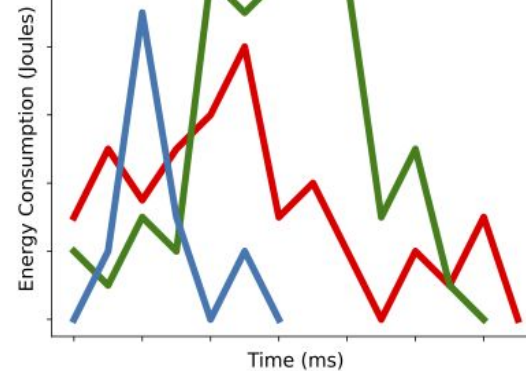


Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation

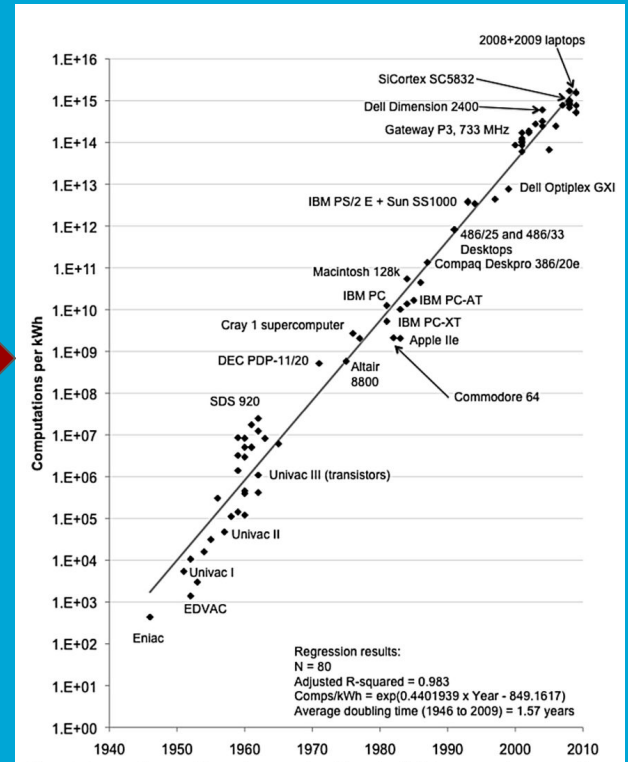
Different scheduling techniques and failures

+



Dynamic-Real World Workloads

Koomey's Law



Dynamic Models - A New Approach

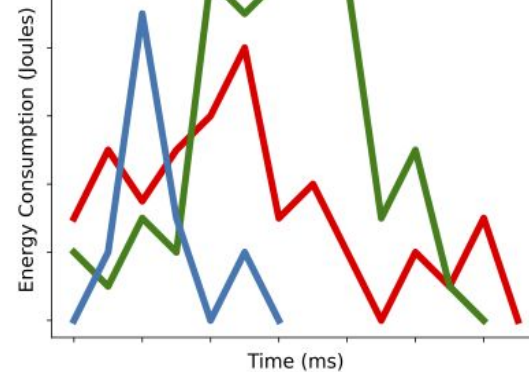


Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation

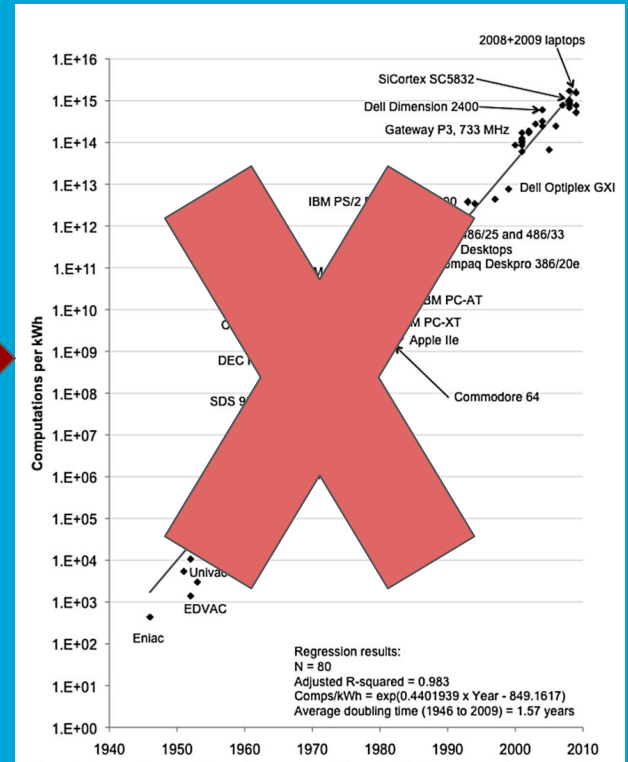
Different scheduling techniques and failures

+



Dynamic-Real World Workloads

Koomey's Law



Dynamic Models - A New Approach

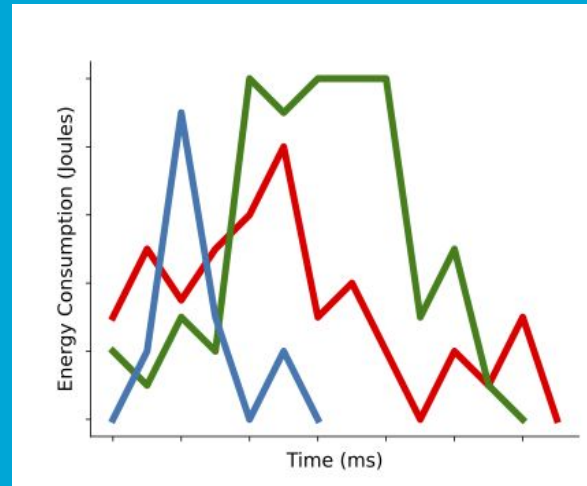


Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation

Different scheduling techniques and failures

+



Dynamic-Real World Workloads

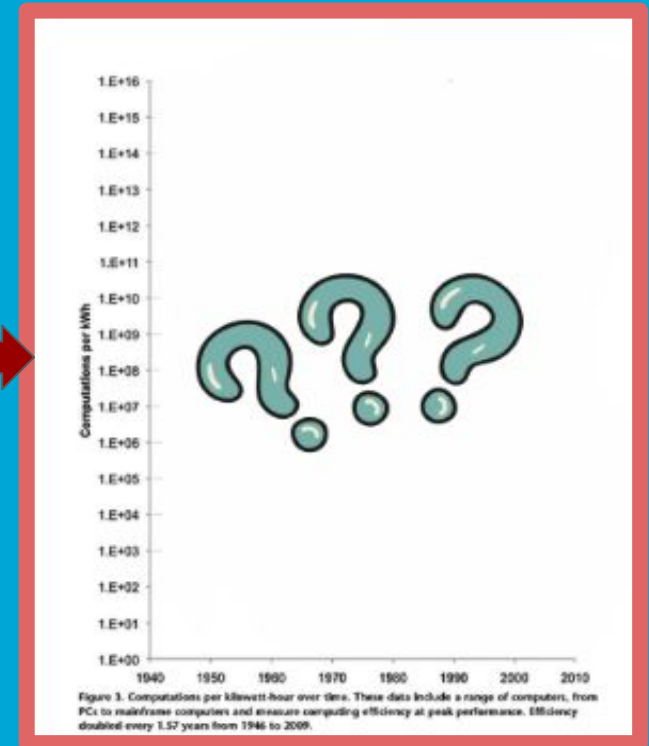
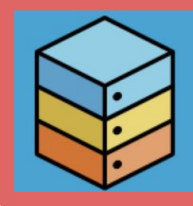


Figure 3. Computations per kilowatt-hour over time. These data include a range of computers, from PCs to mainframe computers and measure computing efficiency at peak performance. Efficiency doubled every 1.57 years from 1946 to 2009.

Dynamic Models - A New Approach

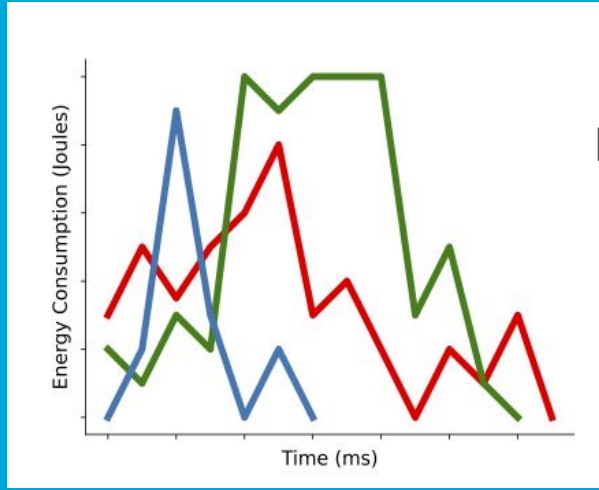


Koomey's static models are insufficient to understand real-world energy usage

Simulating Infrastructures from 1990s to 2020s through Simulation

Different scheduling techniques and failures

+



Dynamic-Real World Workloads

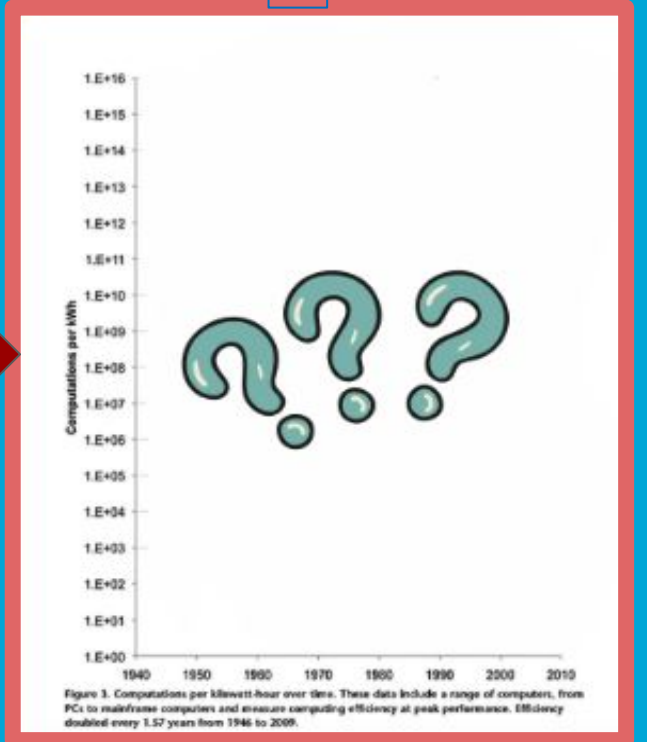


Figure 3. Computations per kilowatt-hour over time. These data include a range of computers, from PCs to mainframe computers and measure computing efficiency at peak performance. Efficiency doubled every 1.57 years from 1946 to 2009.

Research Questions

RQ1 How to select relevant computer systems and workloads to build a taxonomy of exemplary infrastructures per decade from 1950s to 2023?

RQ2 How to model the operation of data centers from different time periods to analyze the evolution of energy consumption by such infrastructure?

RQ3 How did energy consumption evolve from the 1950s to 2023 for the execution of realistic workloads in [contemporary] data centers?

Research Questions

RQ1 How to select relevant computer systems and workloads to build a taxonomy of exemplary infrastructures per decade from 1950s to 2020s.

How to Select Systems?

RQ2 How to model the operation of data centers from different time periods to analyze the evolution of energy consumption by such infrastructure?

RQ3 How did energy consumption evolve from the 1950s to 2023 for the execution of realistic workloads in [contemporary] data centers?

Research Questions

RQ1 How to select relevant computer systems and workloads to build a taxonomy of exemplary infrastructures per decade from 1950s to 2020s.

How to Select Systems?

RQ2 How to model such systems in different periods to analyze the evolution of energy consumption by such infrastructure?

How to Model Chosen Systems?

RQ3 How did energy consumption evolve from the 1950s to 2023 for the execution of realistic workloads in [contemporary] data centers?

Research Questions

RQ1 How to select relevant computer systems and workloads to build a taxonomy of exemplary infrastructures per decade from 1950s to 2020s.

How to Select Systems?

RQ2 How to model such systems over long periods to analyze the evolution of energy consumption by such infrastructure?

How to Model Chosen Systems?

RQ3 How did energy consumption change over time during the execution of realistic workloads in [contemporary] data centers?

How to Analyse the Findings?

CELEBRATING 25 YEARS OF TOP500

51 LISTS | 10421 SYSTEMS

125 VENDORS | 2853 SITES

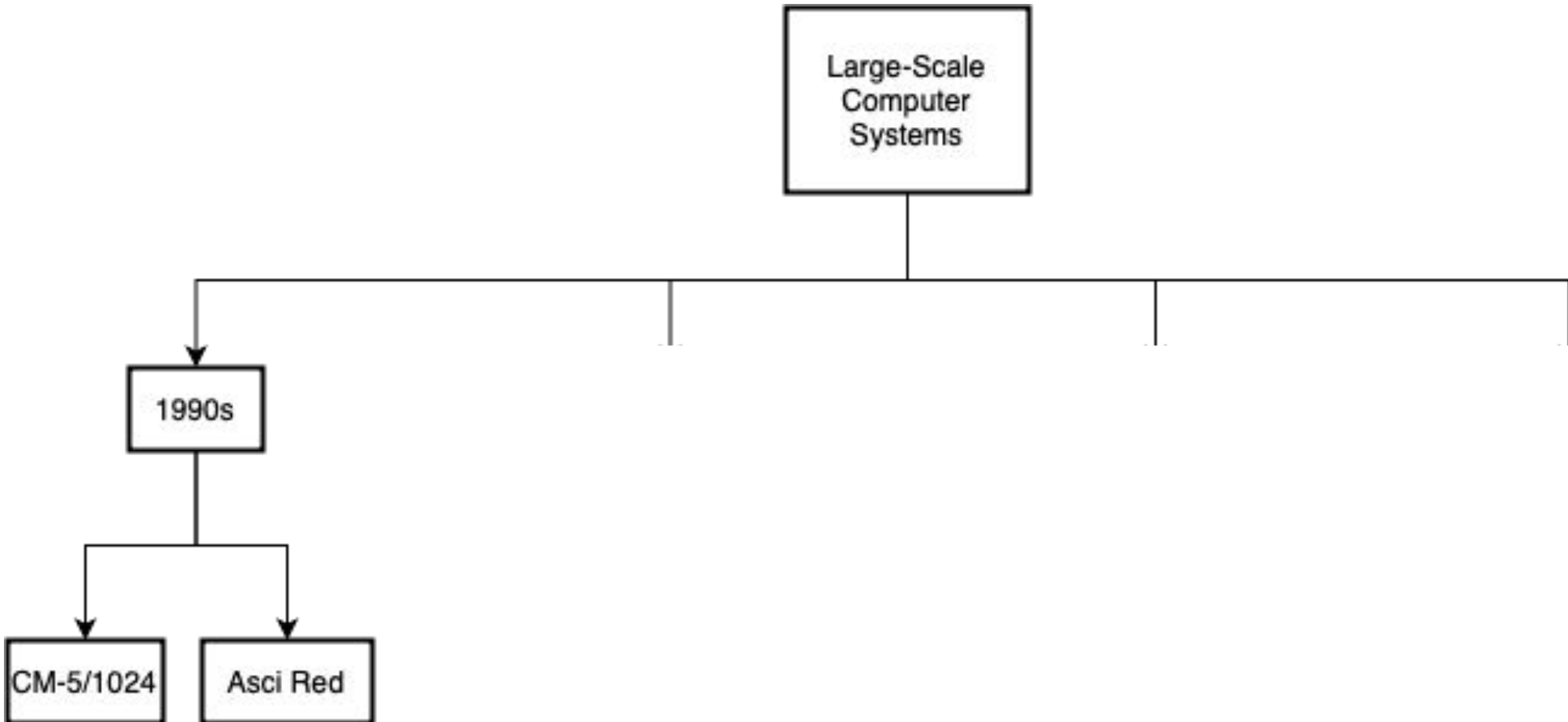
59 COUNTRIES

FROM 425 MFLOP/S TO 122 PFLOP/S

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	8,699,904	1,206.00	1,714.81	22,786
	Aurora - HPE Cray EX - Intel Exascale Compute Blade, Xeon CPU Max 9470 52C 2.4GHz, Intel Data Center GPU Max, Slingshot-11, Intel DOE/SC/Argonne National Laboratory United States	9,264,128	1,012.00	1,980.01	38,698
	Eagle - Microsoft NDv5, Xeon Platinum 8480C 48C 2GHz, NVIDIA H100, NVIDIA Infiniband NDR, Microsoft Azure Microsoft Azure United States	2,073,600	561.20	846.84	
4	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899

