



# BEYOND THE AI ENERGY WALL OPTIMAL SERVER DESIGN & OPERATION

**EPFL**

**BABAK FALSAFI**

PROFESSOR, EPFL

PRESIDENT, SWISS DATACENTER EFFICIENCY ASSOCIATION



BEYOND THE AI ENERGY WALL  
~~OPTIMAL SERVER DESIGN & OPERATION~~

**SINGLE-CORE PERFORMANCE WON'T  
BREAK THE WALL**



**BABAK FALSAFI**

PROFESSOR, EPFL

PRESIDENT, SWISS DATACENTER EFFICIENCY ASSOCIATION

# OUR DIGITAL UNIVERSE

DATA  
CENTER

EFFICIENCY



Fueled by:

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

# DATA CENTERS ARE THE BACKBONE

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



2.4 km

200 m

Boydton, VA (300 MW)

# CLOUDS TURN DATA INTO VALUE / SERVICE

DATA CENTER

EFFICIENCY

Edge Cloud

Enterprise Cloud

Public Cloud

Users & Devices



Temporal/Sensitive/Local Data

Persistent/Global Data

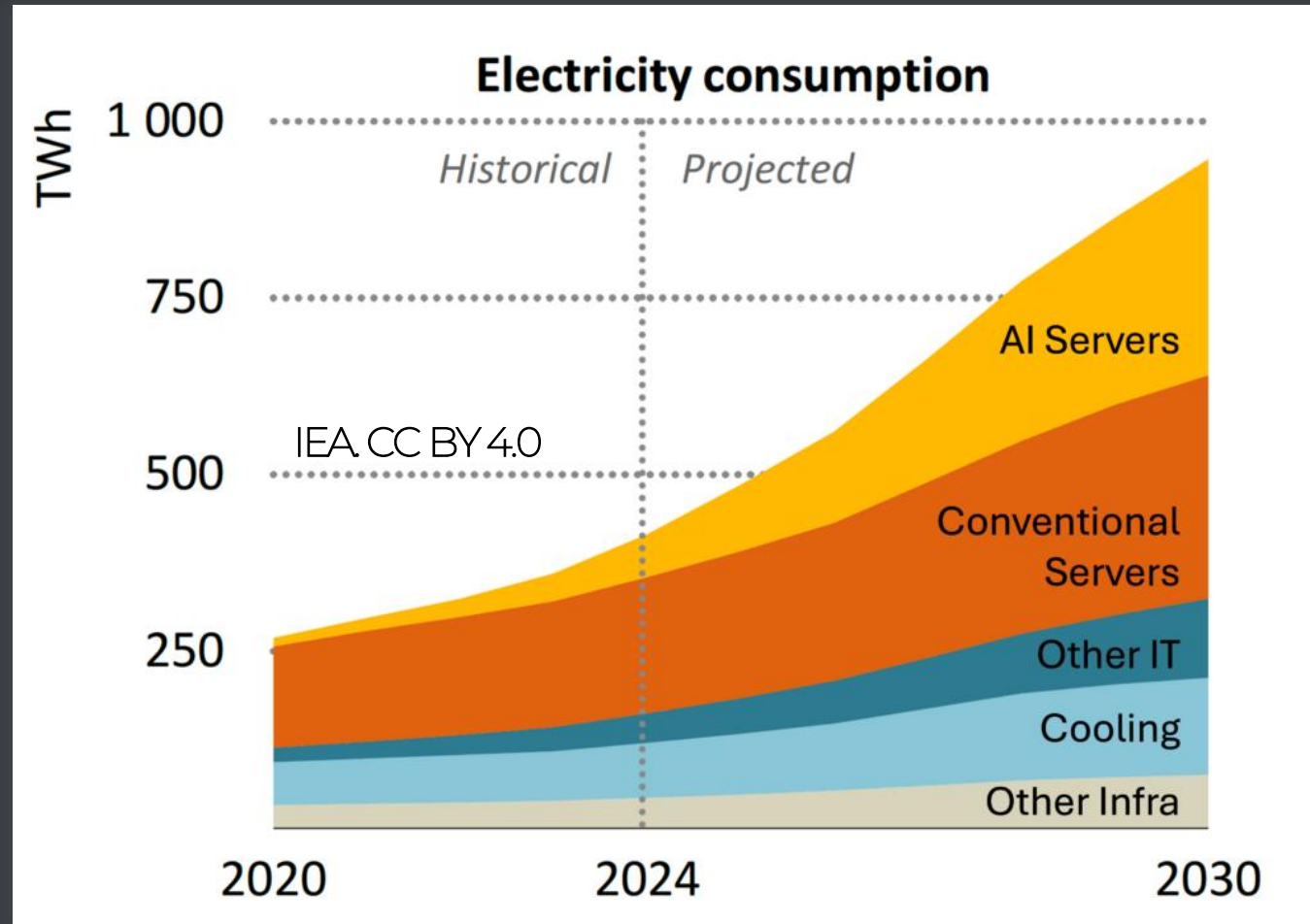


# CLOUD IS AN ECONOMIC MODEL





# ENERGY GROWTH PROJECTIONS



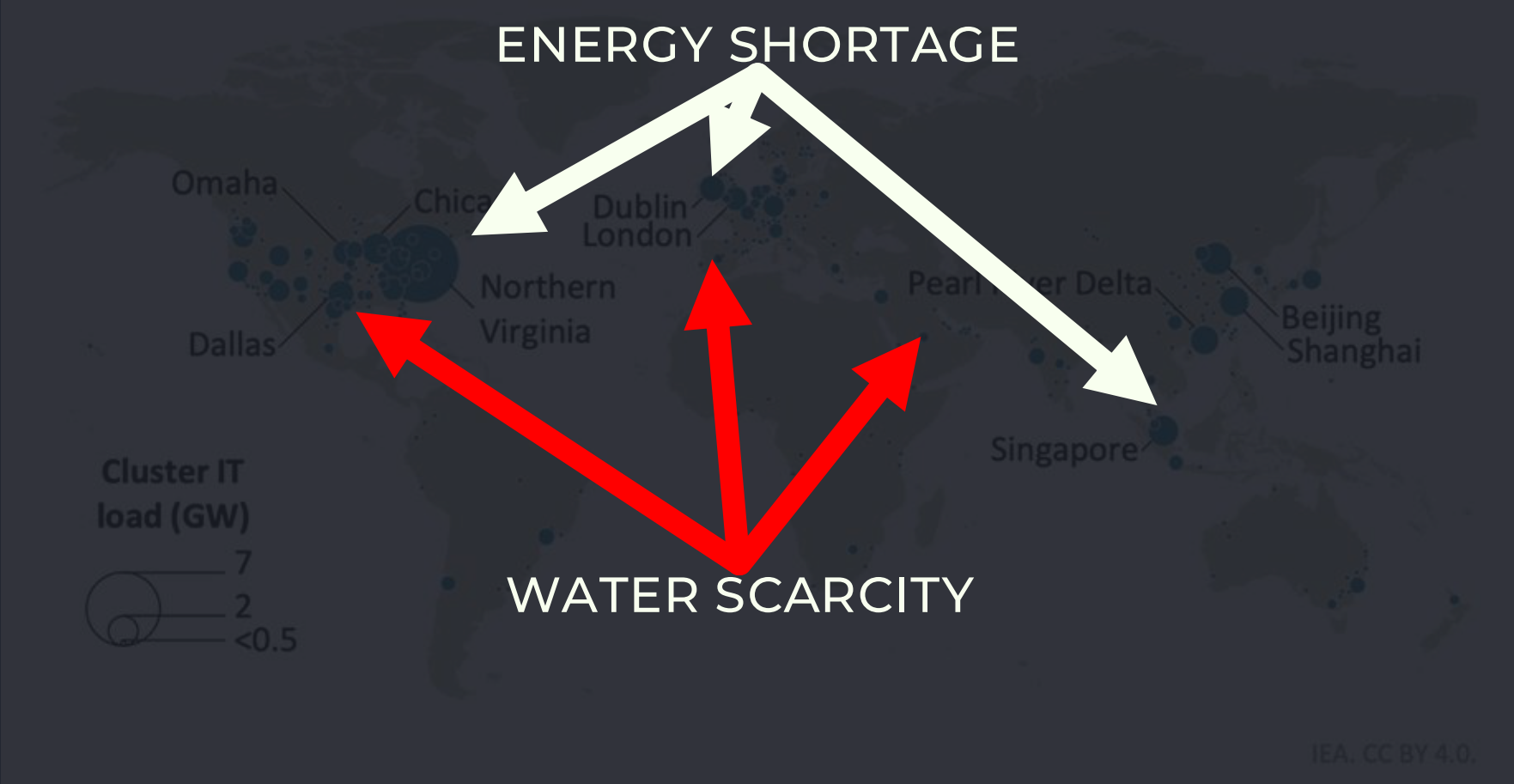
