

# AI IN THE POST-MOORE ERA

Babak Falsafi



[parsa.epfl.ch](http://parsa.epfl.ch)

**EPFL**

# OUR DIGITAL UNIVERSE



Fueled by:

- Data volume
- Data growth rate
- Monetization of data
- Data's impact on GDP
- ....now AI

# DATACENTERS: THE BACKBONE OF OUR DIGITAL UNIVERSE



- 100s of thousands of commodity or home-brewed servers
  - Consuming 10s to 100s MW
- Centralized to exploit economies of scale
- Network fabric w/  $\mu$ -second connectivity
- Often limited by ingress
  - Electricity
  - Network
  - Cooling



Boydton DC, 300MW

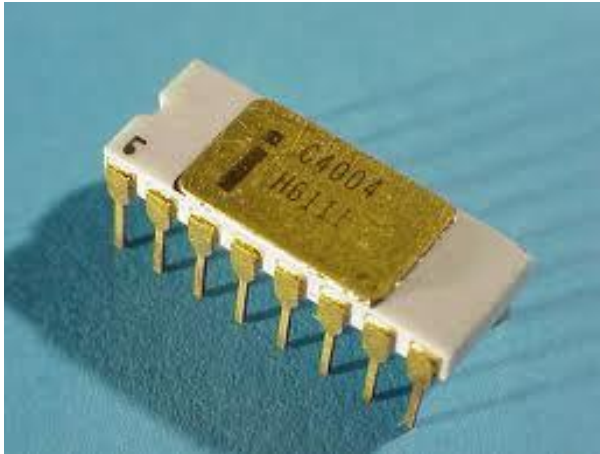
# CLOUDS AT VARIOUS SCALES



# UNIVERSE MADE POSSIBLE BY MOORE'S LAW

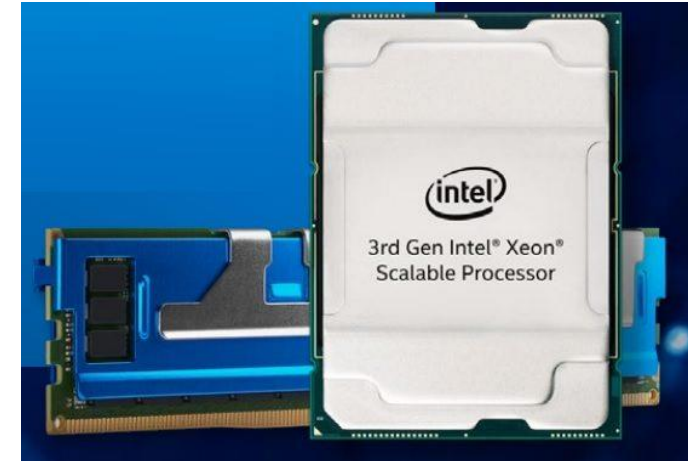


1971  
Intel 4004



92,000 ops/s  
1 Watt

2021  
Intel Ice Lake



1,200,000,000,000 ops/s  
270 Watts



# MOORE'S LAW: EXPONENTIAL DENSITY & EFFICIENCY



1971

Intel 4004



2021

Intel Ice Lake



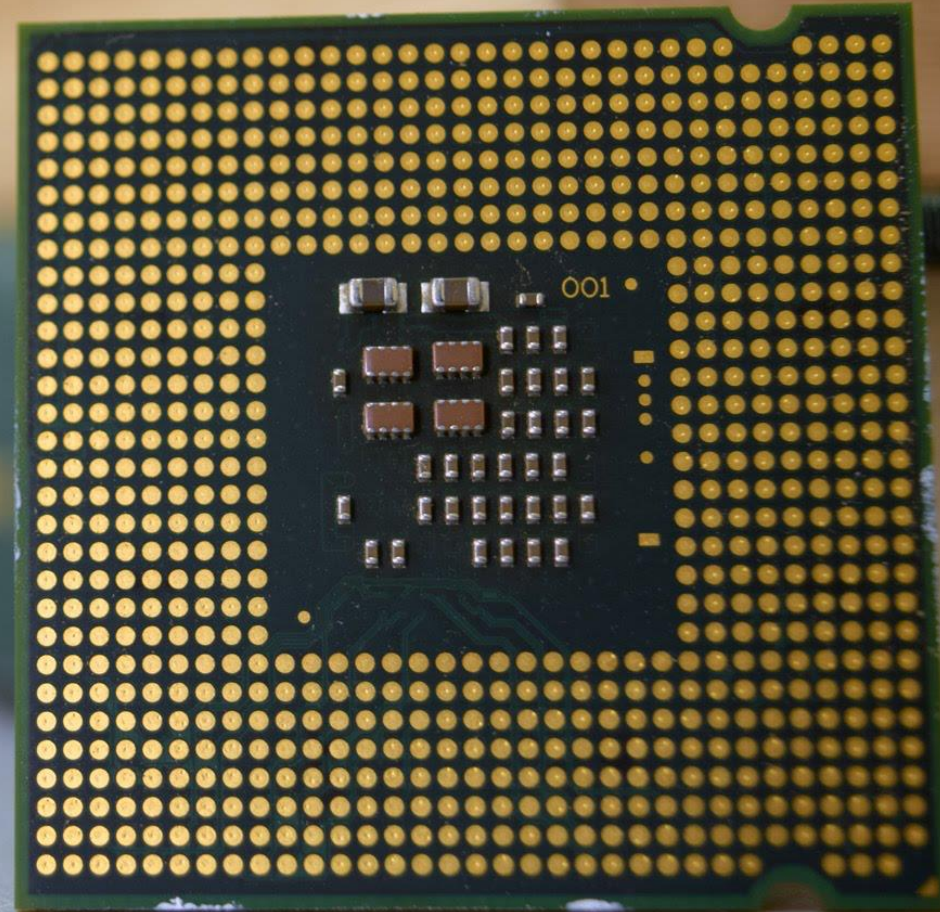
In 50 years:  
13 million times faster  
48 thousand times more efficient