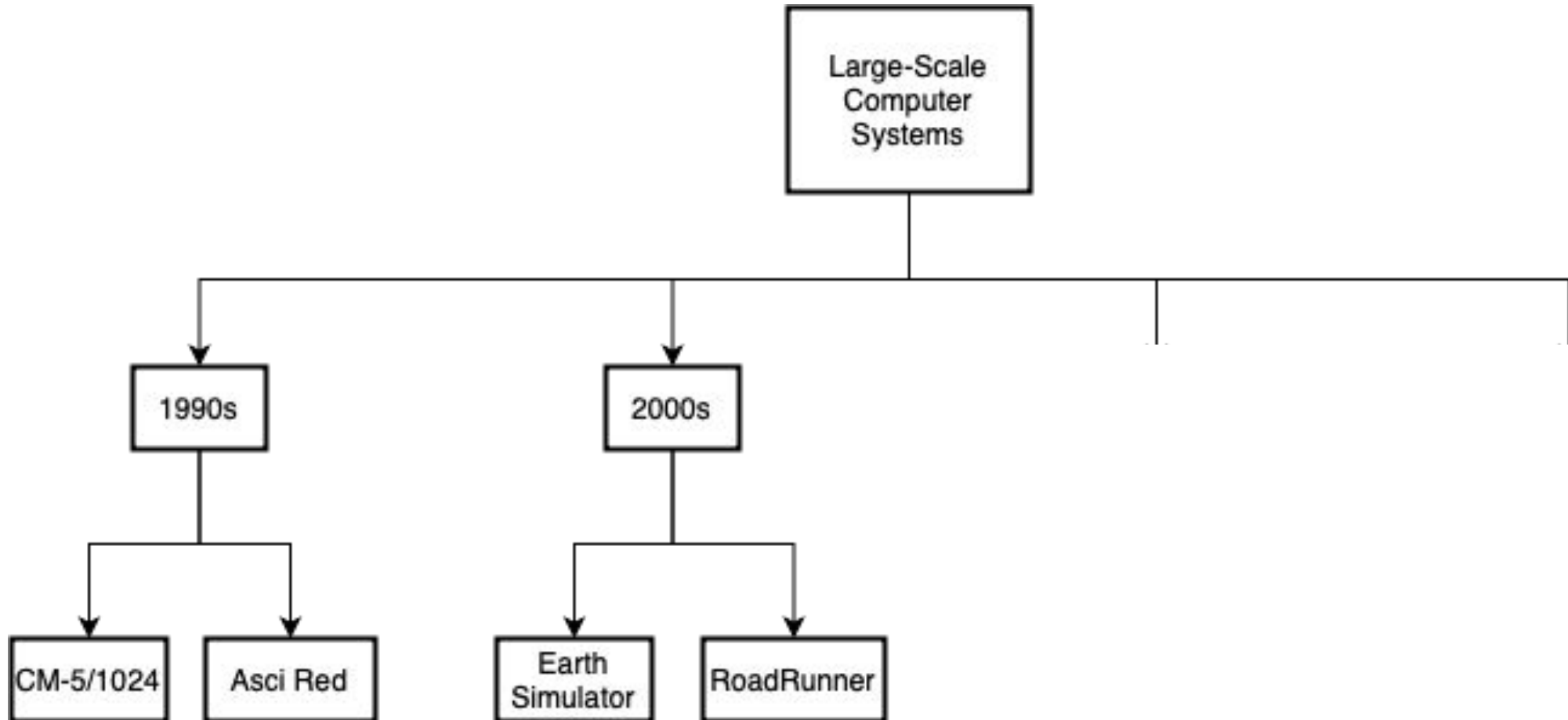
Taxonomy of Large-Scale Computer Systems from the 1990s to the 2020s

RQ1



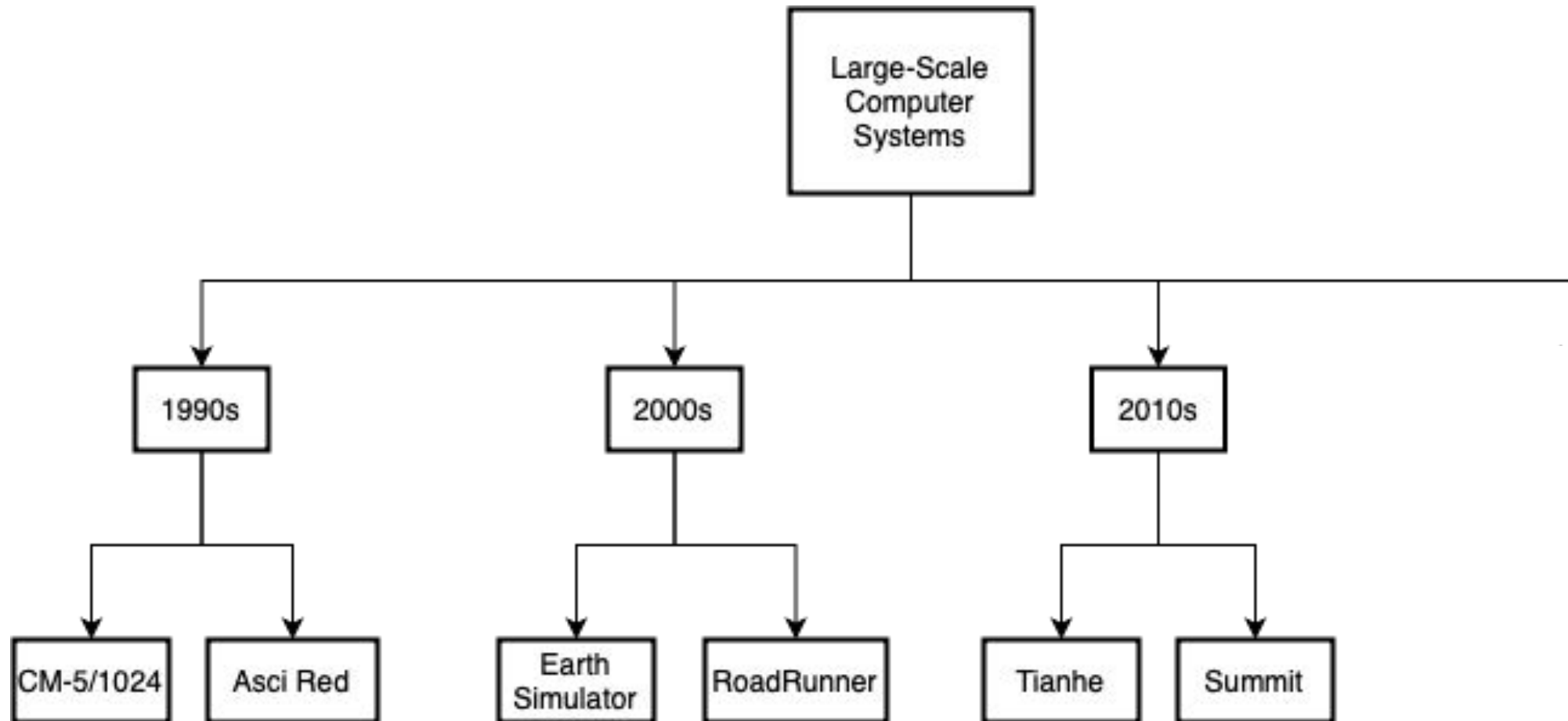
Taxonomy of Large-Scale Computer Systems from the 1990s to the 2020s

RQ1



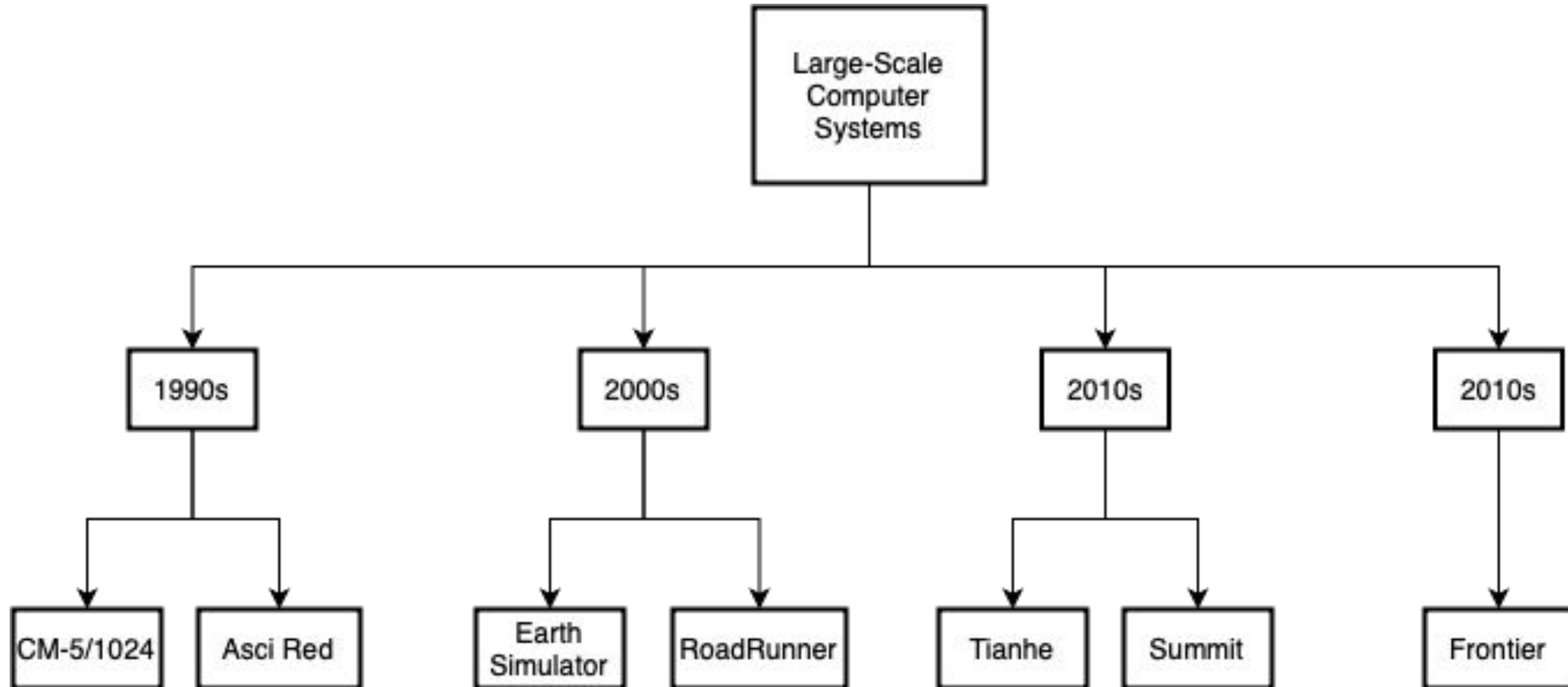
Taxonomy of Large-Scale Computer Systems from the 1990s to the 2020s

RQ1



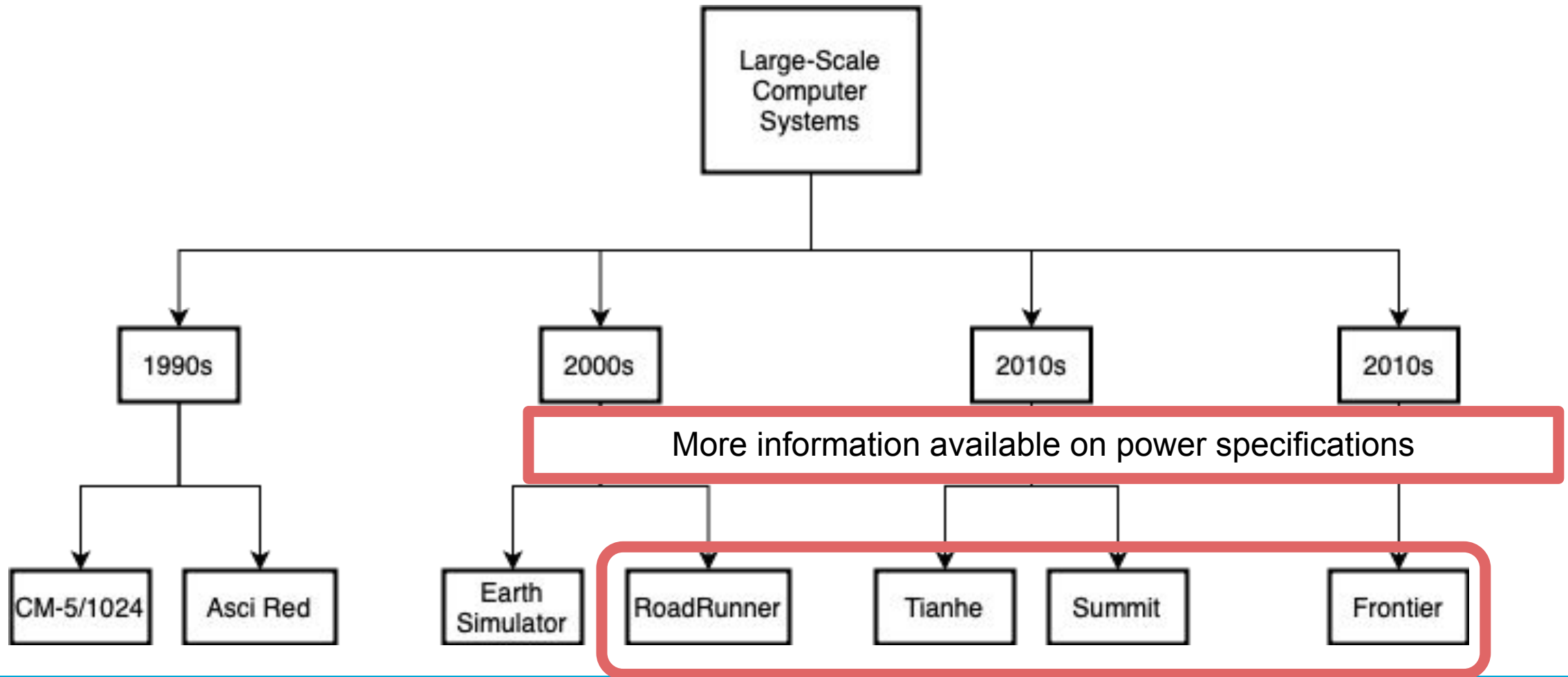
Taxonomy of Large-Scale Computer Systems from the 1990s to the 2020s

RQ1



Taxonomy of Large-Scale Computer Systems from the 1990s to the 2020s

RQ1



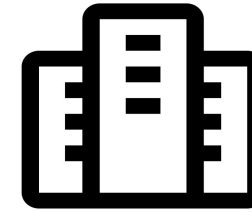
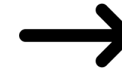
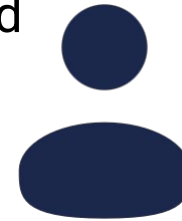
Modelling System Topologies and Workloads