# THE SKEWED DEMOGRAPHICS OF DATACENTERS

DATA CENTER

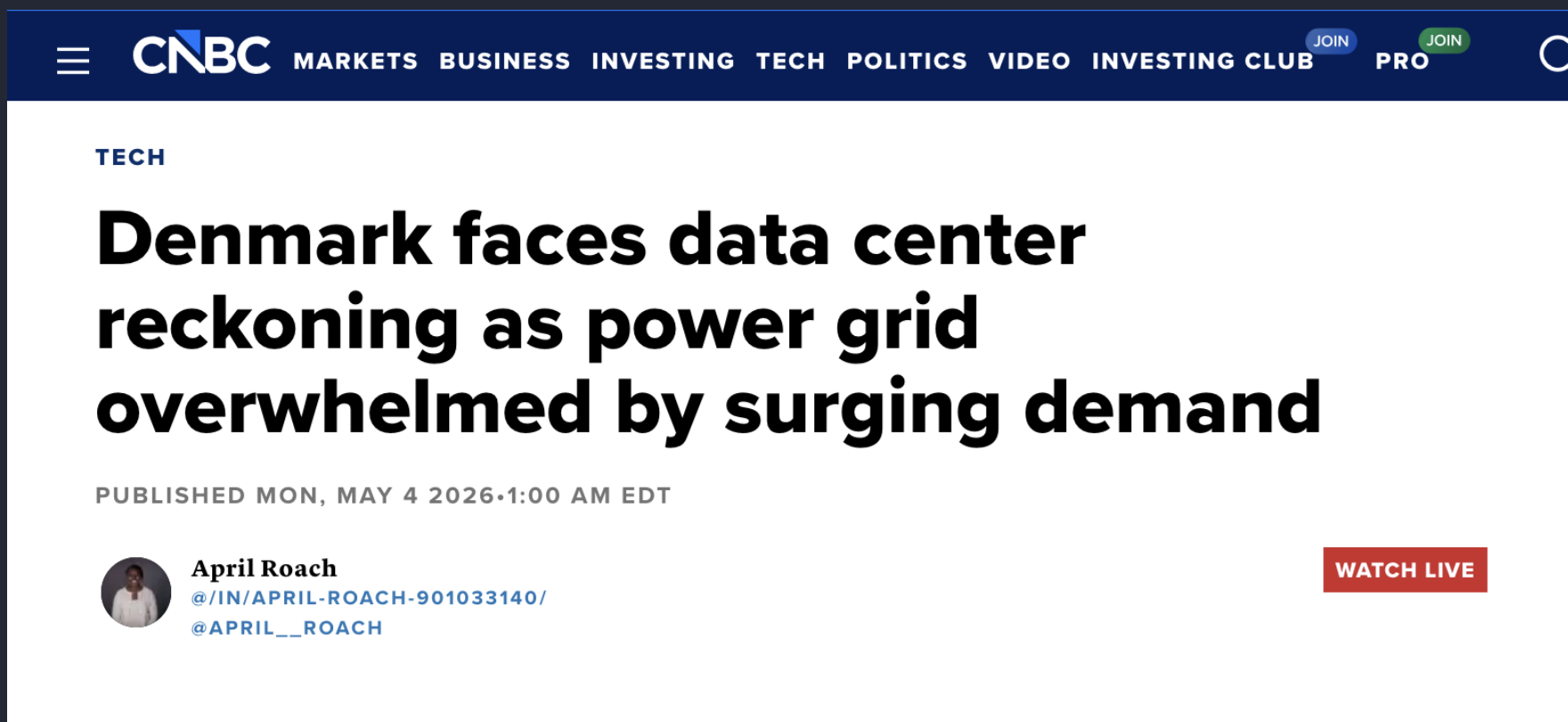
EFFICIENCY



# THE SKEWED DEMOGRAPHICS OF DATACENTERS



# NOT JUST A PROBLEM FOR GRIDS IN US A CASE IN POINT



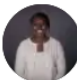
The image shows a screenshot of a CNBC news article. The top navigation bar includes the CNBC logo and various menu items: MARKETS, BUSINESS, INVESTING, TECH, POLITICS, VIDEO, INVESTING CLUB, and PRO. There are also 'JOIN' buttons next to 'INVESTING CLUB' and 'PRO'. A search icon is visible on the right. The article is categorized under 'TECH'. The main headline reads 'Denmark faces data center reckoning as power grid overwhelmed by surging demand'. Below the headline, it says 'PUBLISHED MON, MAY 4 2026•1:00 AM EDT'. The author is identified as April Roach, with her profile picture and social media handles: @/IN/APRIL-ROACH-901033140/ and @APRIL\_\_ROACH. A red 'WATCH LIVE' button is located in the bottom right corner of the article preview.