92,000 ops/s  
1 Watt

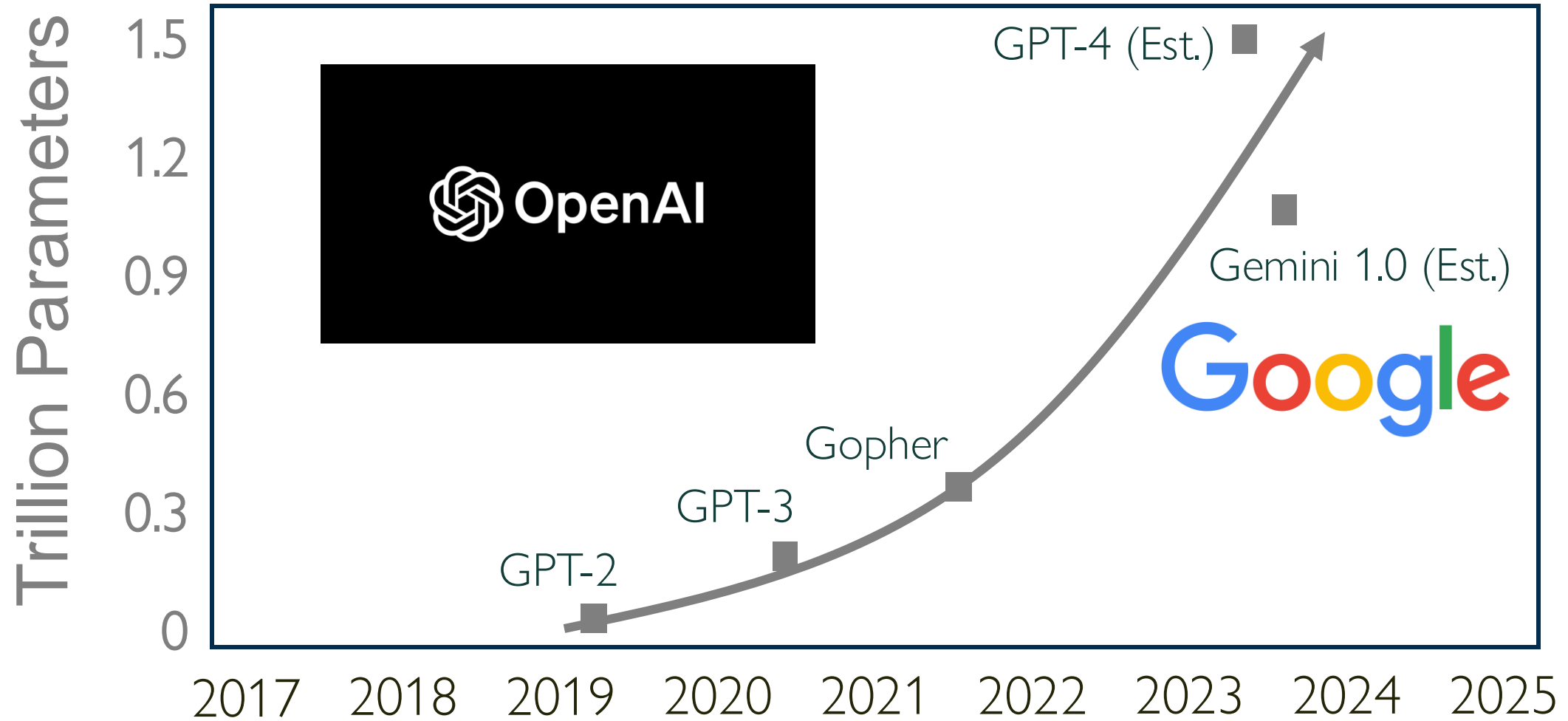
1,200,000,000,000 ops/s  
270 Watts

# LONG LIVE MOORE'S LAW

THE  
END  
OF  
MOORE'S  
LAW

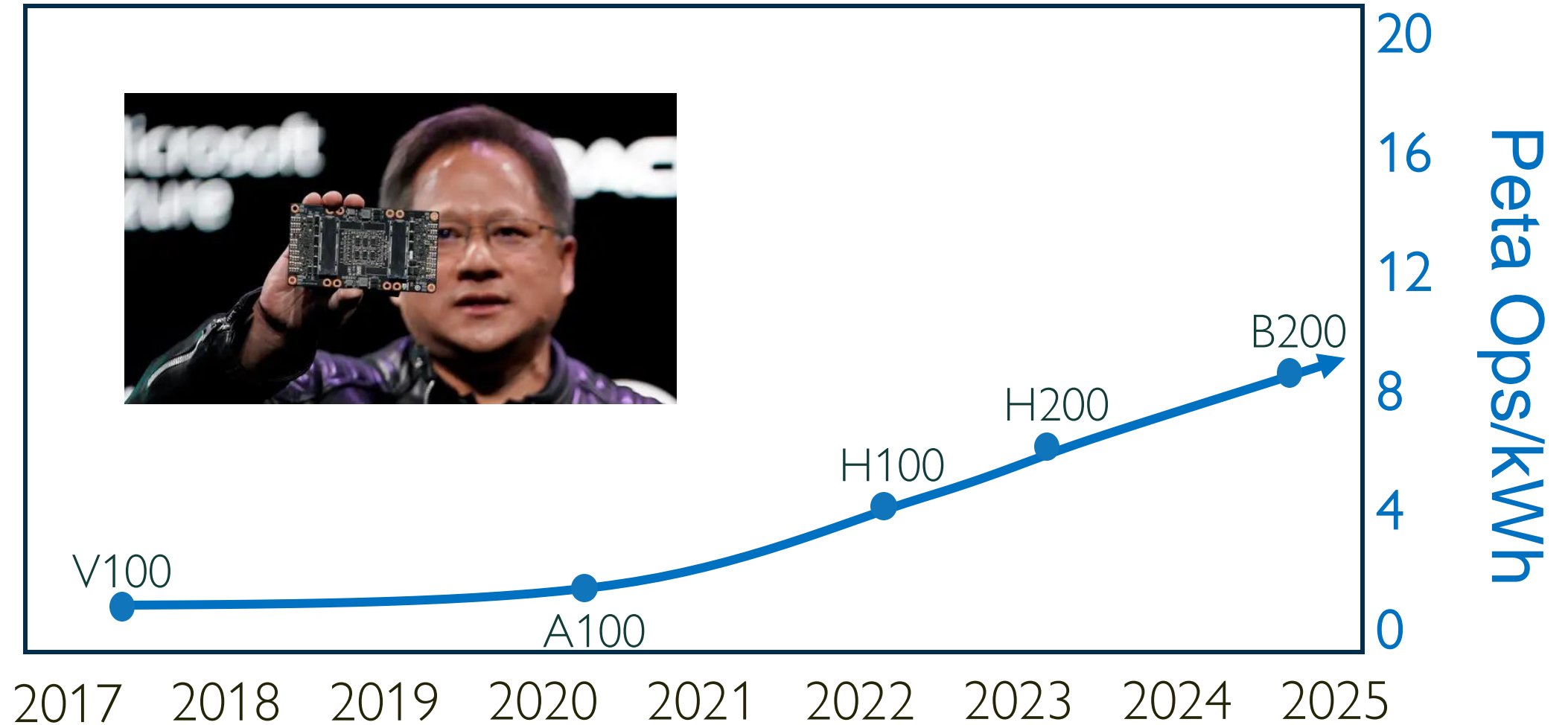


# LLMS' GROWTH

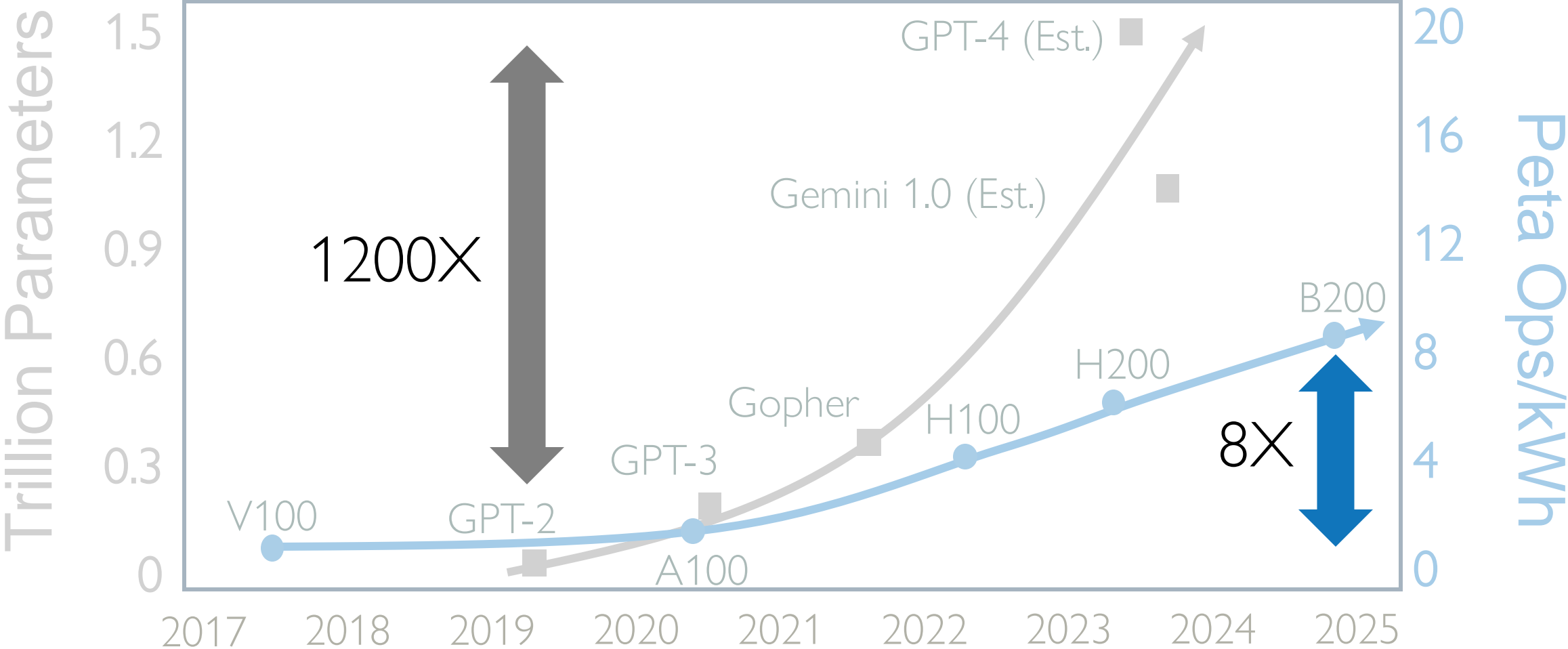




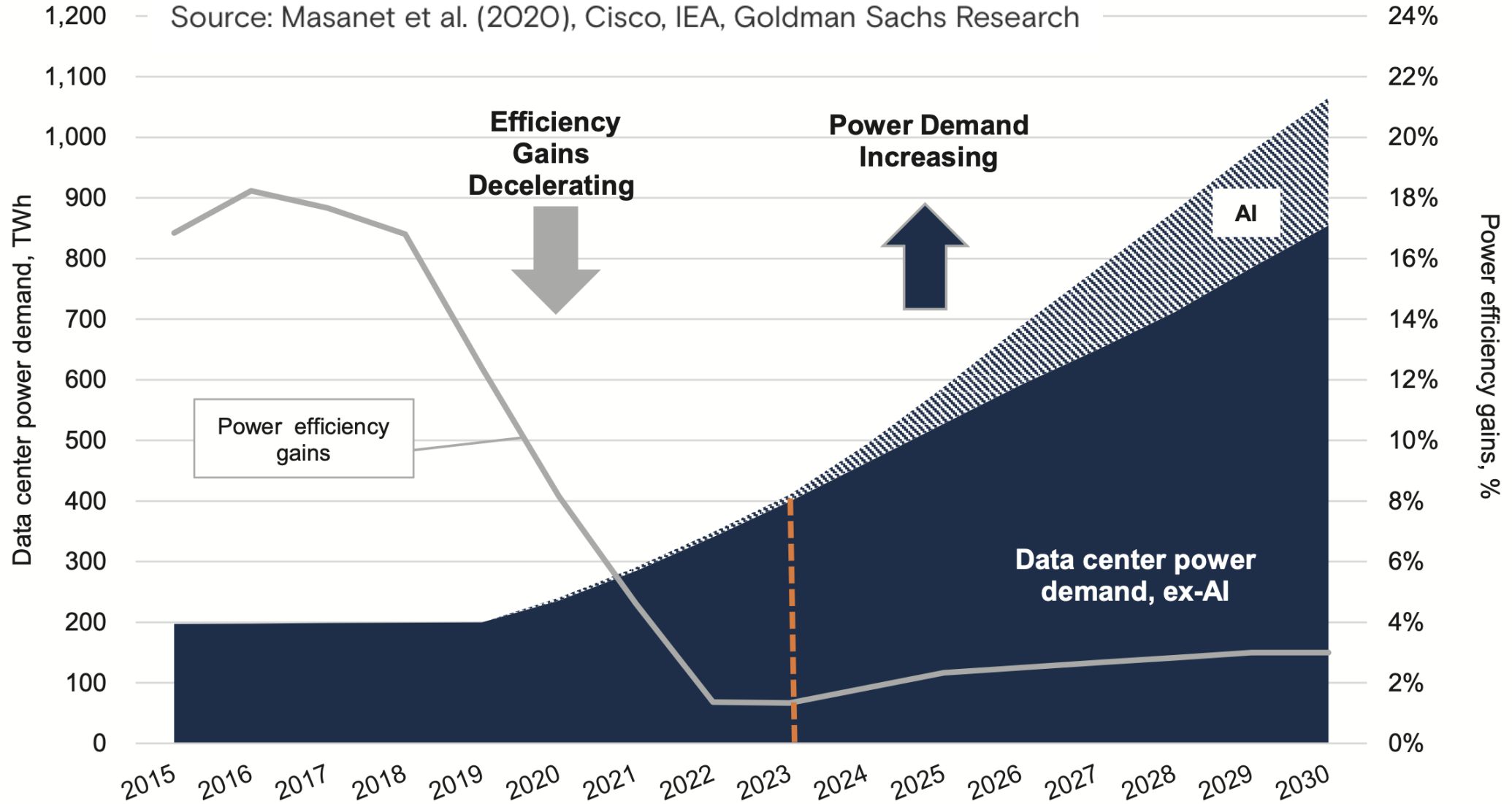