RQ2

Properties	Bitbrains	SURF Lisa
Workload Type	Fixed-Time	Fixed-Work
Average Job Duration	28 days	3.1 hours
Number of Jobs	50	6,295

Workload Characteristics

User Defines Topology and Workload



Analyse Results



Simulate with OpenDC

System Name	Installation Year	Host Count	Core Count	Core Speed (MHz)	Peak Power (W)
Frontier	2021	9,472	64	2,000	240
Summit	2018	4,608	44	3,070	190
Tianhe	2013	17,792	24	2,200	115
RoadRunner	2008	11,340	9	3,200	90
Earth Simulator	2002	640	8	1,000	-
Asci Red	1997	4,536	2	200	-
CM-5/1024	1993	1,024	1	32	-

System Topologies

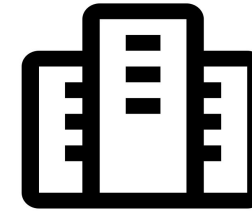
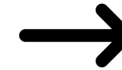
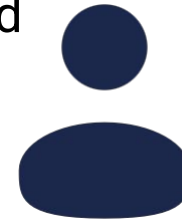
Modelling System Topologies and Workloads

RQ2

Properties	Bitbrains	SURF Lisa
Workload Type	Fixed-Time	Fixed-Work
Average Job Duration	28 days	3.1 hours
Number of Jobs	50	6,295

Workload Characteristics

User Defines Topology and Workload



Analyse Results

Simulate with OpenDC

System Name	Installation Year	Host Count	Core Count	Core Speed (MHz)	Peak Power (W)
Frontier	2021	9,472	64	2,000	240
Summit	2021	4,608	44	3,070	190
Tianhe	2013	17,792	24	2,200	115
RoadRunner	2008	11,340	9	3,200	90
Earth Simulator	2002	640	8	1,000	-
Asci Red	1997	4,536	2	200	-
CM-5/1024	1993	1,024	1	32	-

System Topologies

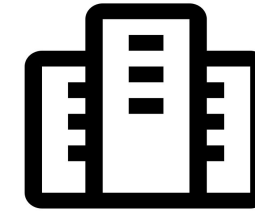
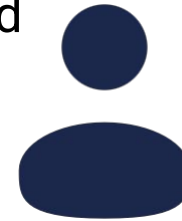
Modelling System Topologies and Workloads

RQ2

Properties	Bitbrains	SURF Lisa
Workload Type	Fixed-Time	Fixed-Work
Average Job Duration	28 days	3.1 hours
Number of Jobs	50	6,295

Workload Characteristics

User Defines Topology and Workload



Analyse Results

Simulate with OpenDC

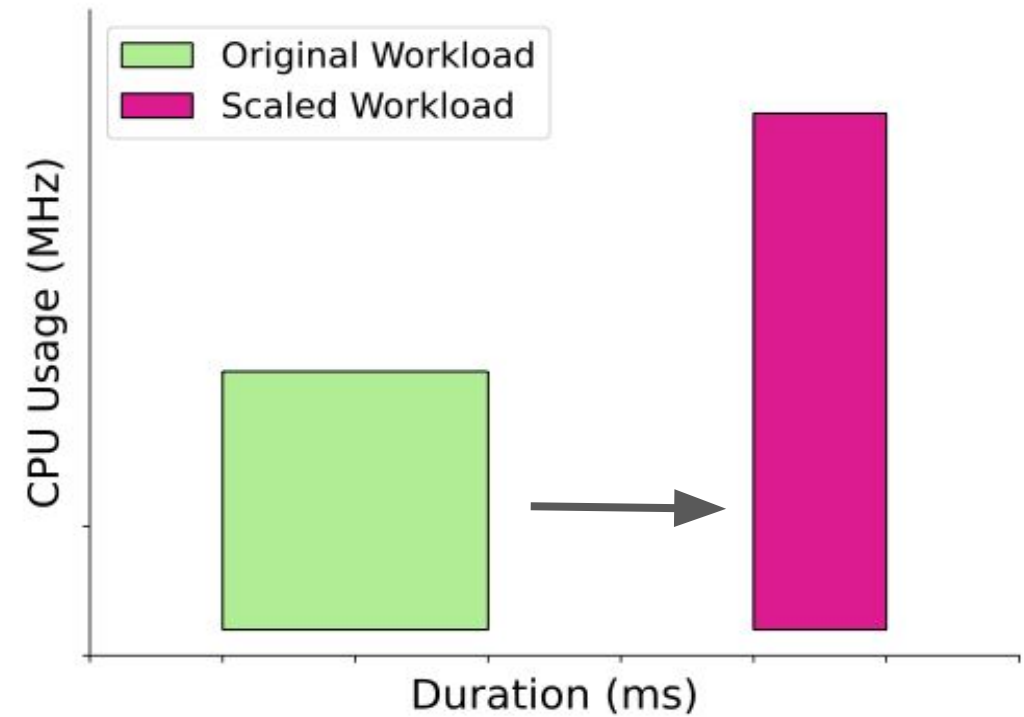
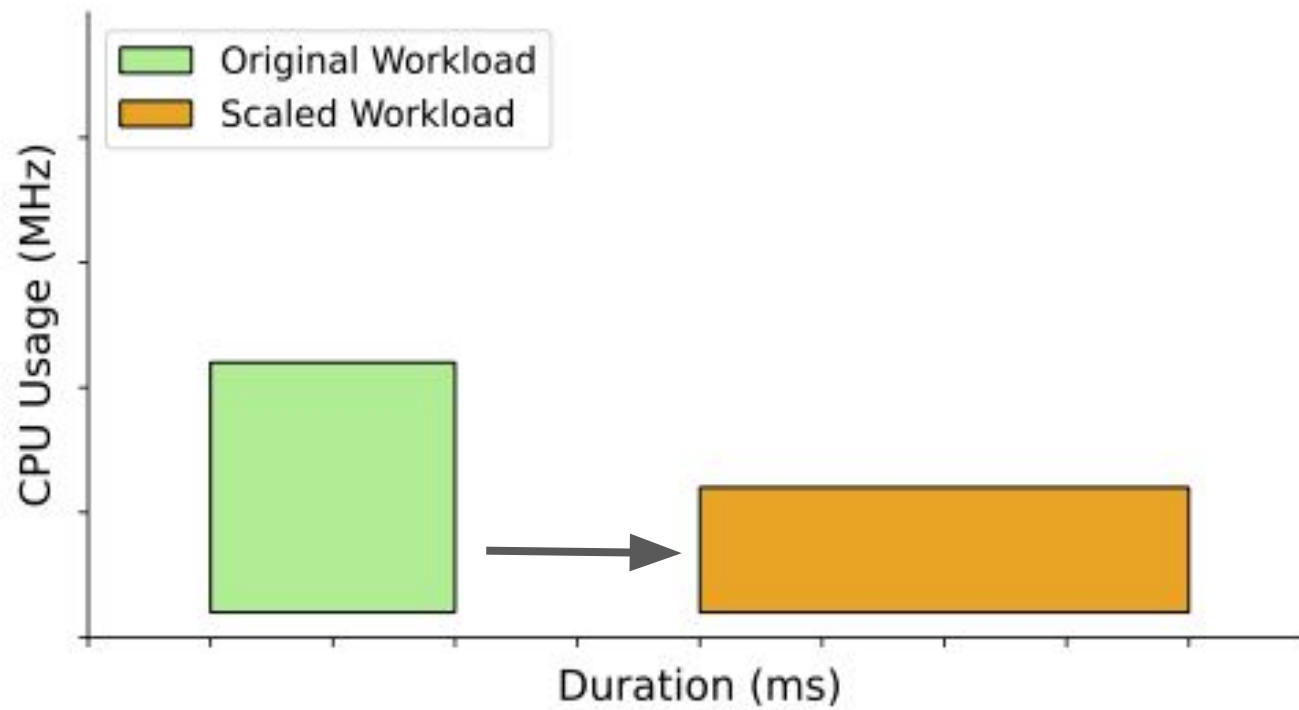
System Name	Installation Year	Host Count	Core Count	Core Speed (MHz)	Peak Power (W)
Frontier	2021	9,472	64	2,000	240
Summit	2018	4,608	44	3,070	190
Tianhe	2013	17,792	24	2,200	115
RoadRunner	2008	11,340	9	3,200	90
Earth Simulator	2002	640	8	1,000	-
Asci Red	1997	4,536	2	200	-
CM-5/1024	1993	1,024	1	32	-

System Topologies

Scaling Workloads for Diverse Configurations

RQ2

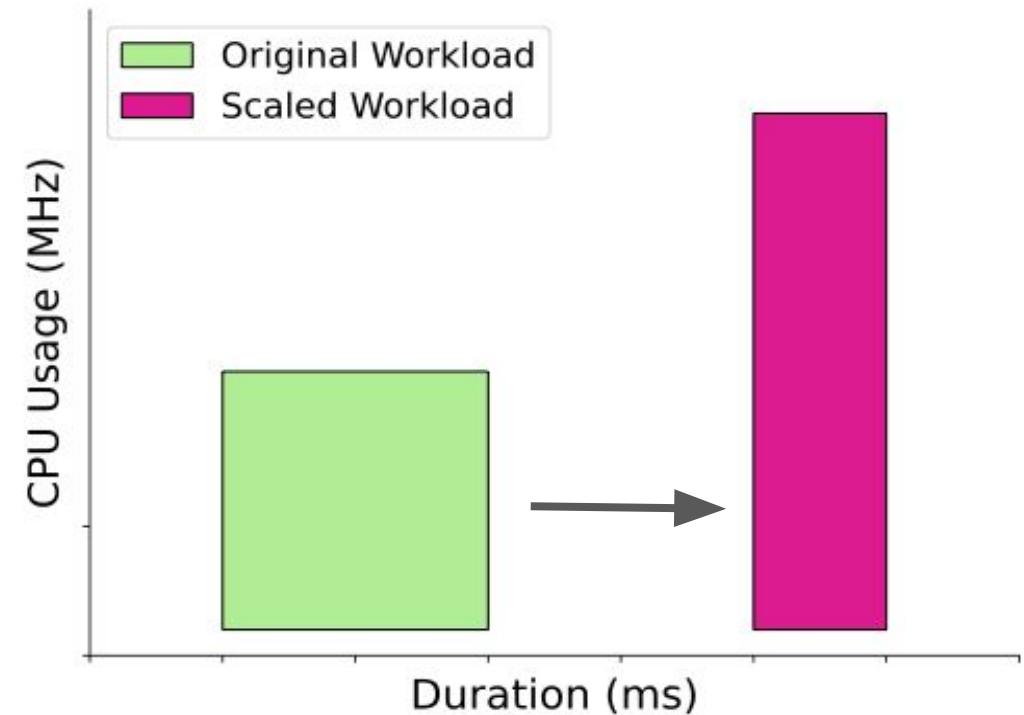
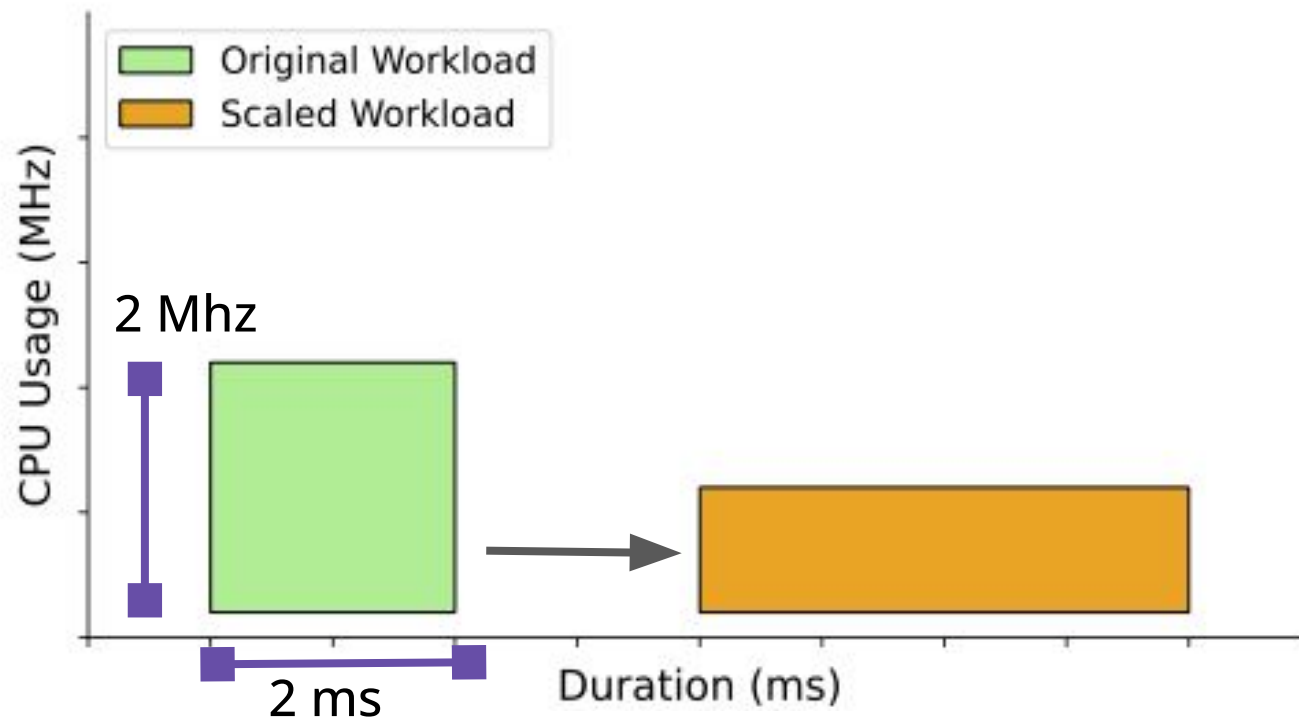
Fixed-Work