≡ **CNBC** MARKETS BUSINESS INVESTING TECH POLITICS VIDEO INVESTING CLUB JOIN PRO JOIN 🔍

**TECH**

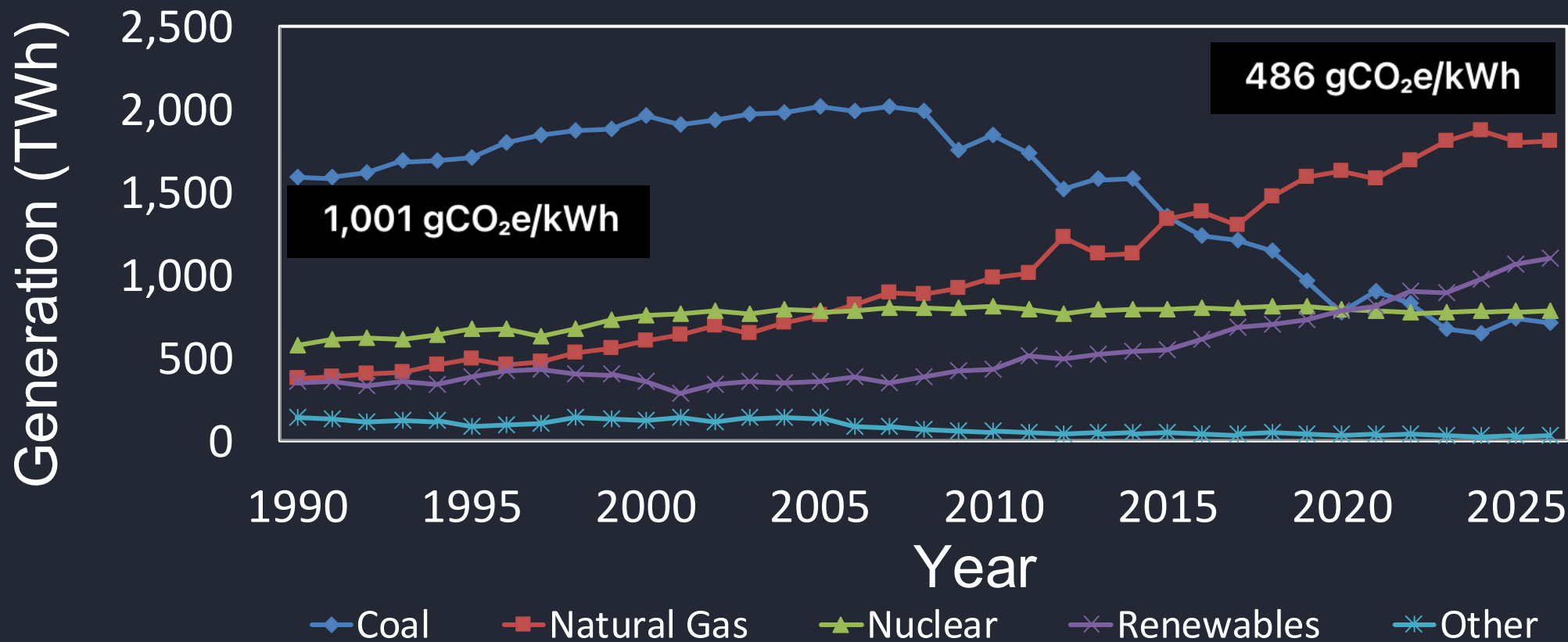
## Denmark faces data center reckoning as power grid overwhelmed by surging demand

PUBLISHED MON, MAY 4 2026•1:00 AM EDT

 **April Roach**  
[@/IN/APRIL-ROACH-901033140/](#)  
[@APRIL\\_\\_ROACH](#)

**WATCH LIVE**

# ENERGY INFORMATION ADMINISTRATION US ELECTRICITY GENERATION (LCA)



# EMISSIONS DOMINATED BY OPERATION

DATA  
CENTER

EFFICIENCY

EMBODIED  
EMISSIONS  
Scope 3



35%

OPERATIONAL  
EMISSIONS  
Scope 1 & Scope 2

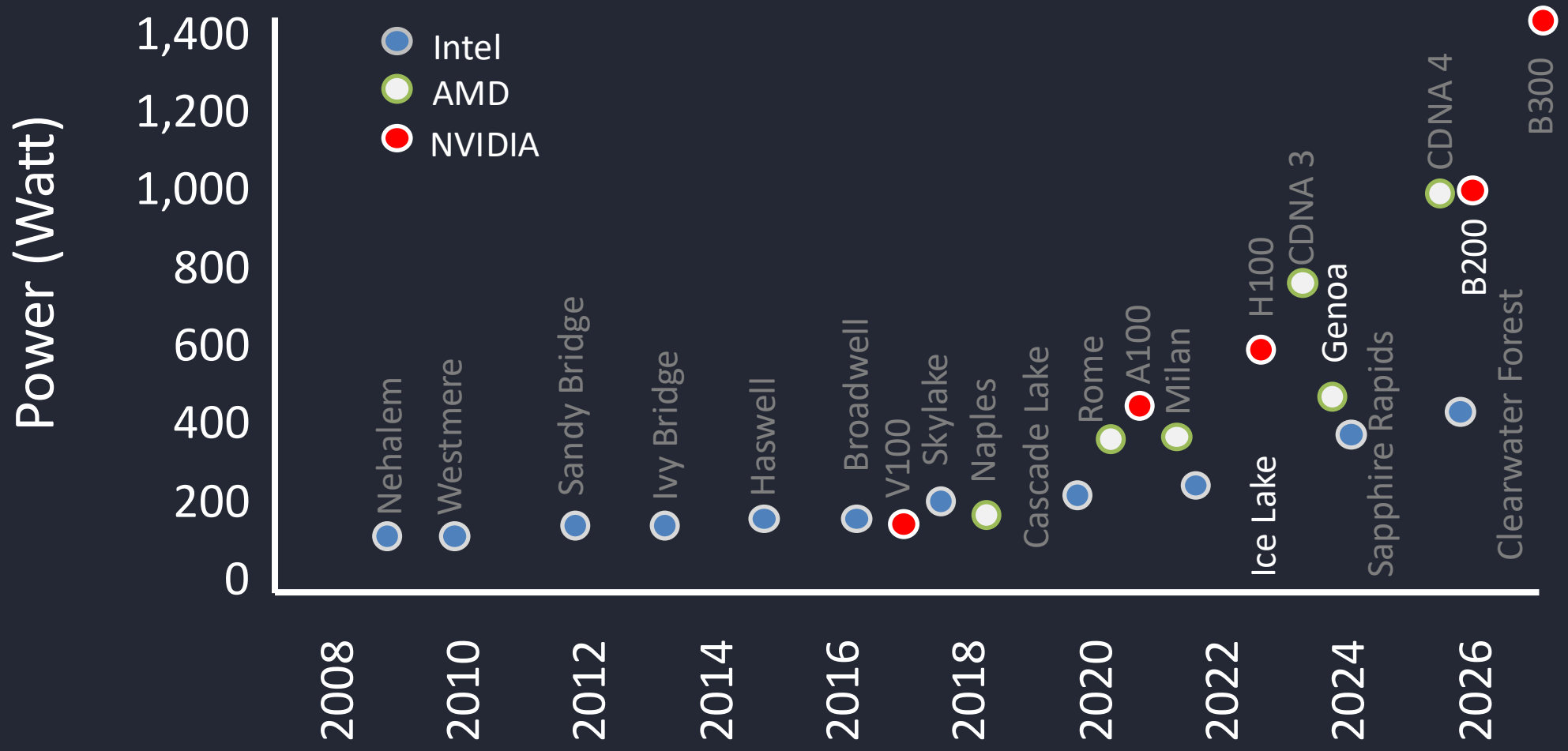


65%

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

© Schneider Electric (2023)