# NVIDIA CHIP EFFICIENCY



# CATCH ME IF YOU CAN!



# GROWTH IN DATACENTER ENERGY



# OPERATIONAL VS. EMBODIED EMISSIONS



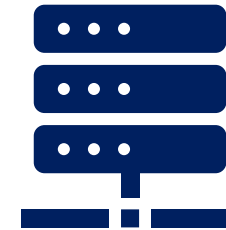
“ The **use stage** GHG emissions in 2020 relating to electricity use account for the **majority of total GHG emissions**. ”

- *Malmodin et al. (2020)*

## OPERATIONAL EMISSIONS

*Scope 1 & Scope 2*

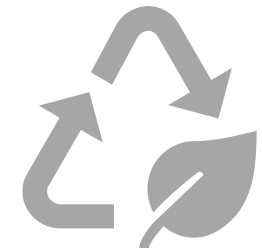
95 million tons CO<sub>2</sub>



## EMBODIED EMISSIONS

*Scope 3*

31 million tons CO<sub>2</sub>



76%

24%



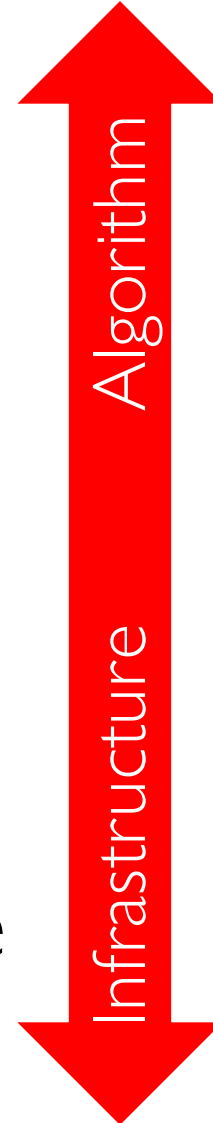
# POST-MOORE DATACENTERS



Design for “ISA”