Scaling Workloads for Diverse Configurations

RQ2

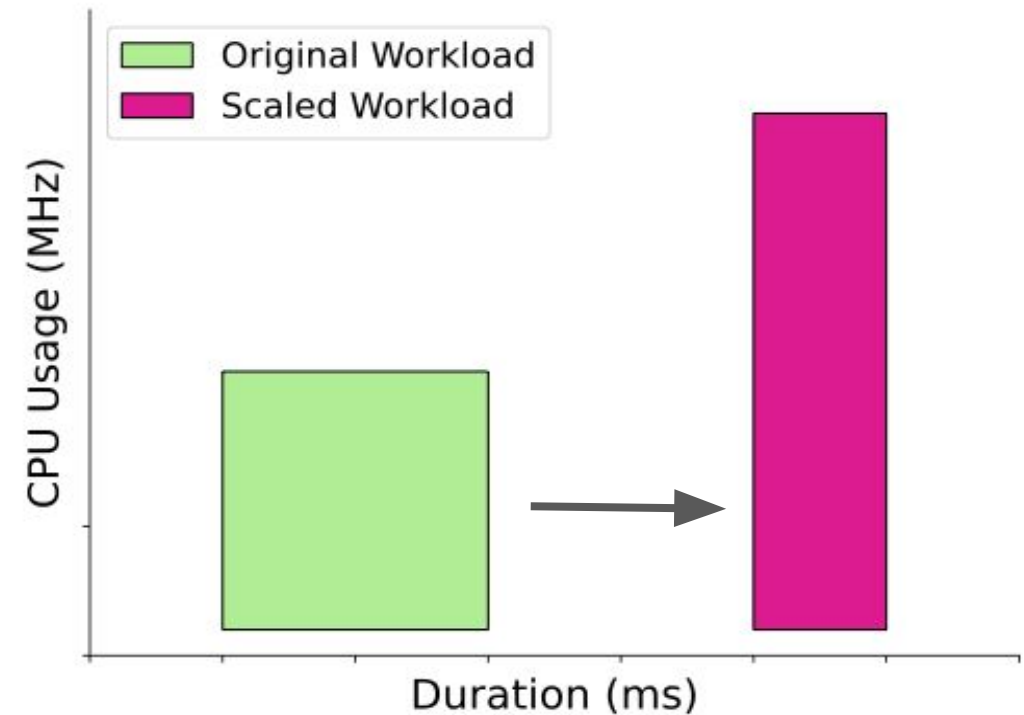
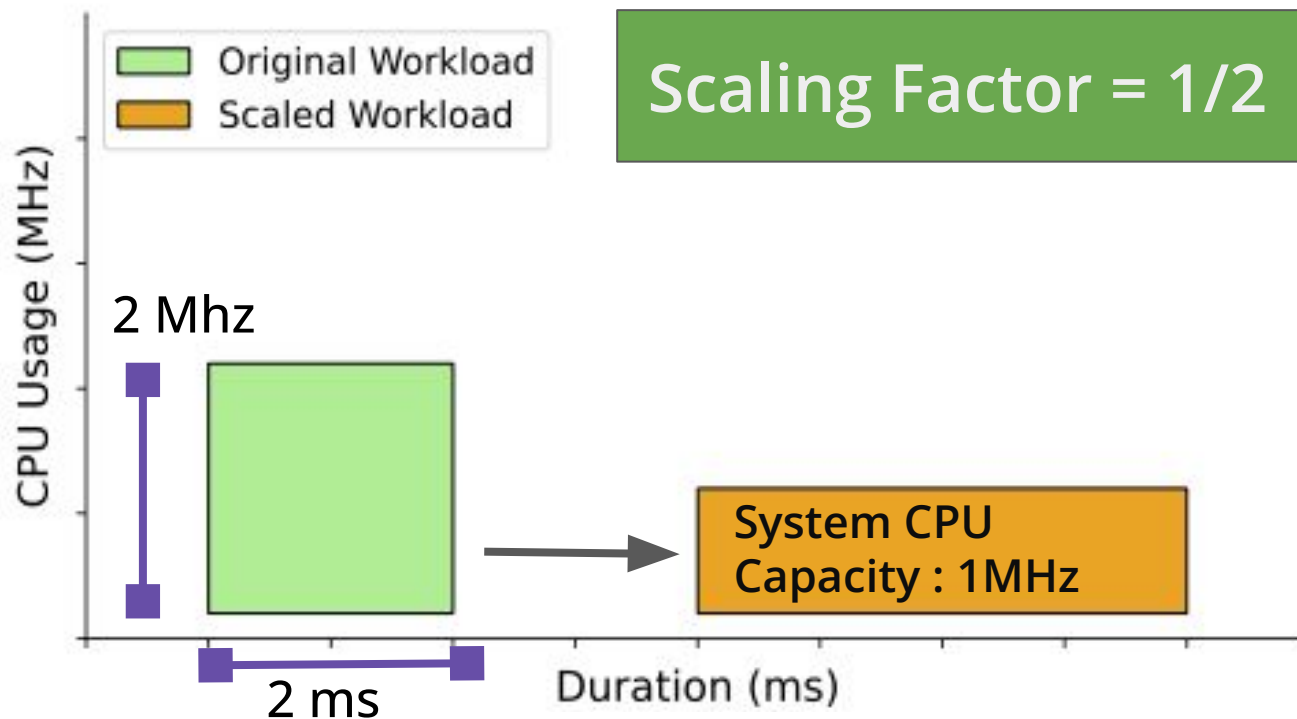
Fixed-Work



Scaling Workloads for Diverse Configurations

RQ2

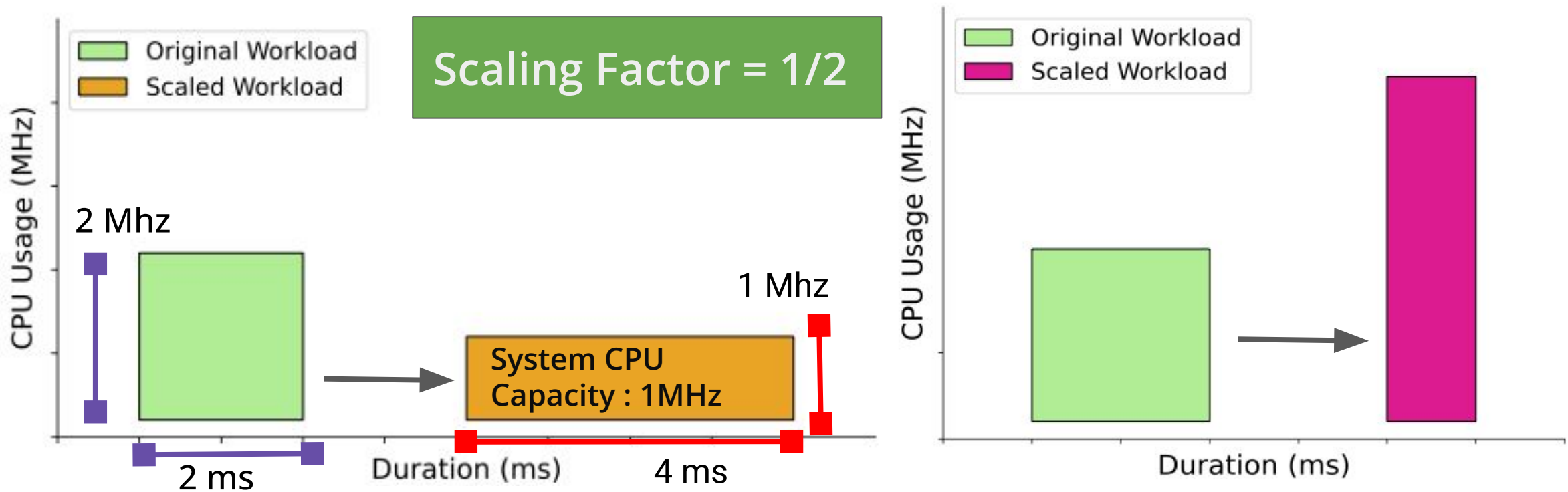
Fixed-Work



Scaling Workloads for Diverse Configurations

RQ2

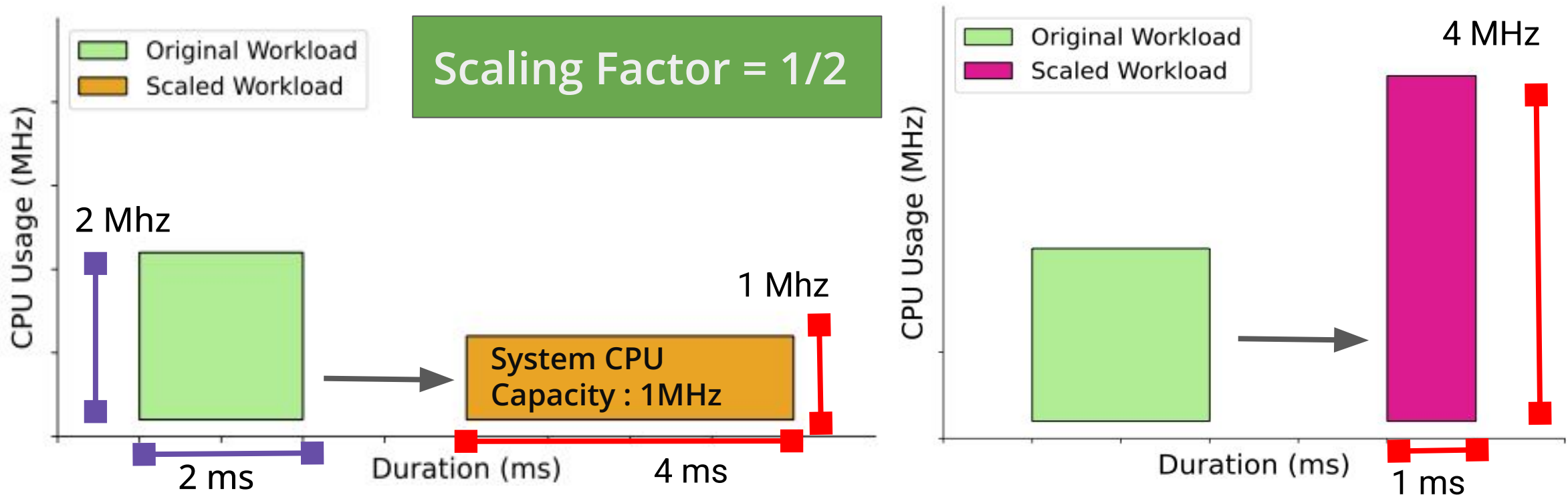
Fixed-Work



Scaling Workloads for Diverse Configurations

RQ2

Fixed-Work



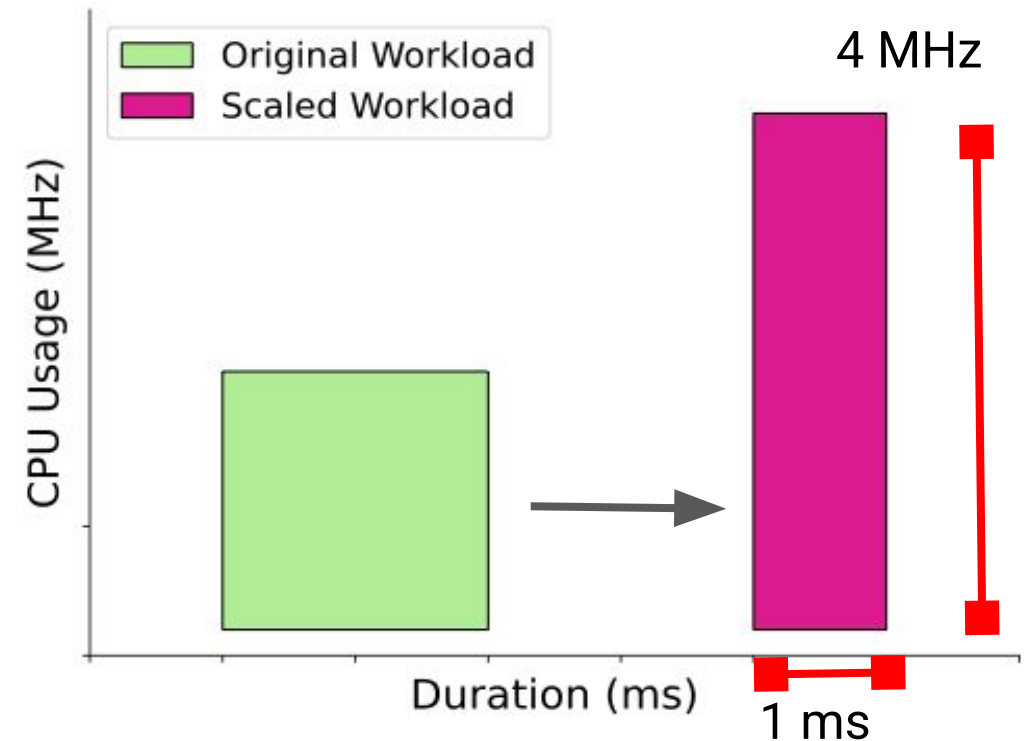
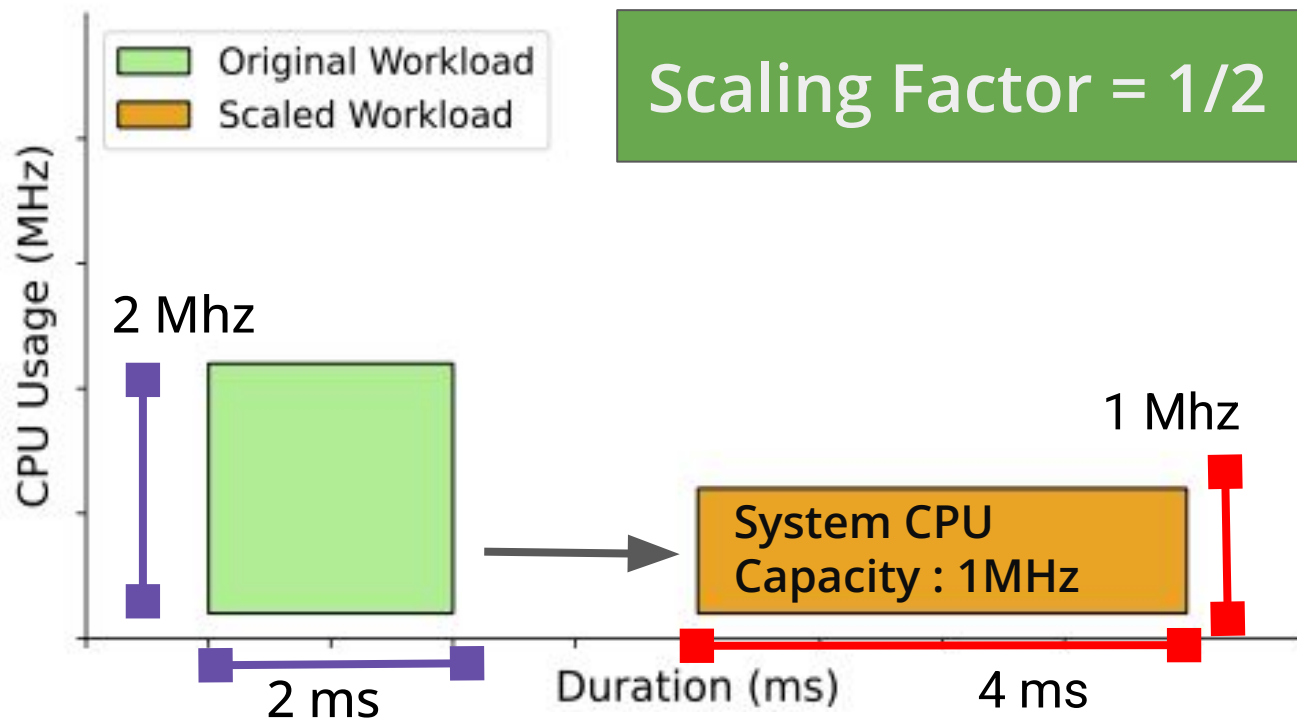
Scaling Workloads for Diverse Configurations