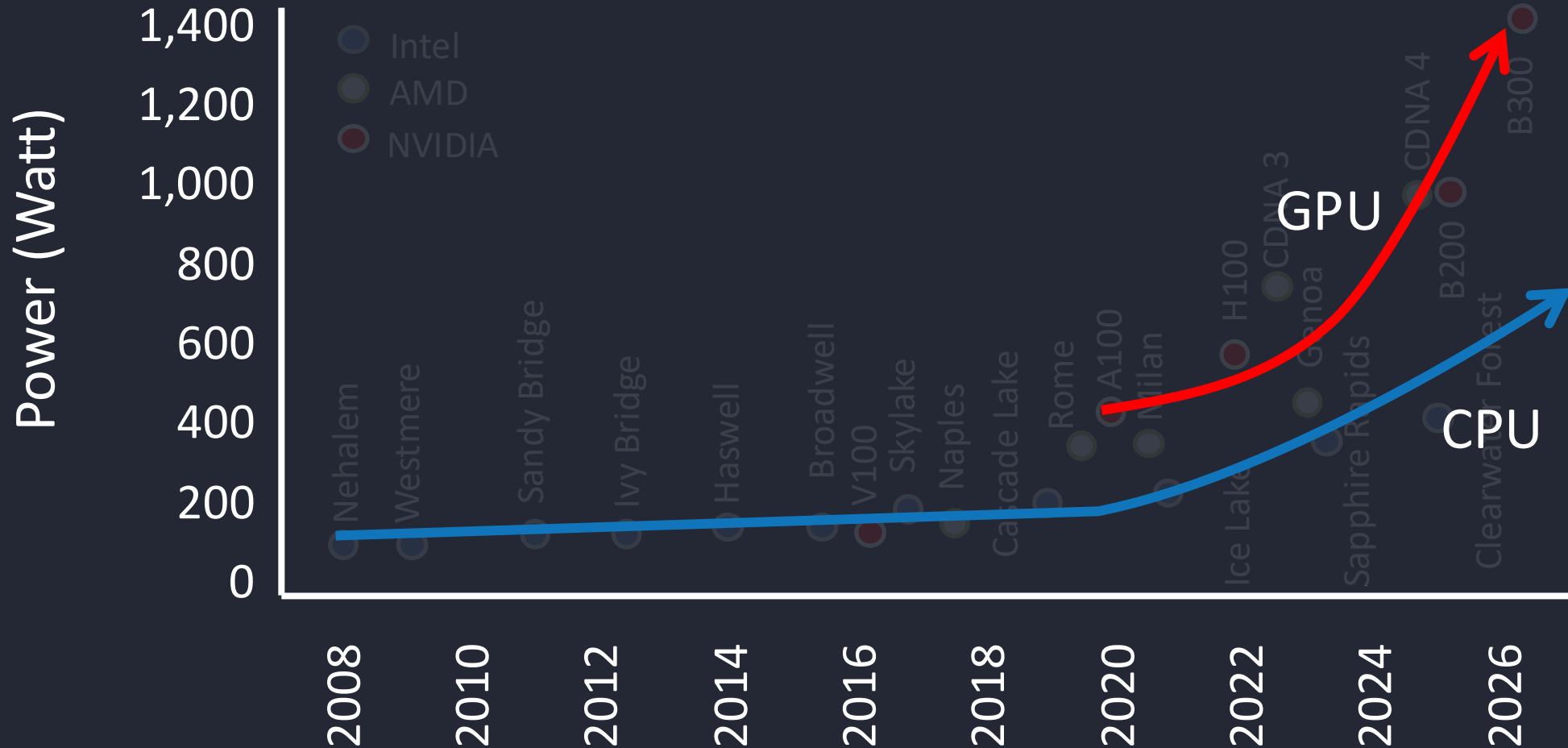
# OPERATIONAL EMISSIONS GOING UP CHIP THERMAL DESIGN POWER



# OPERATIONAL EMISSIONS GOING UP CHIP THERMAL DESIGN POWER

DATA CENTER

EFFICIENCY



# JEVONS PARADOX REBOUND W/ CAPPED ENERGY

Globally 2500 GW are stuck in queues [Forbes, June'26]

- Operators are energy constrained
- Efficiency → more compute for a given energy budget (e.g., 1GW facility)
- As AI demand goes up, the rebound for the compute market is mediated through grids and energy policy



# AI ALSO HAS PROS

DATA  
CENTER

EFFICIENCY

- Helps optimize **energy exploration**, supply and consumption
- **Balances management** of electricity networks
- Application in transport to **save energy** and cost
- Can optimize **energy management** in buildings
- Is a powerful tool for **scientific discovery**
- ...

# AI ALSO HAS PROS

DATA  
CENTER

EFFICIENCY

AlphaEvolve changed workload distribution to save 0.7% of Google's worldwide computing resources.