- Integration
  - reduce data movement
- Specialization
  - cut resources to analyze data
- Approximation
  - compress data & computation

From algorithms to infrastructure





# CENTER @ EPFL SINCE 2011

## Mission

- Sustainable computing
- IT for sustainability
- Best practices, metrics & methodologies

## Impact

- Server-grade ARM CPU
- Cloud-native network/database stacks
- Liquid-cooling from chip to rack

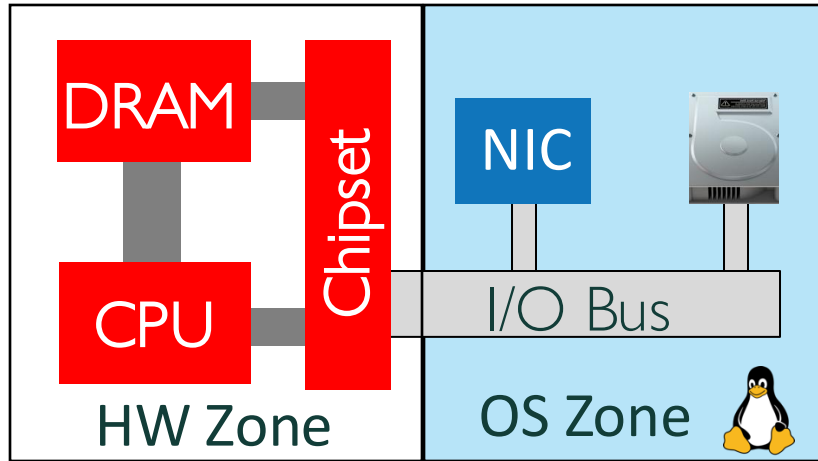


# OUTLINE



- Overview
- Post-Moore Computing
  - Compute infrastructure
  - AI runtime stack
  - Metrics & methodologies
- Summary

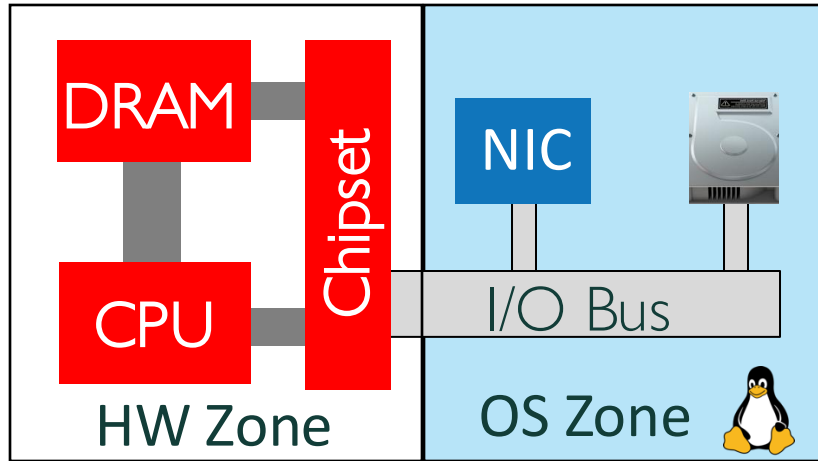
# TODAY'S SERVER = 90'S DESKTOP PC



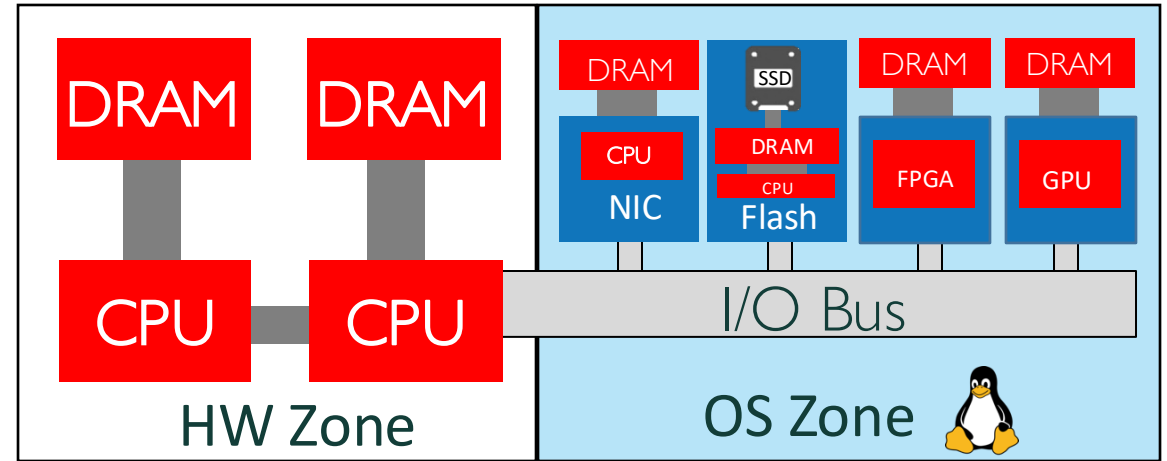
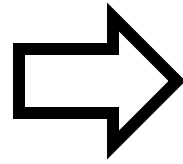
90's Desktop PC