RQ2

Fixed-Work

Resource Consumption = CPU Usage × Duration of the Task

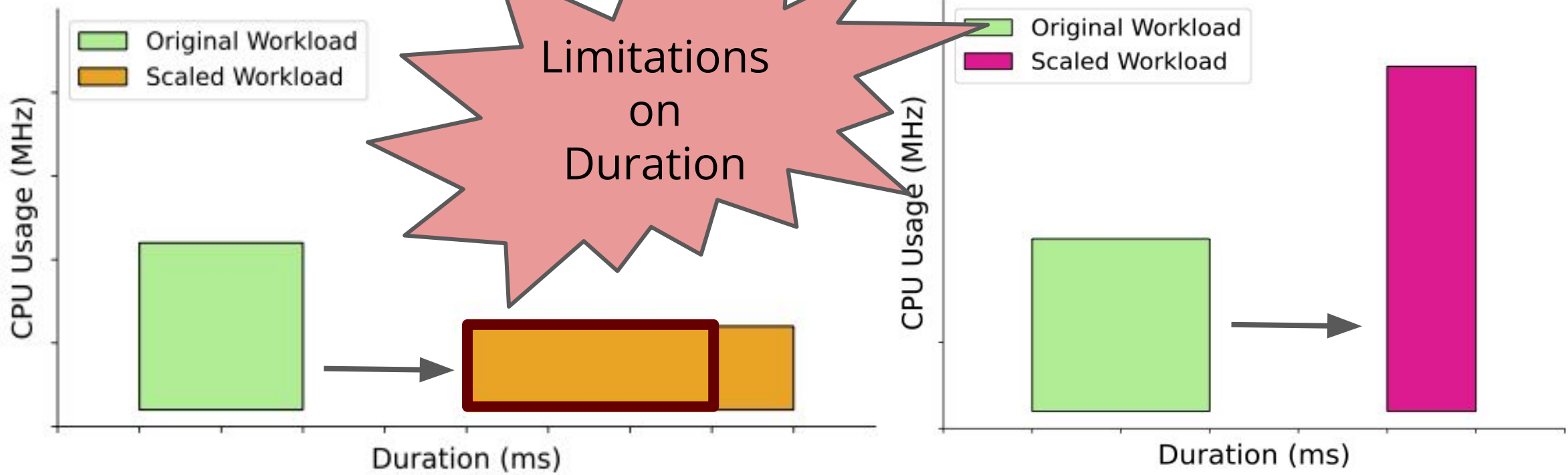
Scaling Factor = 1/2



Scaling Workloads for Diverse Configurations

RQ2

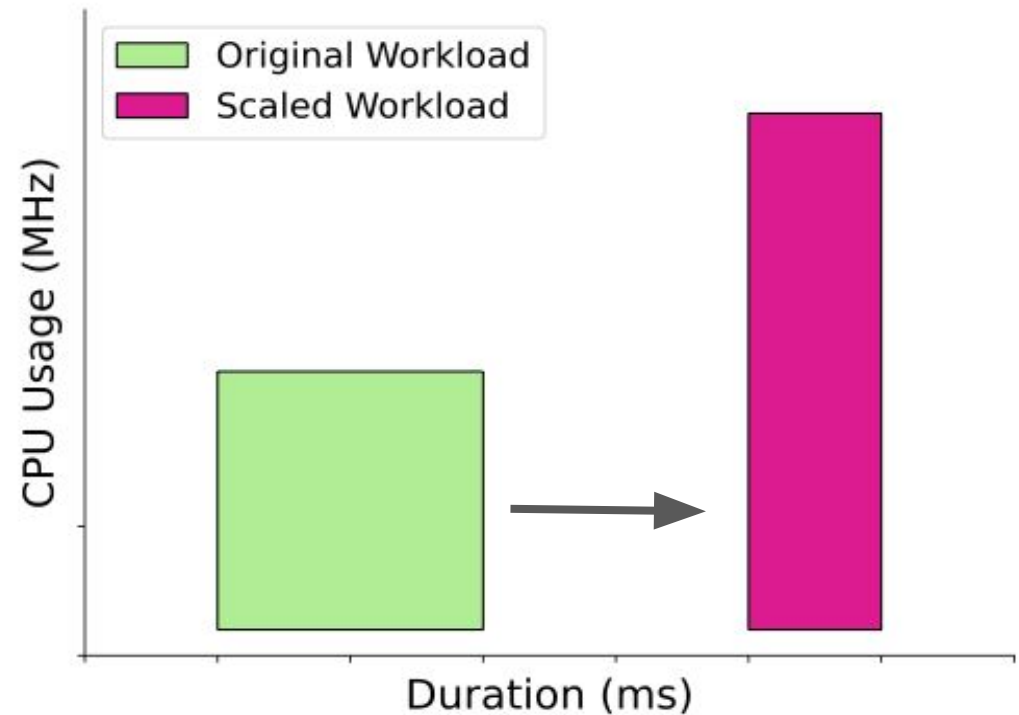
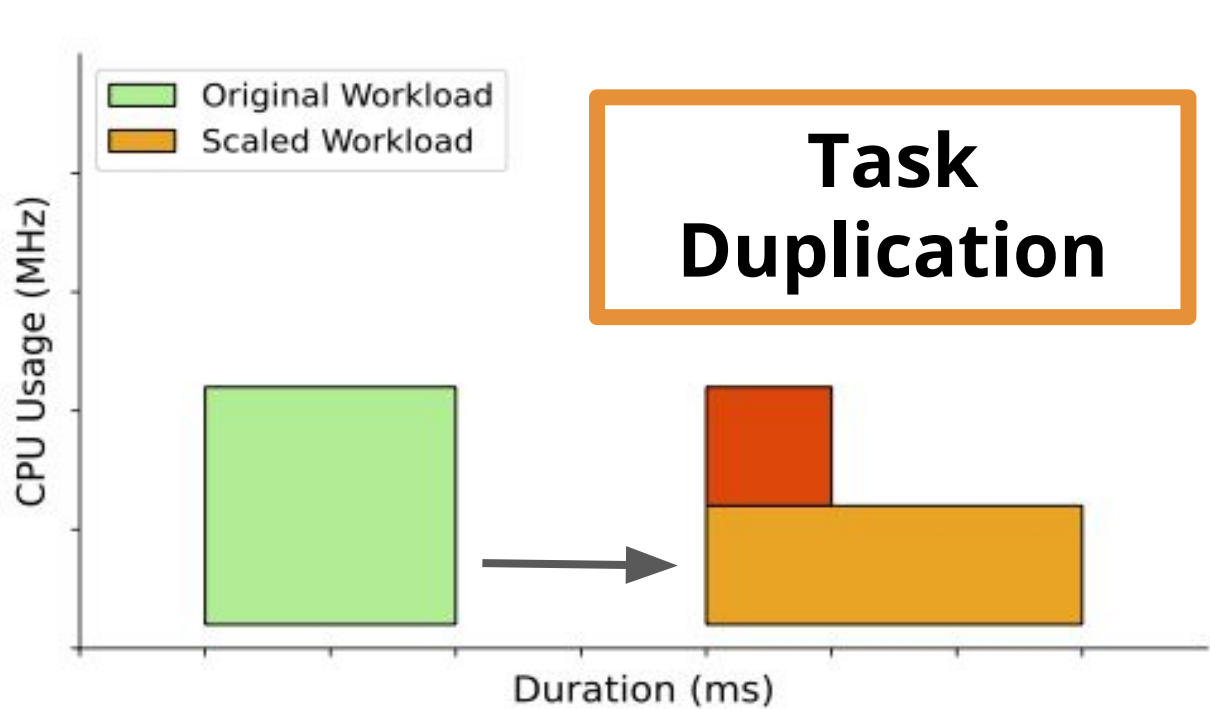
Fixed-Work



Scaling Workloads for Diverse Configurations

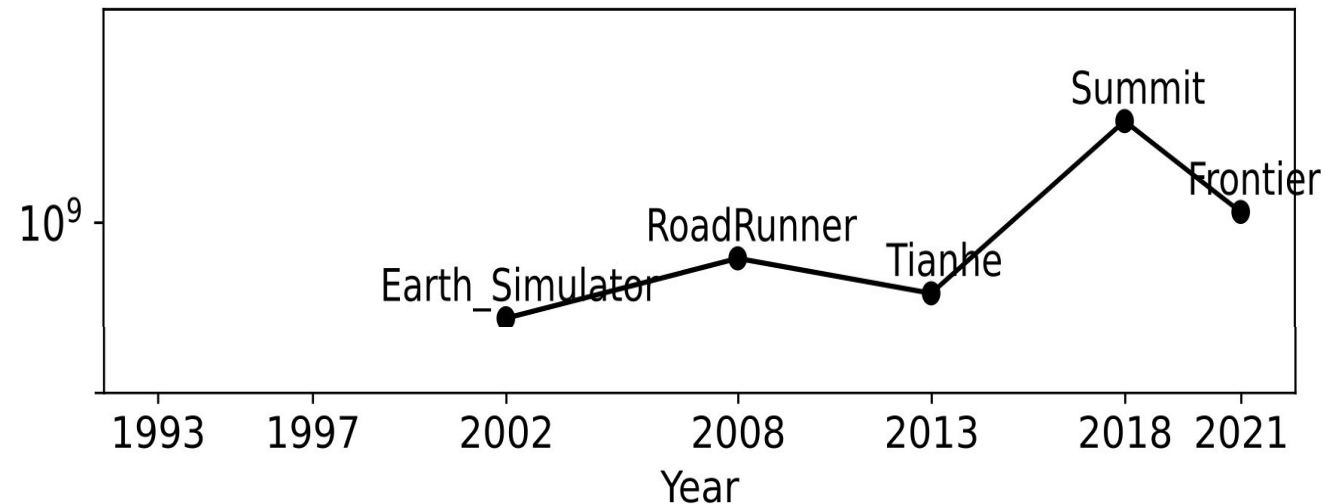
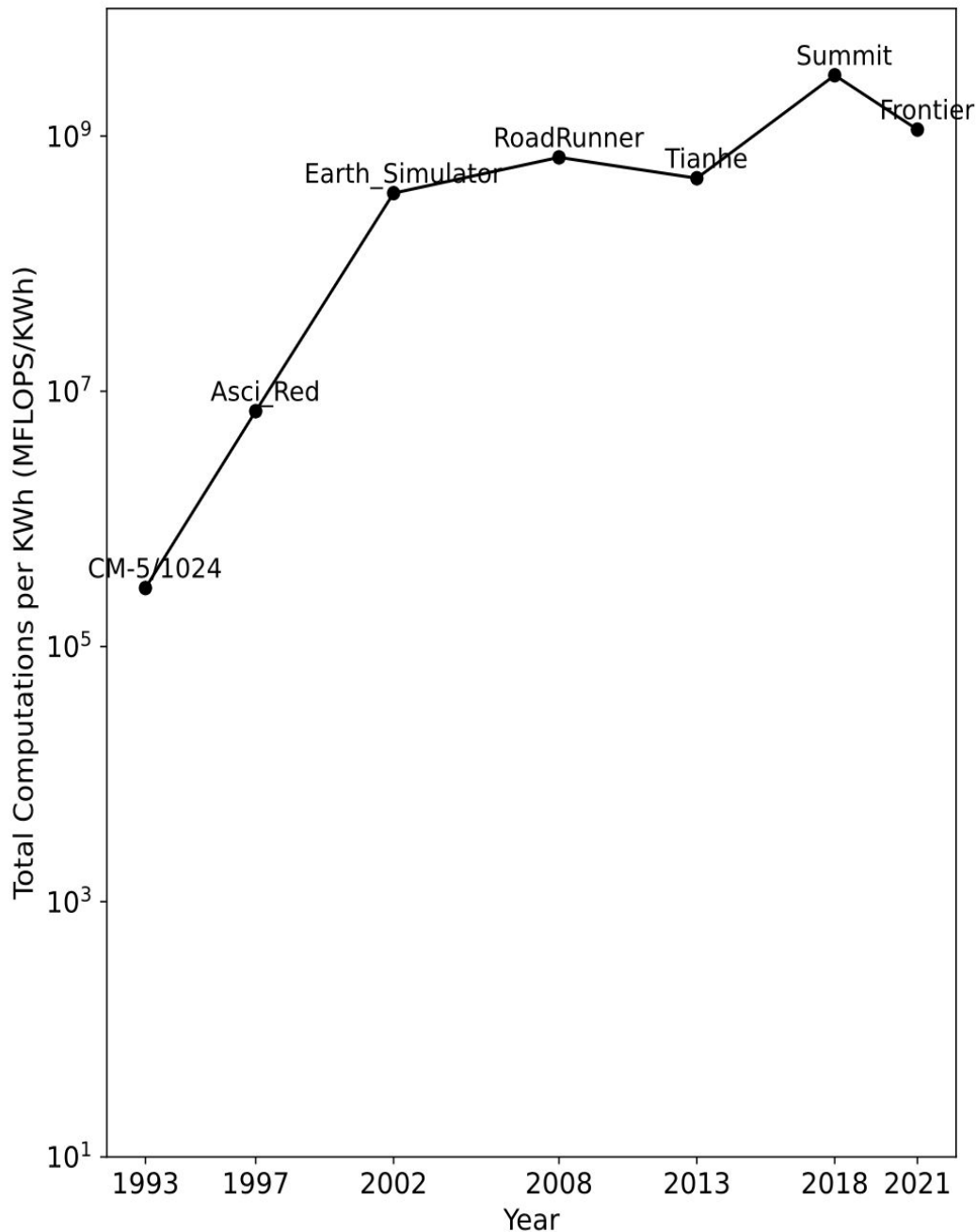
RQ2

Fixed-Work



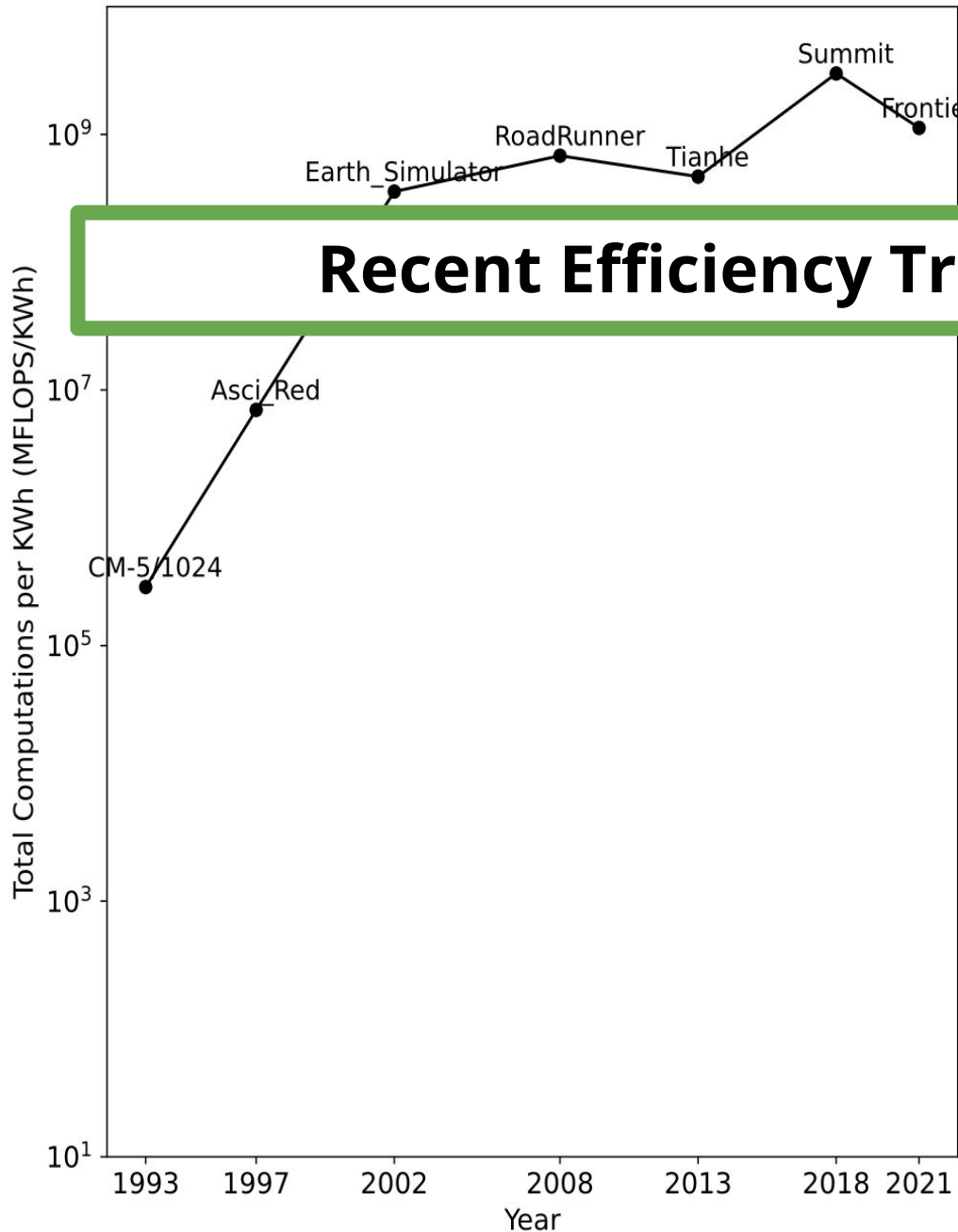
Evolution of Energy Efficiency

- Final results for SURF Lisa Workload
- Validating Koomey's Law
 - no failures, or interferences in this model