May 2025

# OPTIMAL DESIGN & OPERATION

DATA  
CENTER

EFFICIENCY

Metrics



Design



Best Practices



# OPTIMAL DESIGN & OPERATION

DATA  
CENTER

EFFICIENCY

Metrics



Design



Best Practices



# THE INDUSTRY STANDARD

DATA  
CENTER

EFFICIENCY

$$\text{PUE} = \frac{\text{Total DC Power}}{\text{IT Power}}$$

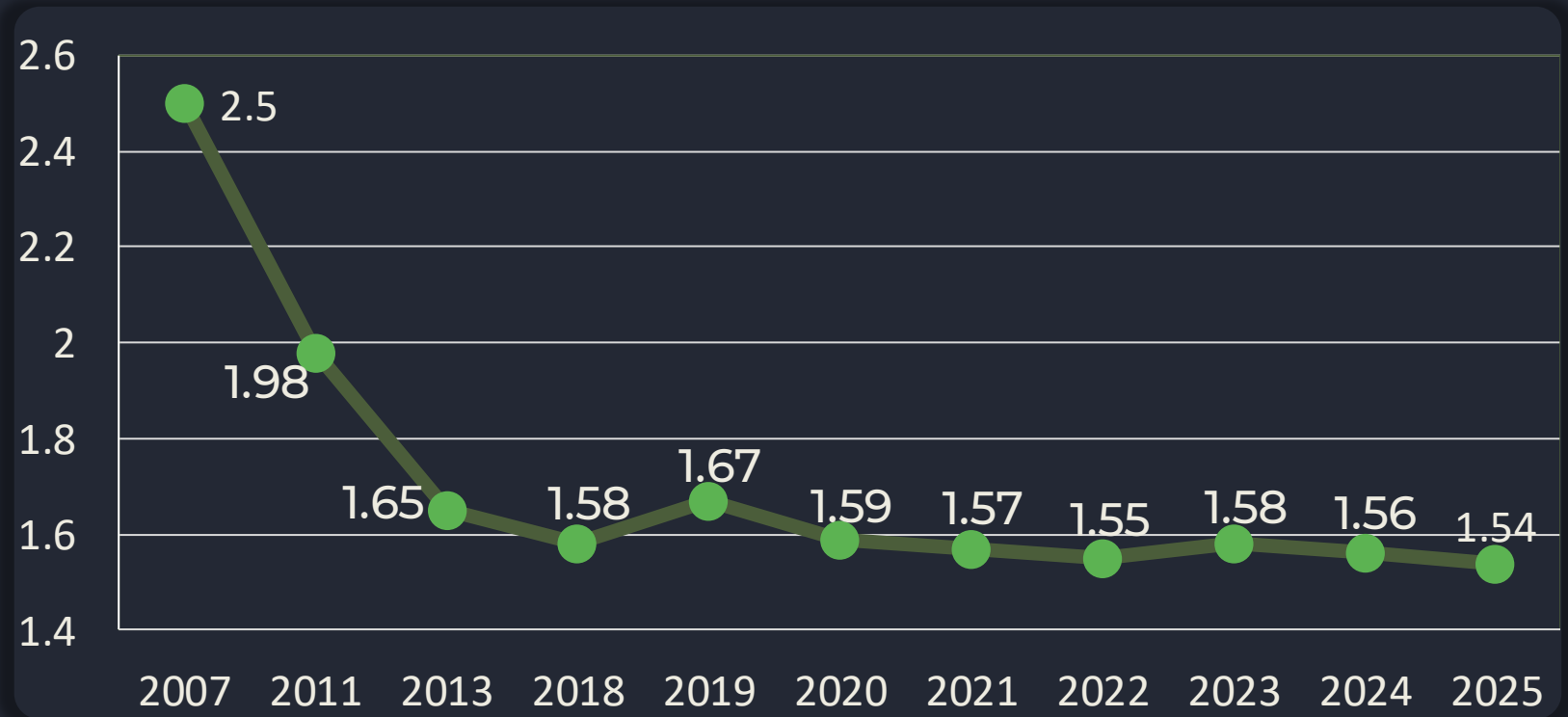


- PUE has been around for two decades
- Easy to calculate, industry-wide adoption, benchmarking
- Global Average (2025): 1.54 (= 65% of the electricity flows into IT)