# TODAY'S SERVER = 90'S DESKTOP PC

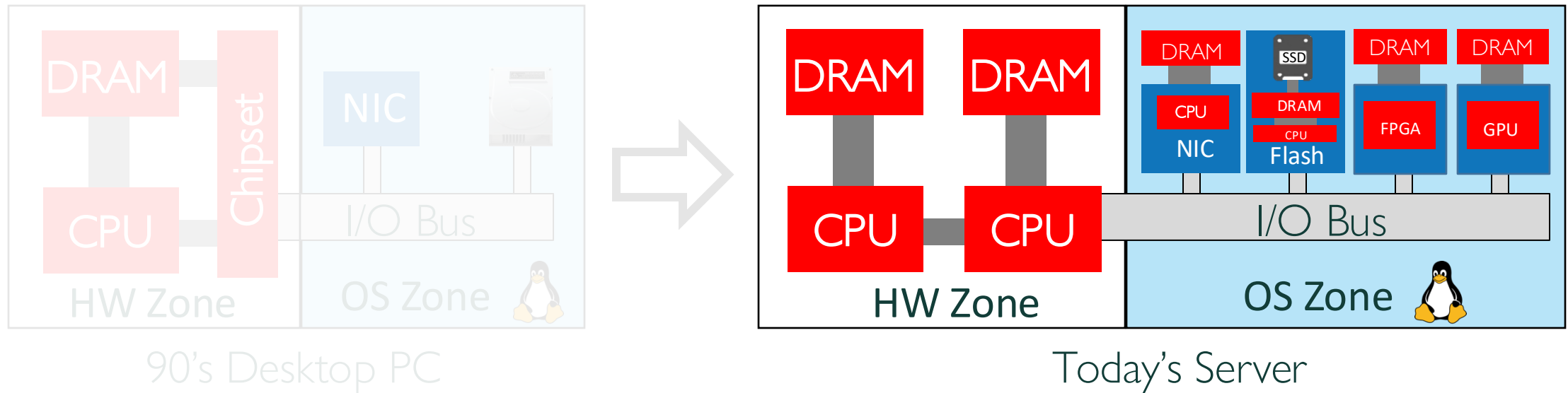


90's Desktop PC



Today's Server

# TODAY'S SERVER = 90'S DESKTOP PC



- Focused on minimizing cost (Google c.a. 2000)
- CPU, memory = **nanosecond** timescale, OS, I/O = **millisecond** timescale
- OS follows legacy interfaces (PCIe) and abstractions (POSIX)
- Silicon fragmented across legacy interfaces

# EFFICIENCY PROBLEMS IN COMPUTING



Hardware/workload mismatch (EPFL, Meta, Google)

Datacenter tax ~ 20% (Google)

- 20,000 threads running per CPU
- Virtualization/containerization
- RPC

Memory wasted (Microsoft)

- 50% of containers do not use their memory
- 20% of memory is stranded

GPU utilization for deep learning < 50% (Microsoft)

# POST-MOORE SERVERS [IEEE Micro'24]



## Server-centric CPU design

- Exploit massive request-level parallelism per service
- Maximize efficiency: throughput/area, throughput/Watt

## Tight integration of CPU, GPU, memory, NIC

- Emerging chip-to-chip standards (UCIe)
- High-bandwidth memory for AI

## Rack-level fabrics

- NVLink, CXL

## Liquid cooling at chip level

- 2x higher power density



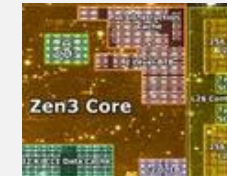
# POST-MOORE CPUS [ISCA'12]



## Today's server CPUs

- ✗ Designed for single-core performance
- ✗ Power-bound  $\rightarrow$   $\frac{1}{2}$  big cores +  $\frac{1}{2}$  memory
- ✗ Run at high frequency  
(power  $\sim$  superlinear w/ performance)

AMD Zen 3  
4.0 mm<sup>2</sup>  
3.7W @ 3 GHz



C	C	C	C
\$	\$	\$	\$
\$	\$	\$	\$
\$	\$	\$	\$
\$	\$	\$	\$
C	C	C	C

## Cloud-native CPUs

- ✓ Custom cores for max area density
- ✓ Higher throughput/Watt at lower frequency
- ✓ Need only memory for per-core working set

ARM N1  
1.4 mm<sup>2</sup>  
0.7W @ 2 GHz



C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
\$	\$	\$	\$	\$	\$	\$	\$
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C
C	C	C	C	C	C	C	C

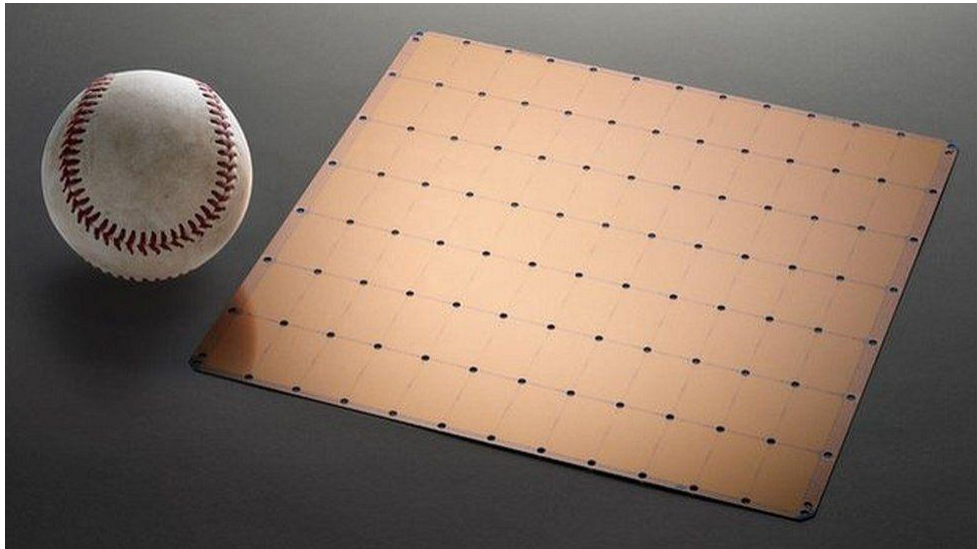
10X higher throughput with SLO!