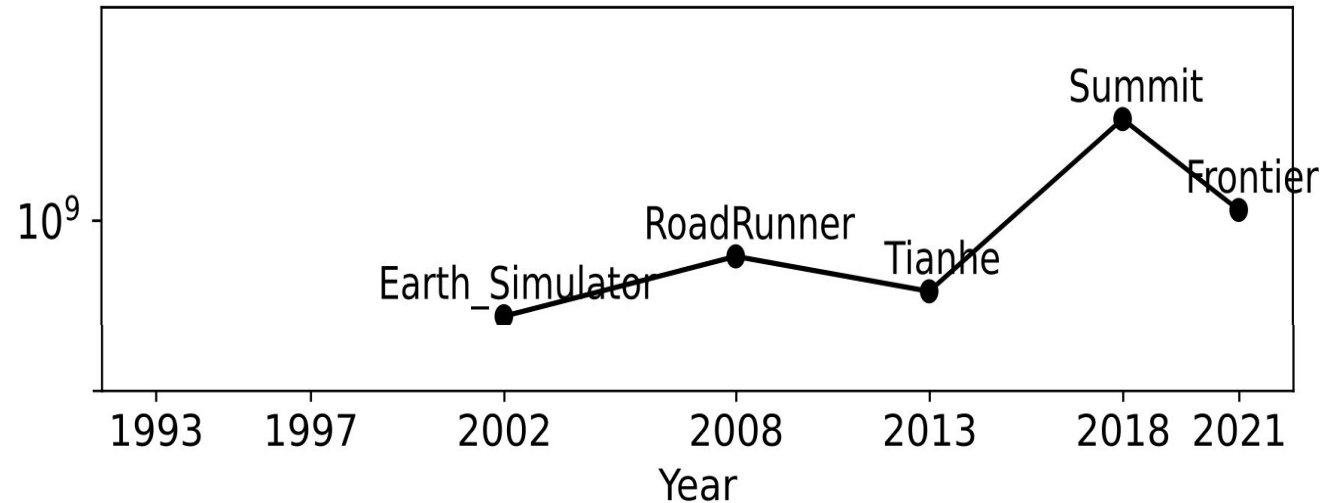


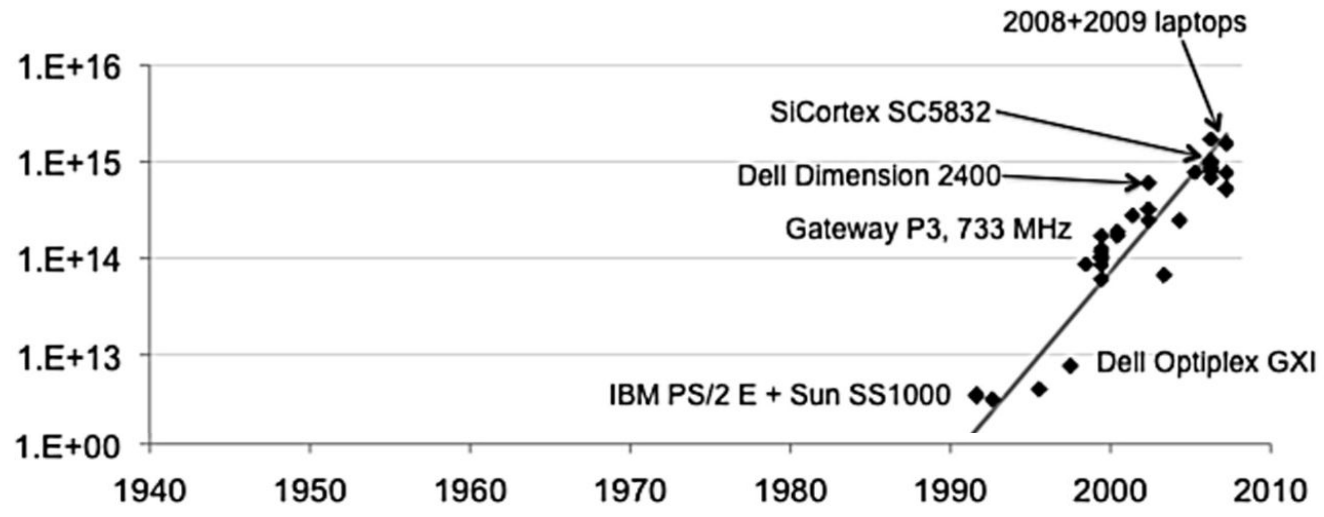
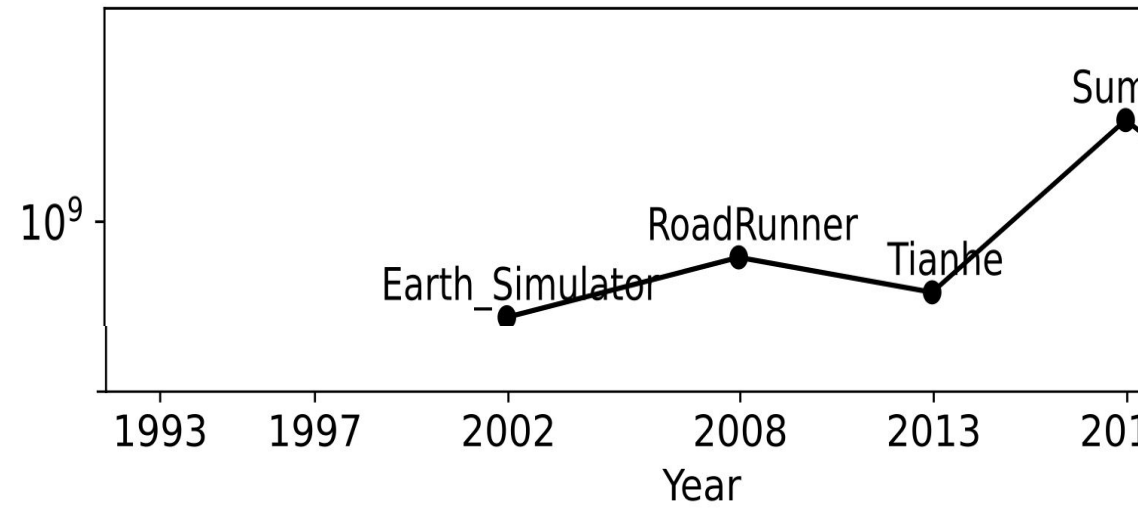
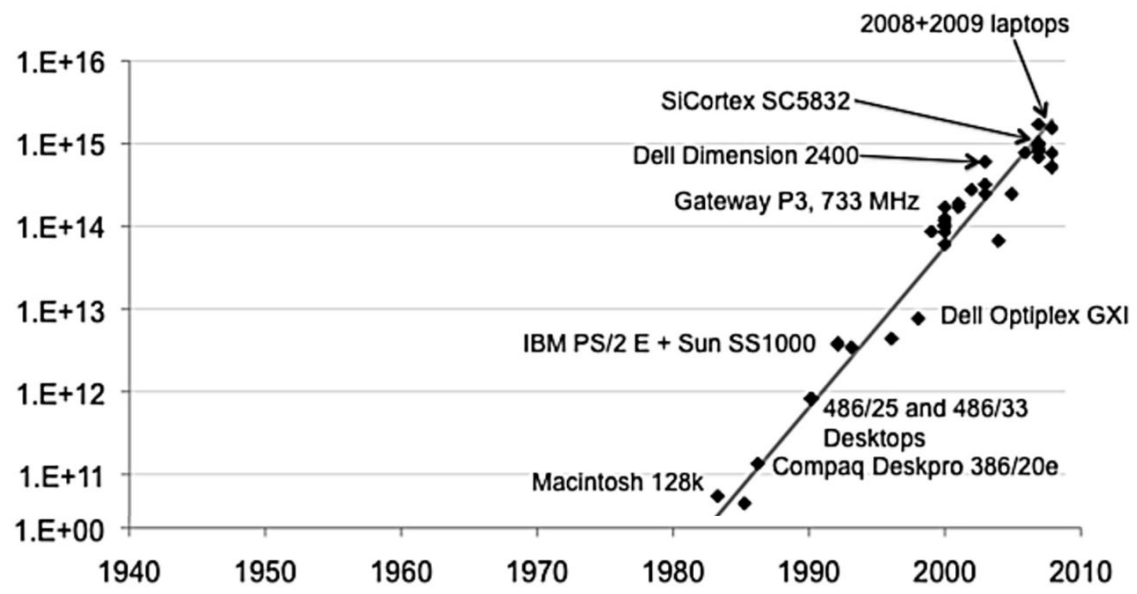
Evolution of Energy Efficiency

Recent Efficiency Trends Appear to be Slowing Down



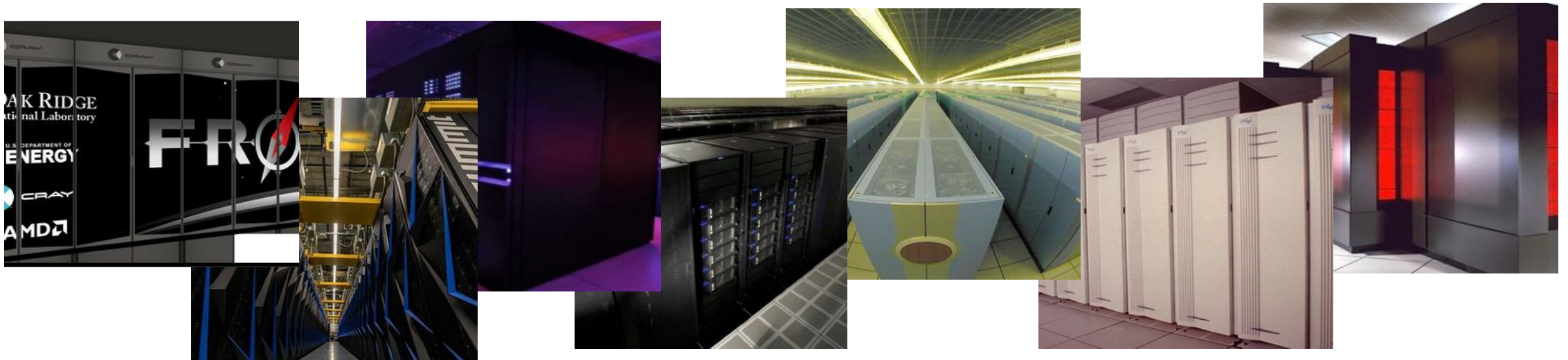
- Workload
- Validating Koomey's Law
 - no failures, or interferences in this model





Future Work / Conclusion

- We are the first to conduct a historical analysis on energy consumption under realistic conditions
- This project lays the groundwork for for more accurate analyses in the future



Future Work / Conclusion

- We are the first to conduct a historical analysis on energy consumption under realistic conditions
- This project lays the groundwork for for more accurate analyses in the future

Future Work:

- More demanding workloads
- Different scaling techniques
- Including more non-linearities

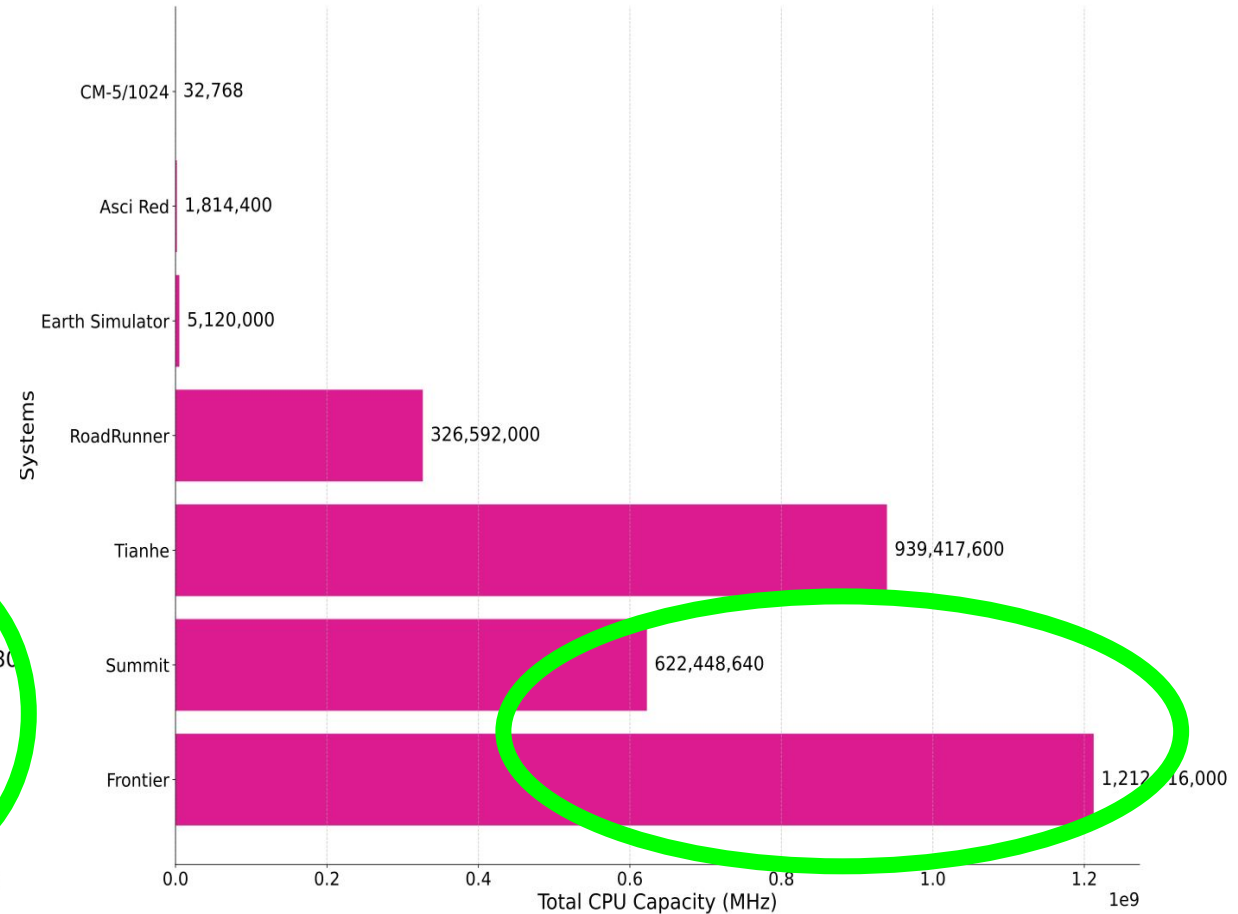
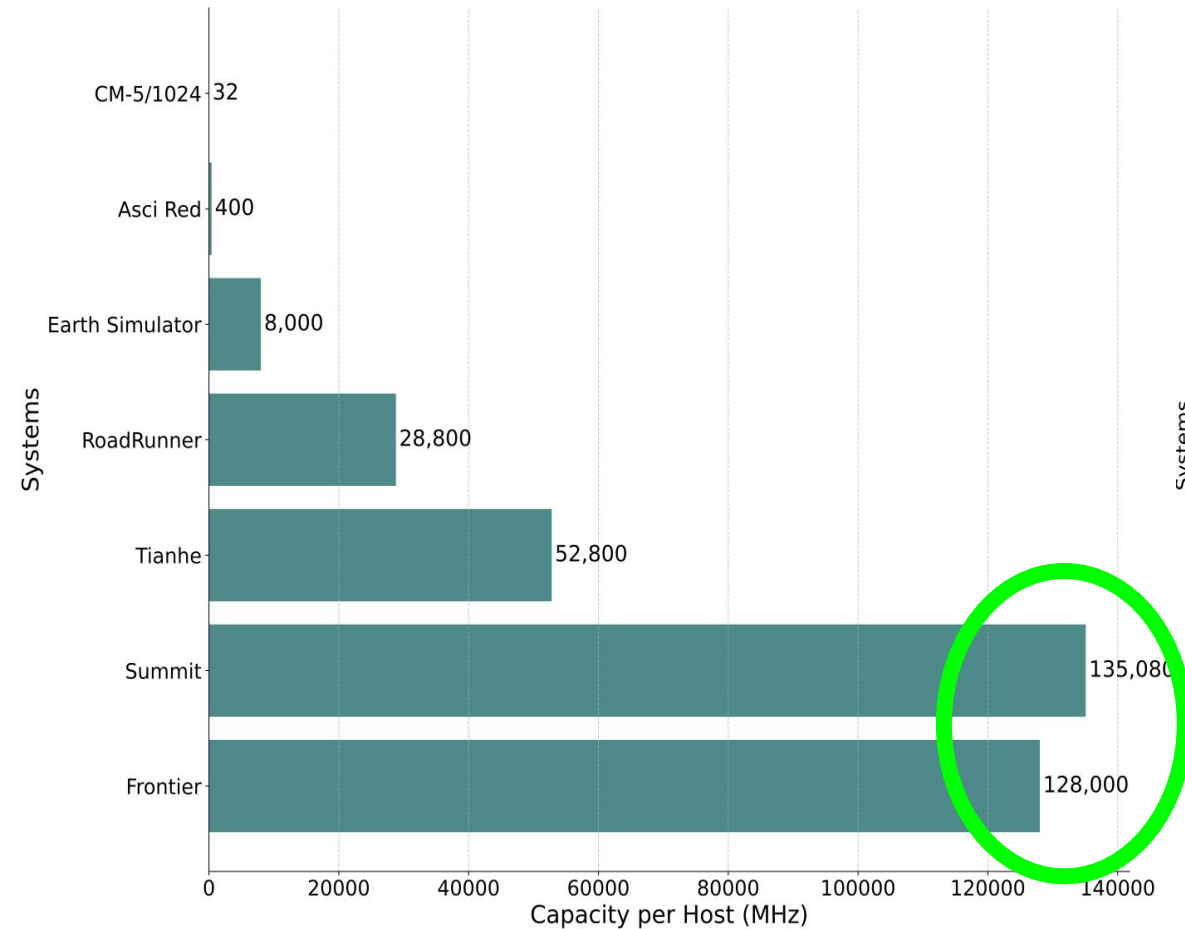