# PUE HAS HIT A WALL

DATA  
CENTER

EFFICIENCY



- Hyperscalers are  $\approx 1.20$
- Co-locators are  $\approx 1.40$

# PUE SAYS LITTLE ABOUT EFFICIENCY

DATA  
CENTER

EFFICIENCY



## **NO IT EFFICIENCY**

Inefficient or underutilized servers make PUE look good



## **NO END-TO-END ENERGY FLOW**

Ignores heat recovery or on-premise renewables



## **NO ENVIRONMENTAL IMPACT**

PUE ignores the energy source and water use

# IT INFRA IS THE BIGGEST ENERGY GUZZLER

DATA CENTER

EFFICIENCY

**PUE = 1.54**

Global Average



65% of the energy flows into IT

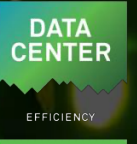
**PUE = 1.20**

Efficiency Leaders

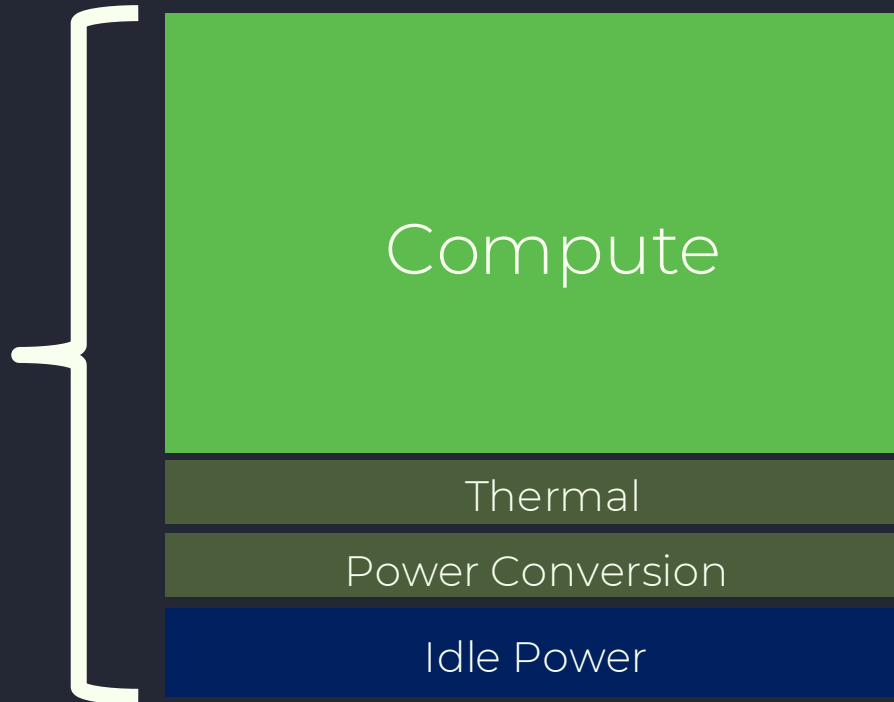


83% of the energy flows into IT

# IT INFRASTRUCTURE: WHERE DOES ENERGY GO?



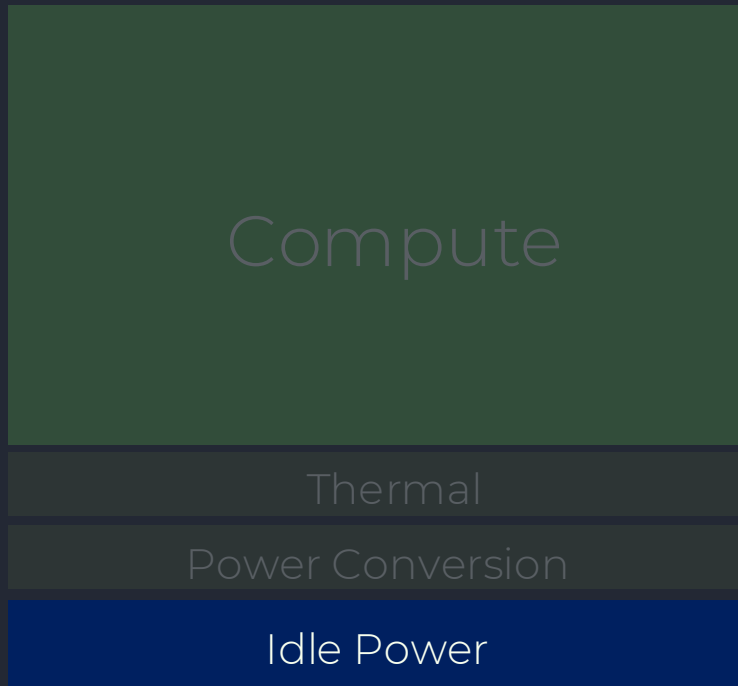
20  
kW



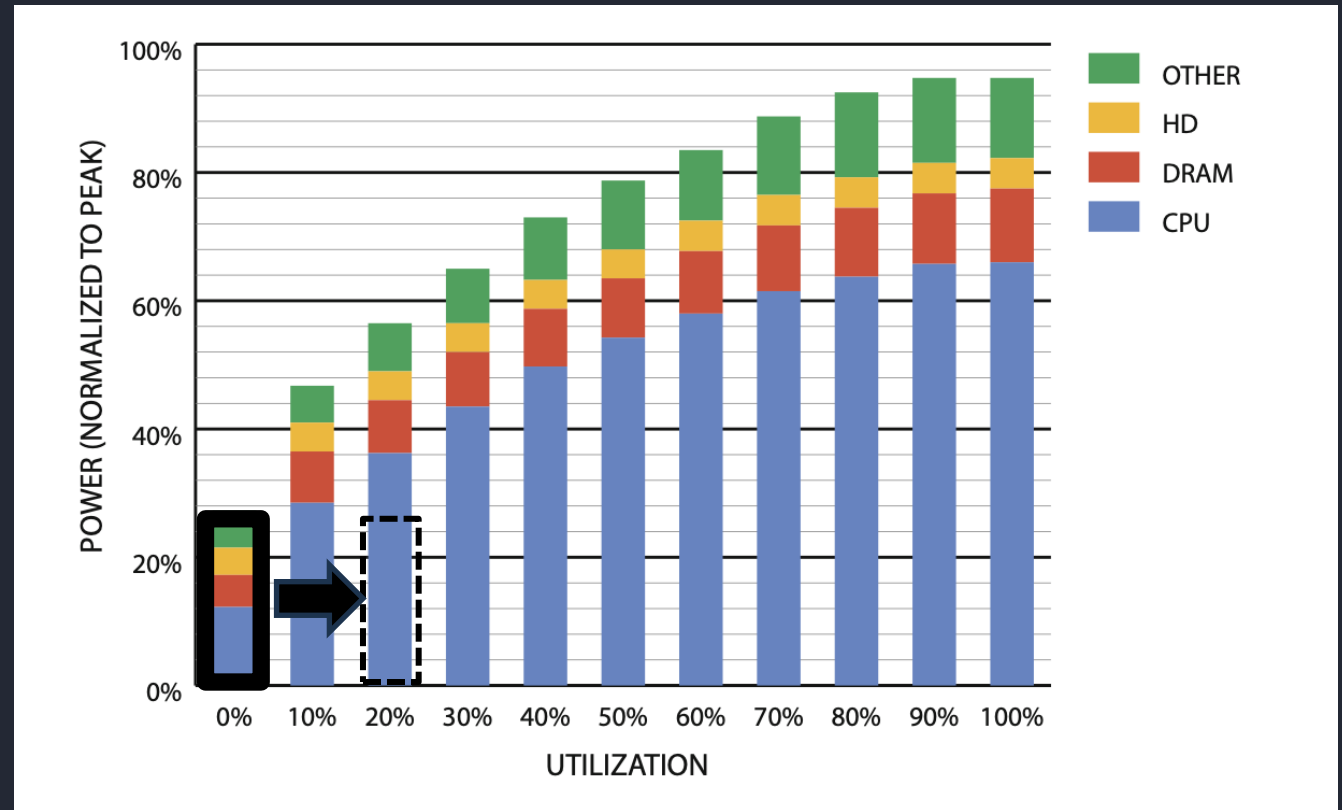
IT

Non-IT

# IT INFRASTRUCTURE ENTERPRISE CUSTOMERS



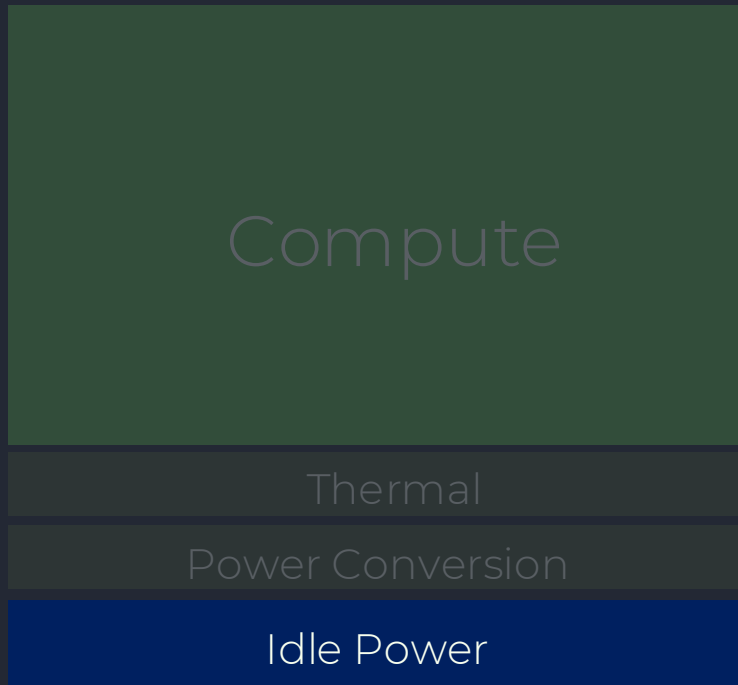
At 20% utilization, Idle Power  $\approx$  half of IT Power