# AI ACCELERATORS



Inference workloads:

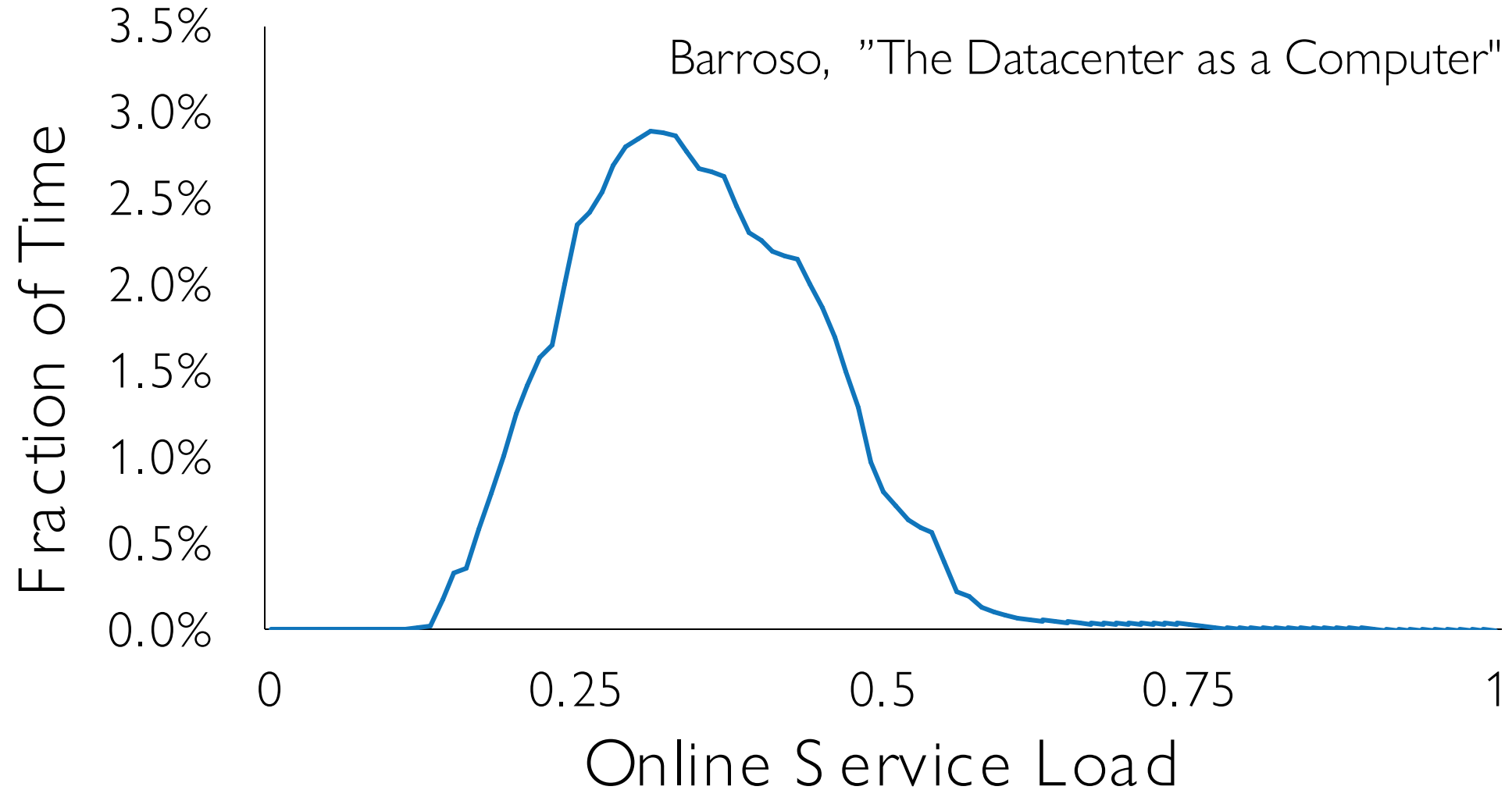
- Online
- Tight latency constraints
- Rely on low-precision arithmetic



Training workloads:

- Offline
- Throughput optimized
- Need high-precision arithmetic

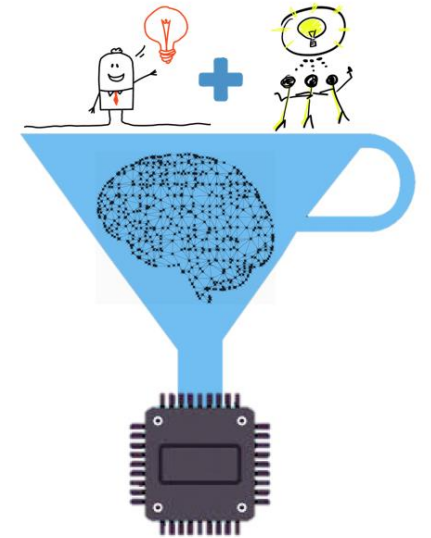
# INFERENCE UTILIZATION



# UNIVERSAL AI ACCELERATORS



- Scaled numeric formats (HBFP [NeurIPS'18], MX)
  - Quantize while maintaining accuracy
  - Use for both for inference and training
  - Inference → helps with outliers
  - Training → don't need FP precision for dot products
  - Work well with sparsification [Harma, ICLR'25]
  - 4-Bit training for transformers [Harma, arXiv'24]