Thank you for Listening

CPU Capacity Per Host

Total CPU Capacity

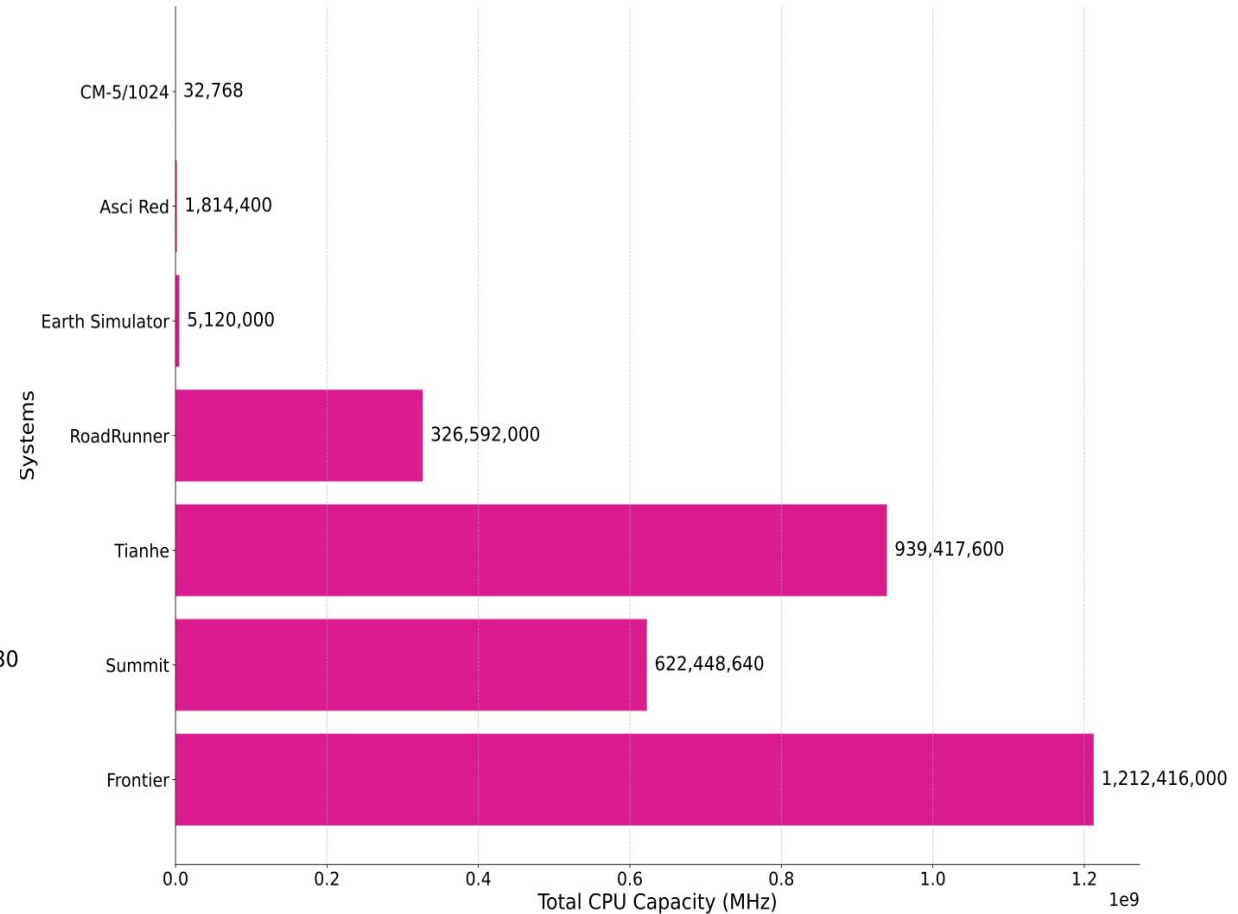
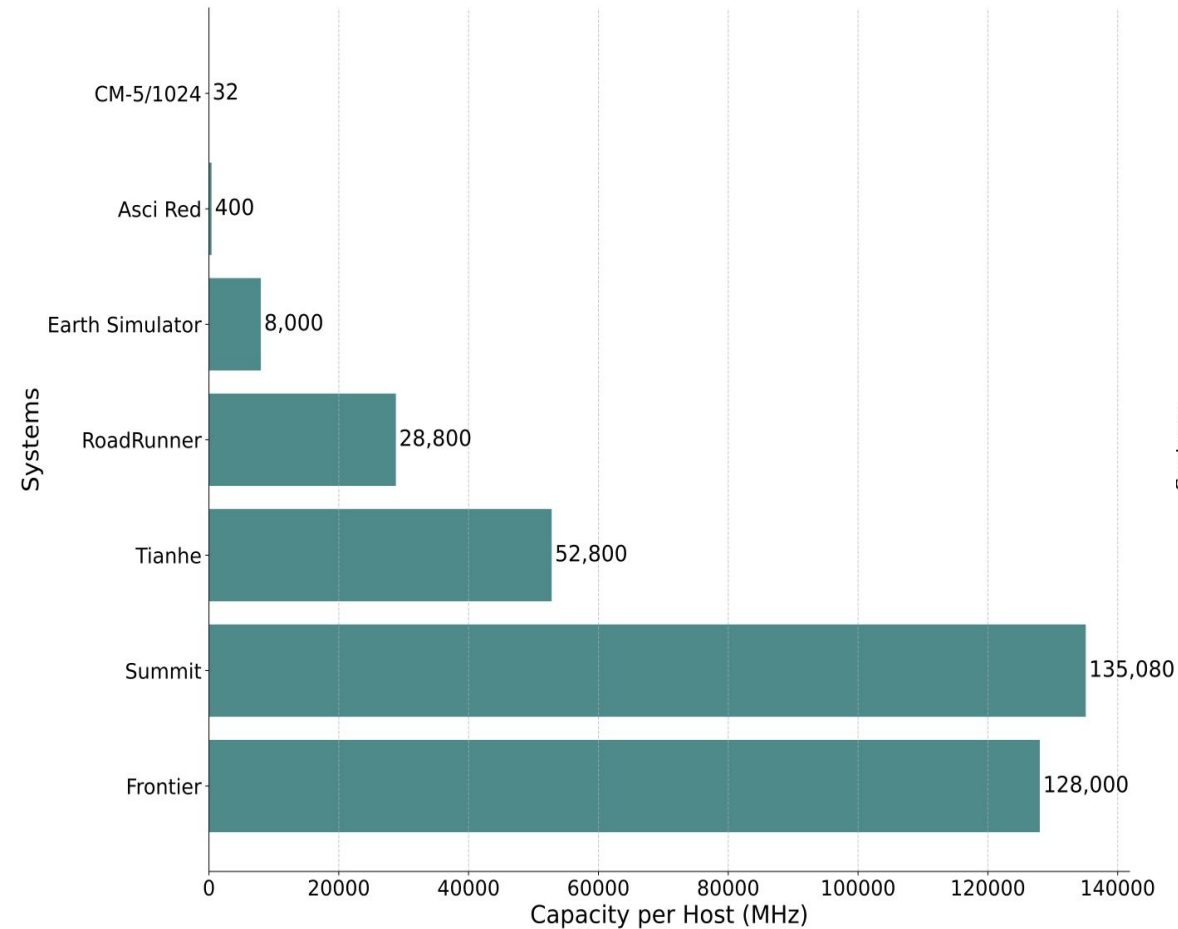
RQ3



CPU Capacity Per Host

Total CPU Capacity

RQ3

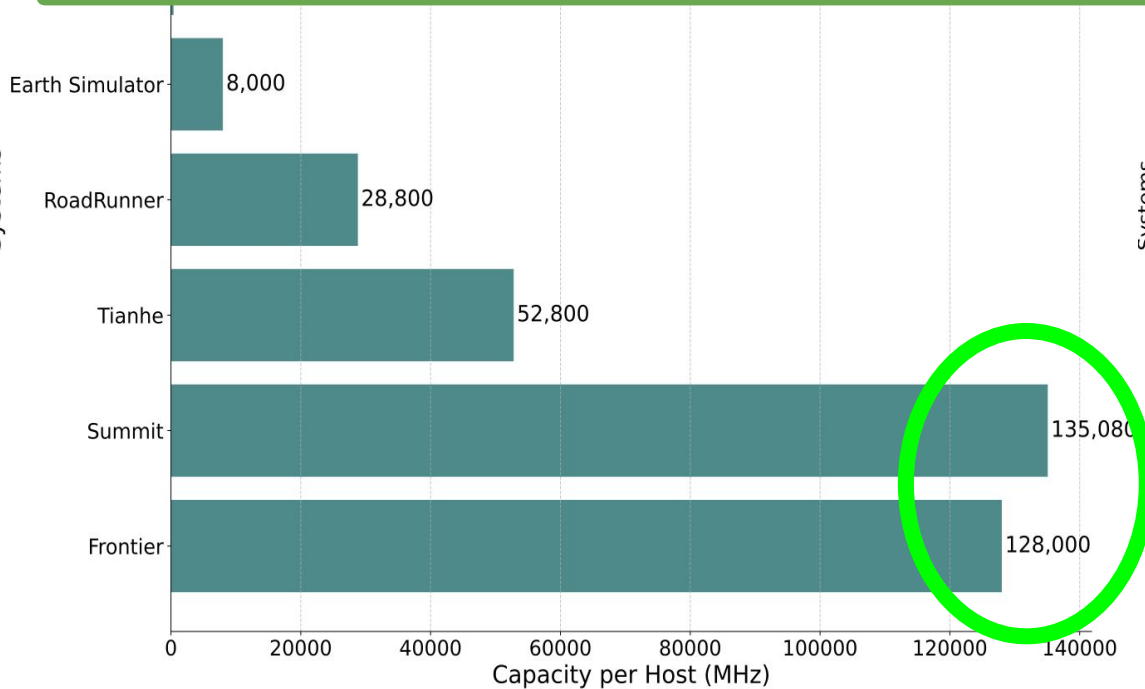


CPU Capacity Per Host

Total CPU Capacity

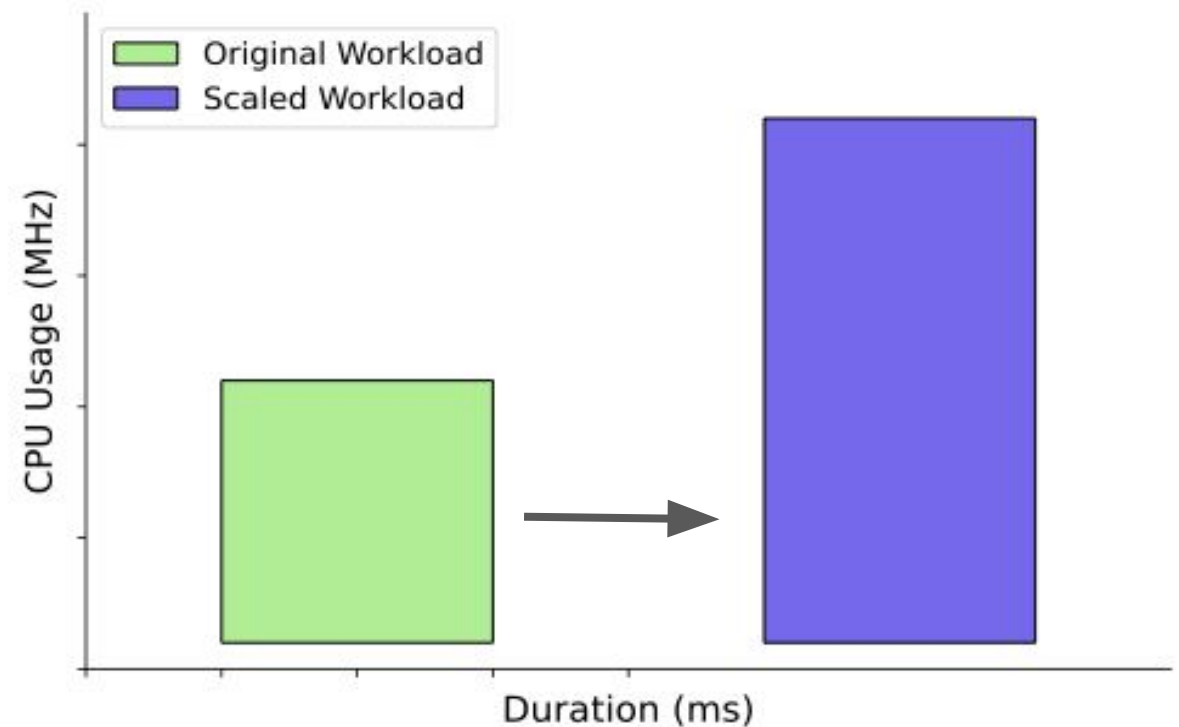
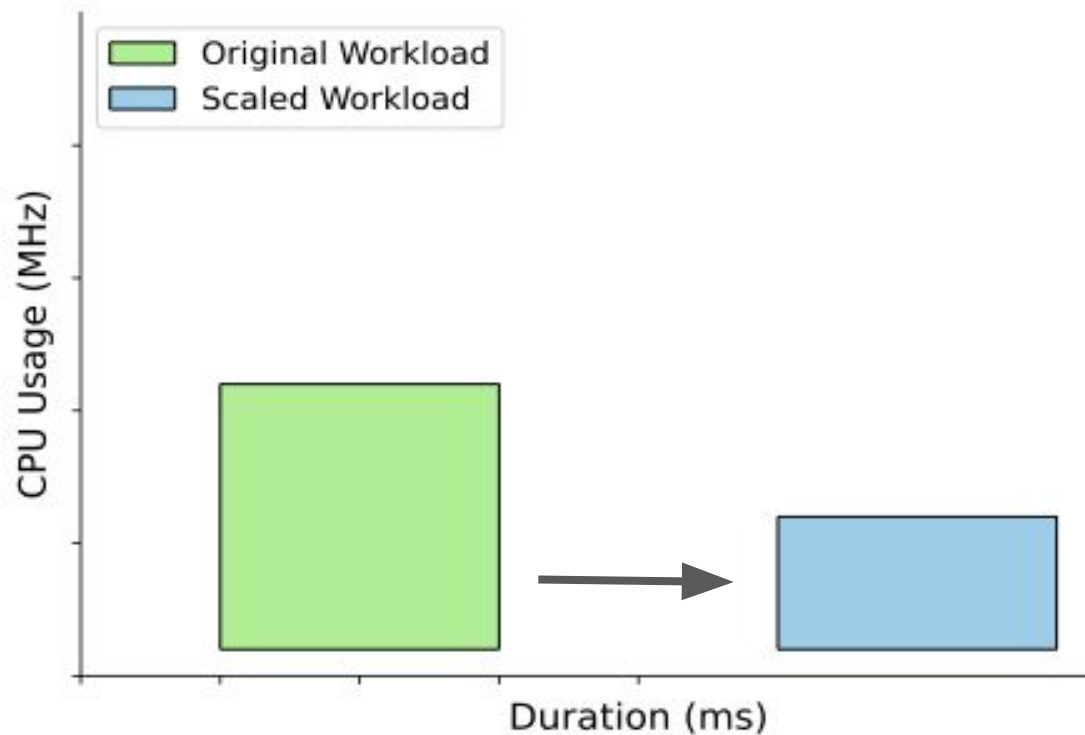
RQ3