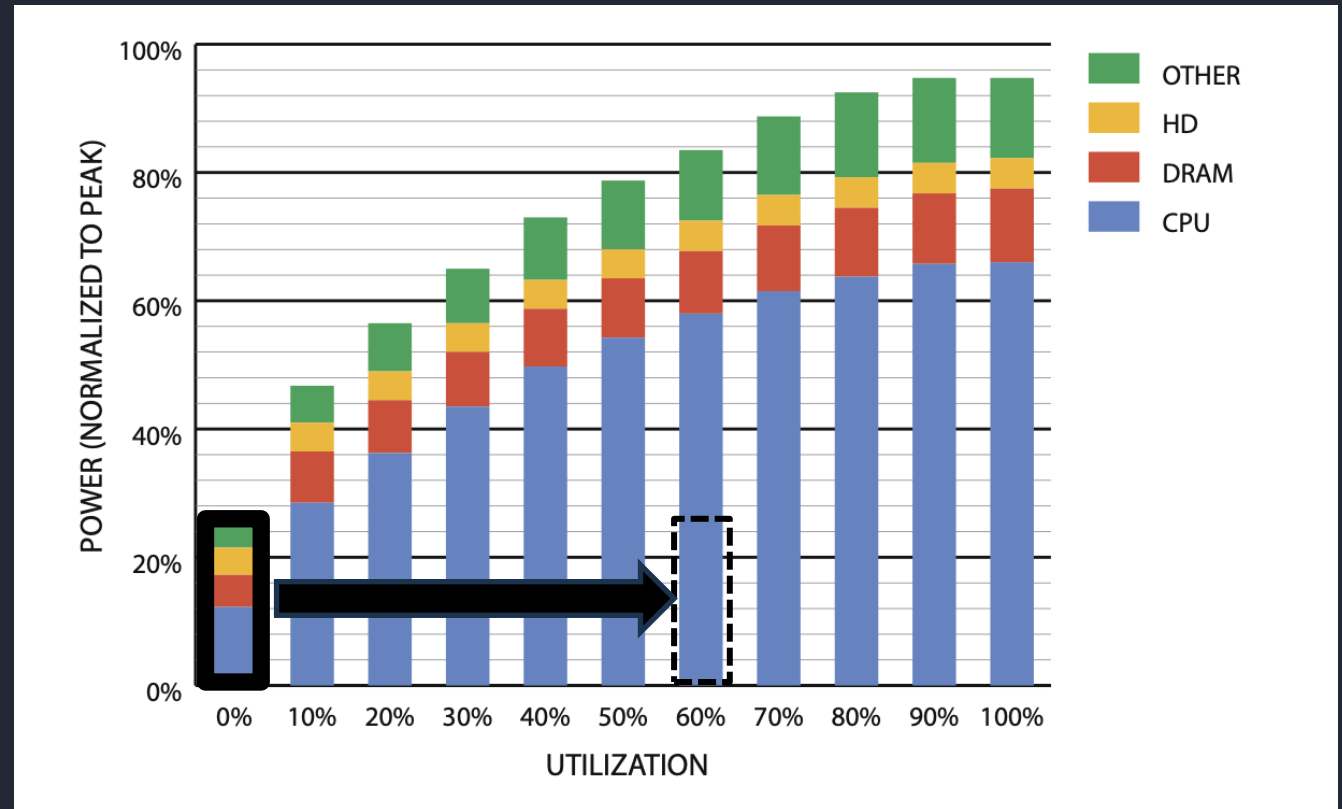
# IT INFRASTRUCTURE HYPERSCALERS

DATA  
CENTER

EFFICIENCY



At 60% utilization, Idle Power  $\approx$  30% of IT Power



# MEASURING POWER & THERMALS

DATA  
CENTER

EFFICIENCY



- Fan power
  - Cubic in fan speed
  - Problem when operating high temp.
  - Shouldn't be part of IT (liquid cooling)
- Power conversion
  - PSU: > 90% efficiency for the cloud
  - Lower efficiency at higher VDC

# MEUSURING COMPUTE EFFICIENCY

DATA  
CENTER

EFFICIENCY

Compute

Thermal

Power Conversion

Idle Power

- Hardware
  - Throughput/W of electricity
  - Throughput/mm<sup>2</sup> of silicon
  - SLO
- End-to-end
  - Need KPIs per workload
  - Output/W, Output/mm<sup>2</sup>

# OPTIMAL DESIGN & OPERATION

DATA  
CENTER

EFFICIENCY

Metrics



Design



Best Practices



# POST-MOORE DATACENTERS DESIGN FOR “ISA”

DATA  
CENTER

EFFICIENCY

- 1. Integration**
  - reduce data movement
- 2. Specialization**
  - cut resources to analyze data
- 3. Approximation**
  - compress data & computation