[parsa.epfl.ch/coltrain](https://parsa.epfl.ch/coltrain)

- Accelerators with prioritized schedulers [Drumond, MICRO'21]
  - Piggy-back fine-tuning jobs on an inference accelerator



# AI RUNTIME



- Enhance utilization [Gao, ICSE'24]
  - Proper batching
  - Overlap CPU-centric tasks
  - Hide data transfer between CPU/GPU
  - Minimize or overlap communication
- Need elasticity (Ana's talk)
  - Today's runtimes are too static (container-based deployment)
  - Allow software to scale up/down GPU, memory, network resources

# METRICS & METHODOLOGIES



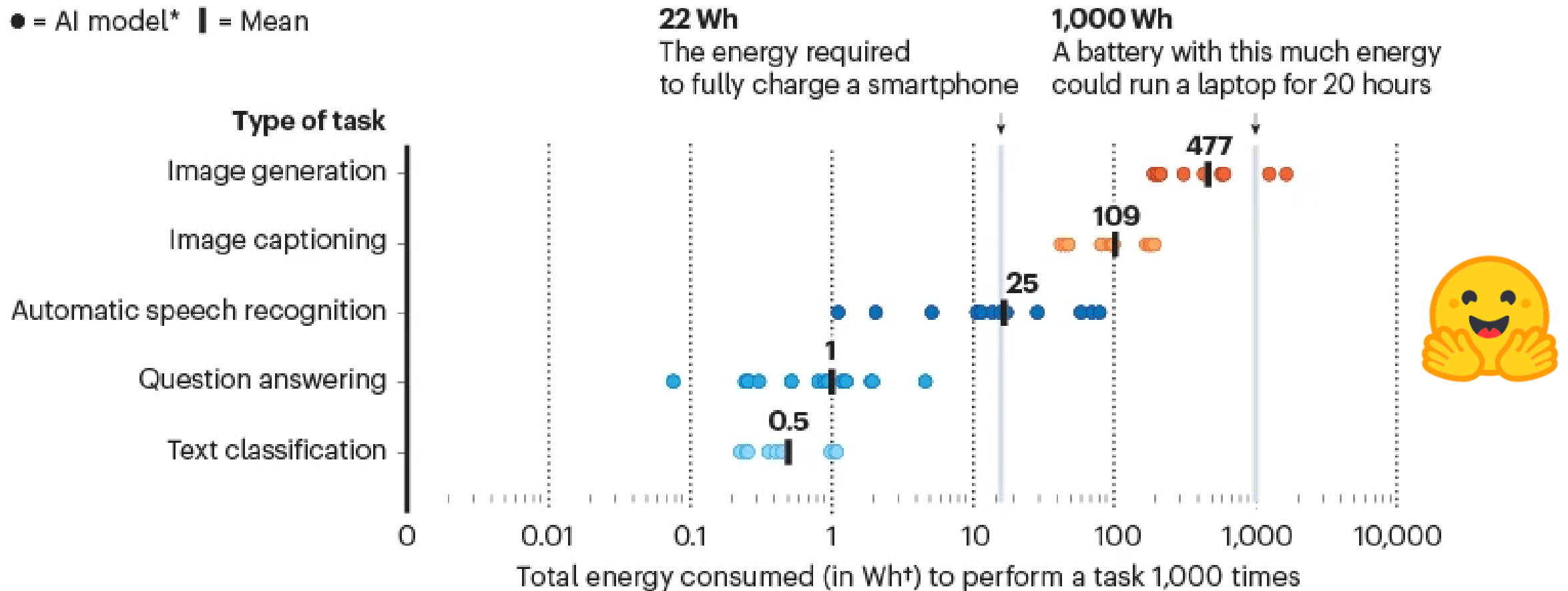
- Design metrics
  - Post-Moore means more accelerators
  - Take into account both operational & embodied emissions
- Operational metrics
  - How much energy do we need for training/inference?
- Operational methodologies
  - How do we measure/monitor our efficiency?



# AI SUSTAINABILITY CLASSIFICATION



● = AI model\* | = Mean



\*Tests conducted on 20 popular open-source models. Each dot represents one model;

†1 Watt-hour represents power consumption of 1 W extended over 1 hour.

©nature

# MEASURE FULL-STACK EFFICIENCY

## DC EFFICIENCY

- electricity w/ renewables, cooling, heat recycling

## IT EFFICIENCY

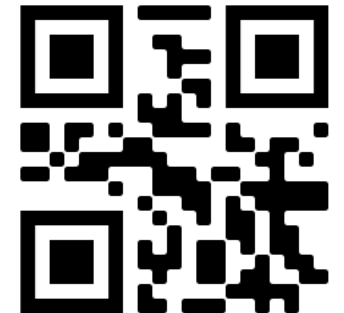
- compute, storage, network and workloads

## CARBON FOOTPRINT

- emissions from input electricity sources



sdea.ch



# SUMMARY



AI's energy requirements grow exponentially

Moore's Law of silicon is dead

Need post-Moore technologies, metrics, best practices

Post-Moore computing:

- Integration + Specialization + Approximation

# THANK YOU!



For more information, please visit us at  
[parsa.epfl.ch](http://parsa.epfl.ch)

# EPFL