Workload demand is key to revealing performance and efficiency differences in large-scale systems.



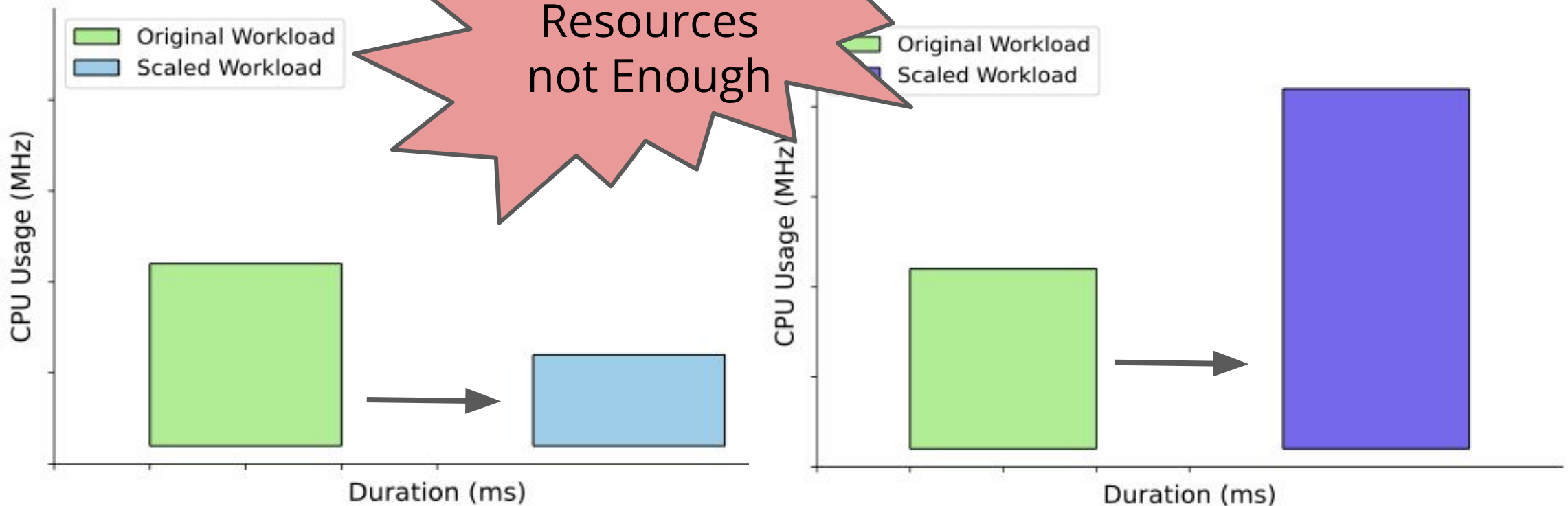
Scaling Workloads for Diverse Configurations

Fixed-Time



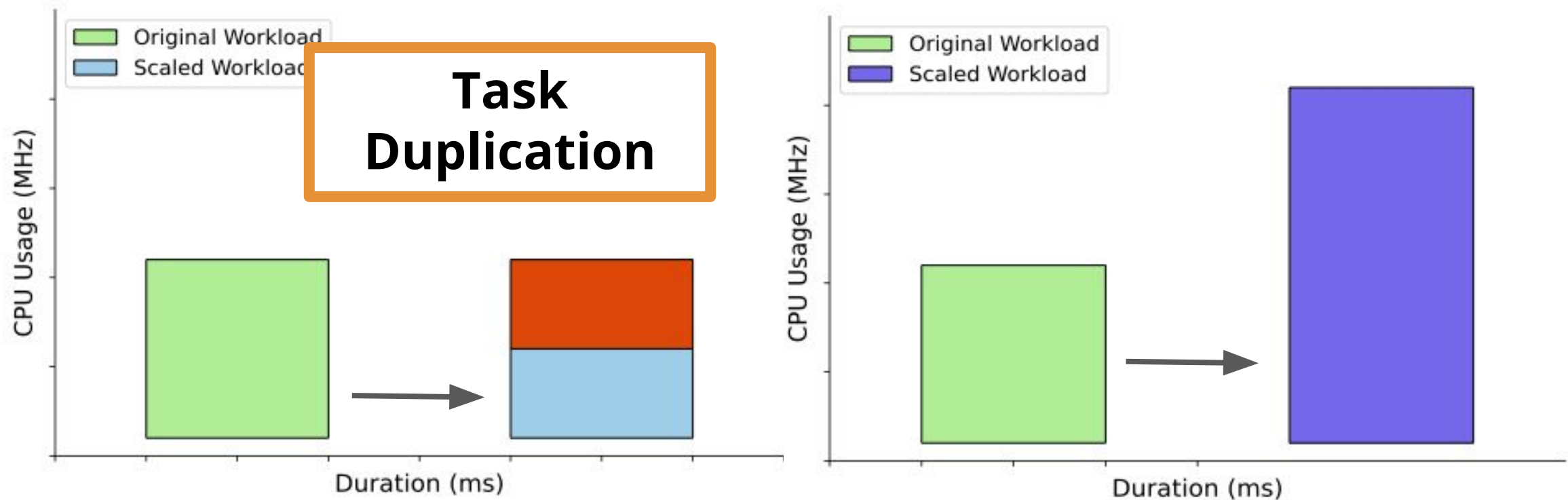
Scaling Workloads for Diverse Configurations

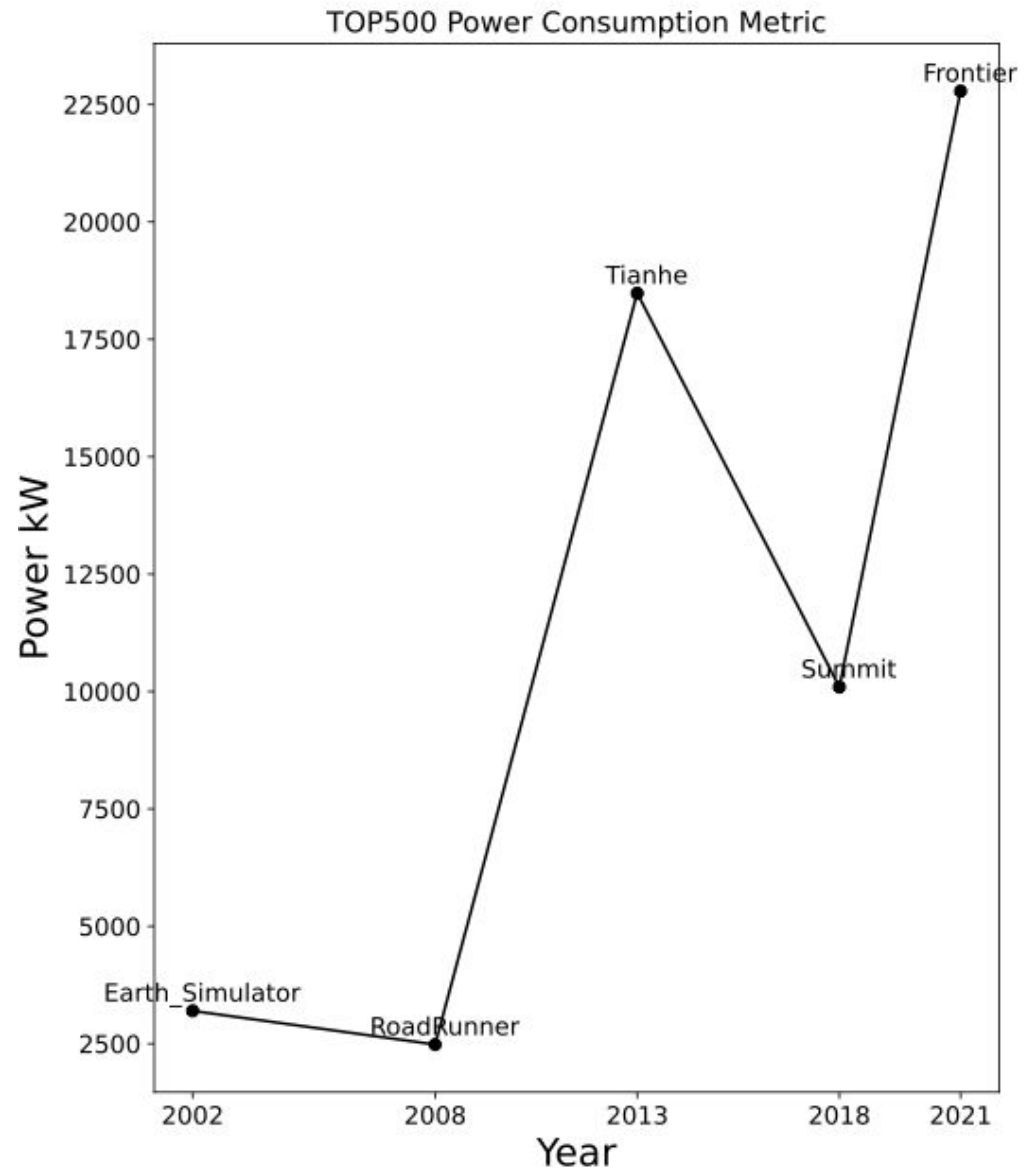
Fixed-Time



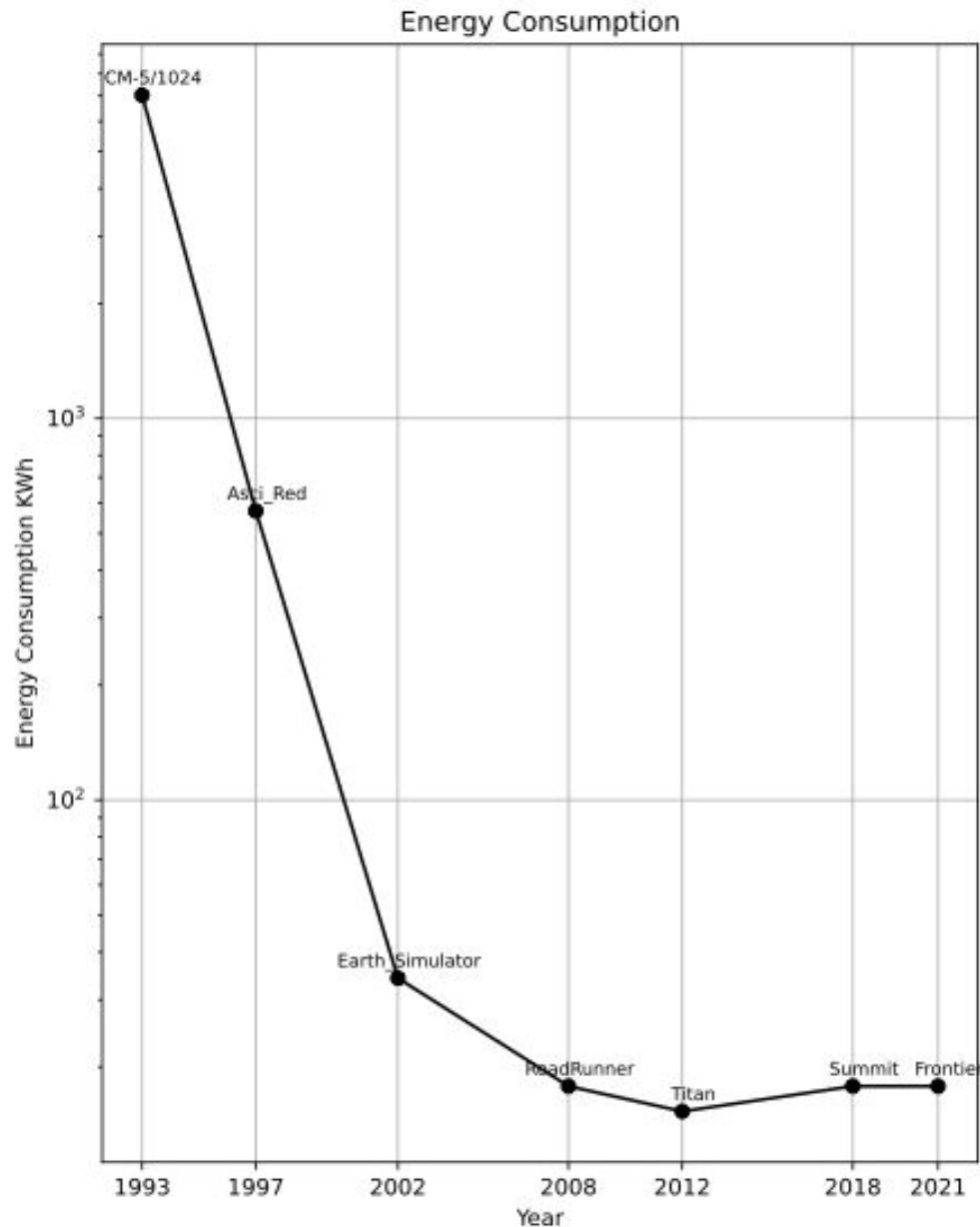
Scaling Workloads for Diverse Configurations

Fixed-Time





Energy Consumption



- No information of power model type and idle power specifications
- Constraints on the accuracy of energy consumption representation