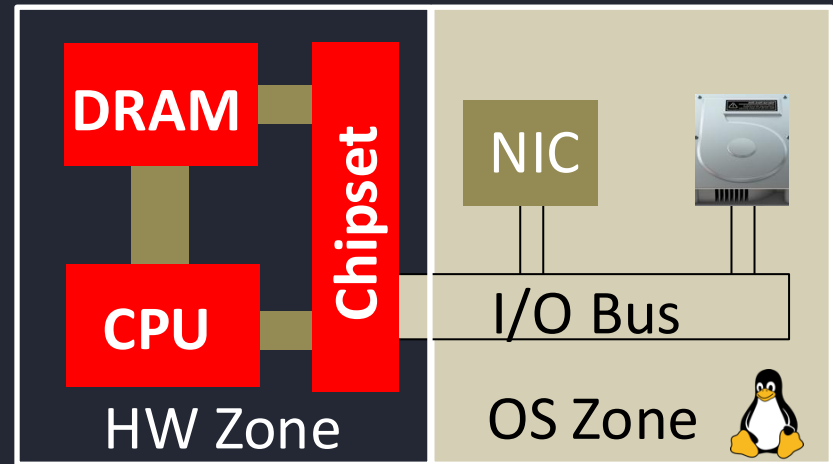
From algorithms to infrastructure



# TODAY'S SERVER ARCHITECTURE BASED ON 90s DESKTOPS

DATA  
CENTER

EFFICIENCY



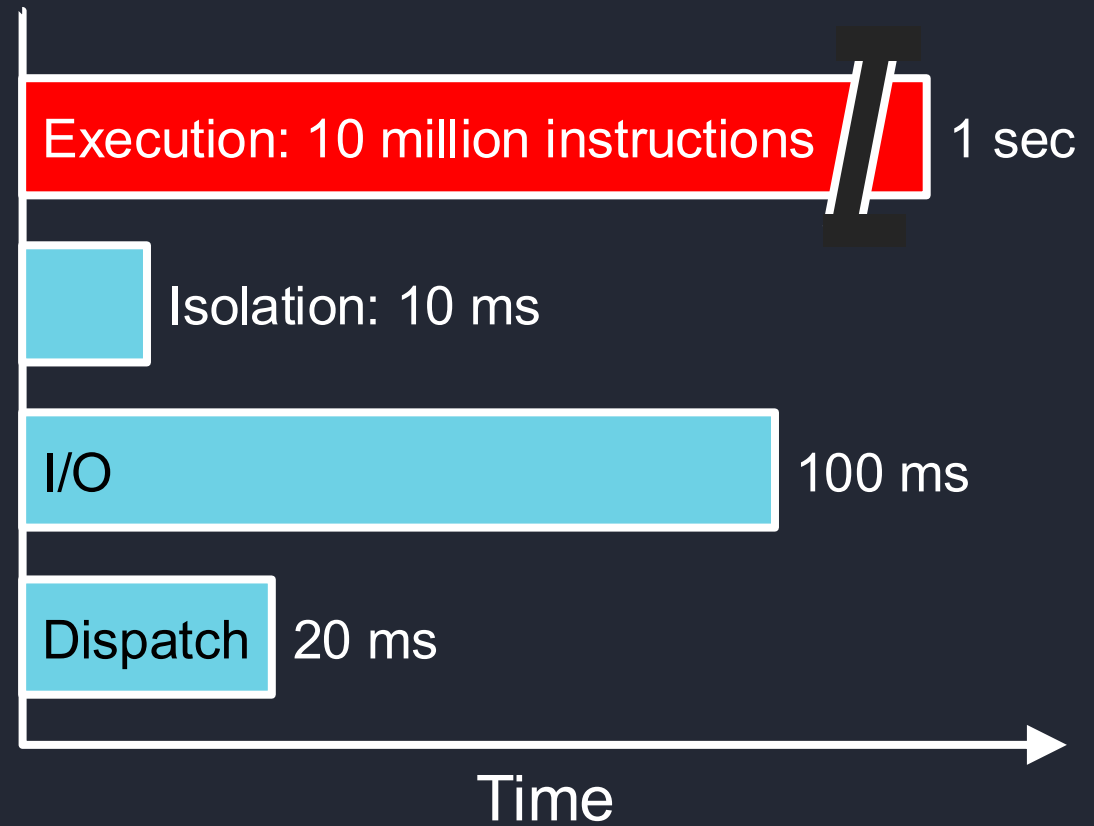
90s' Desktop PC

# 90S DESKTOP APPLICATIONS

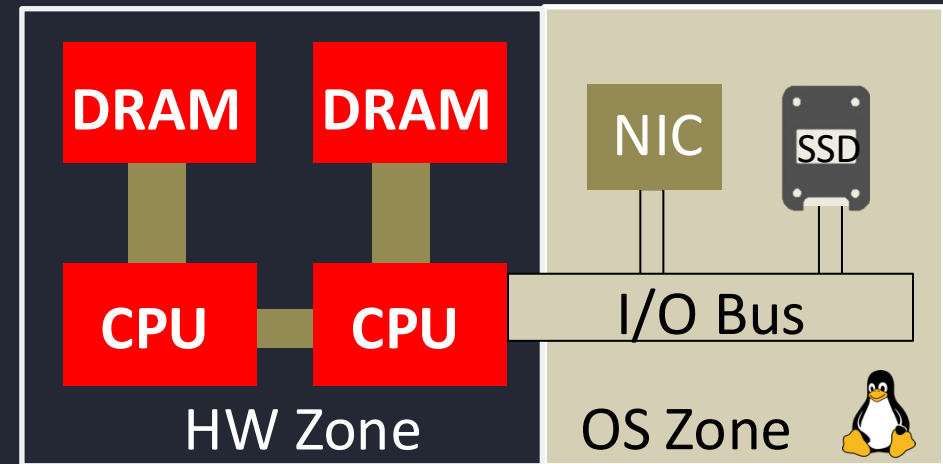


- User-centric
- Monoliths
- Little sharing
- Long-running

# OS ACTIVITY AMORTIZED



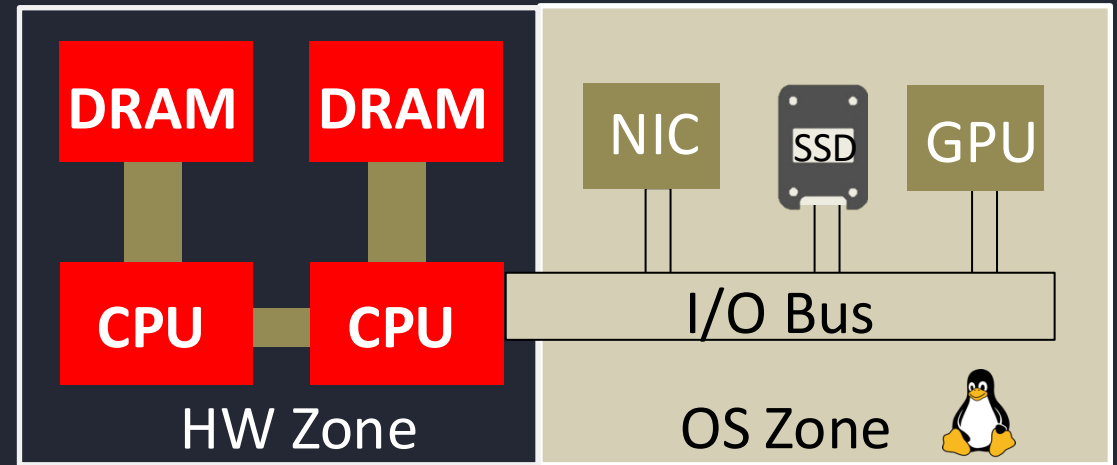
# TODAY'S SERVER ARCHITECTURE BASED ON 90s DESKTOPS



Today's Server

# TODAY'S SERVER ARCHITECTURE BASED ON 90s DESKTOPS

AI ○ ○ ○



Today's AI Inference Server

# TODAY'S SERVER APPLICATIONS

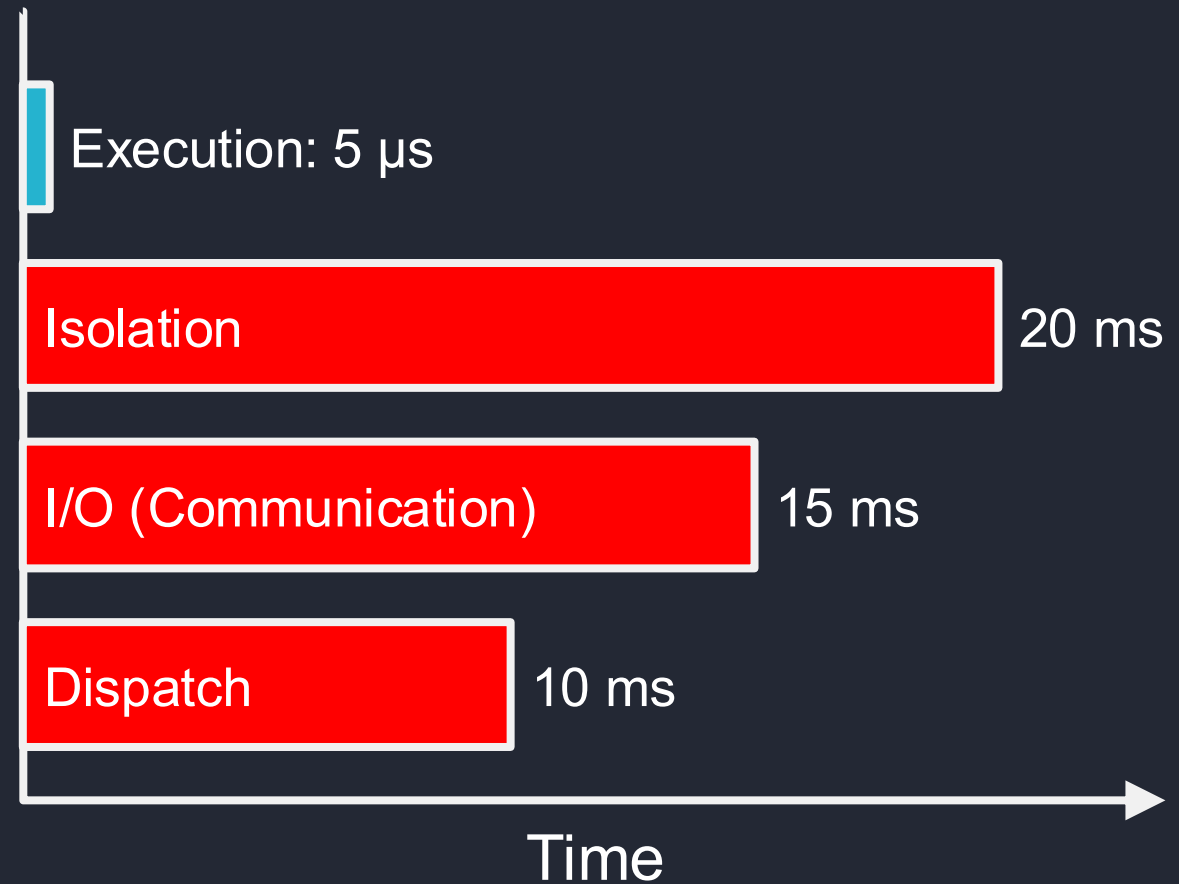
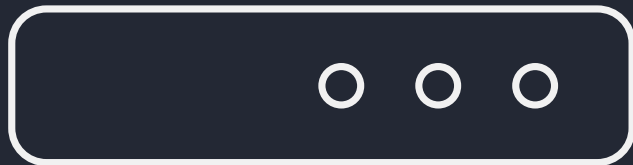
DATA  
CENTER

EFFICIENCY



- Service-centric
- Multi-tenant
- Frequently sharing
- Short-running

# OS IS A BOTTLENECK FOR APPS



# DESIGNING FOR EFFICIENCY COMPUTE

DATA  
CENTER

EFFICIENCY

Compute

Thermal

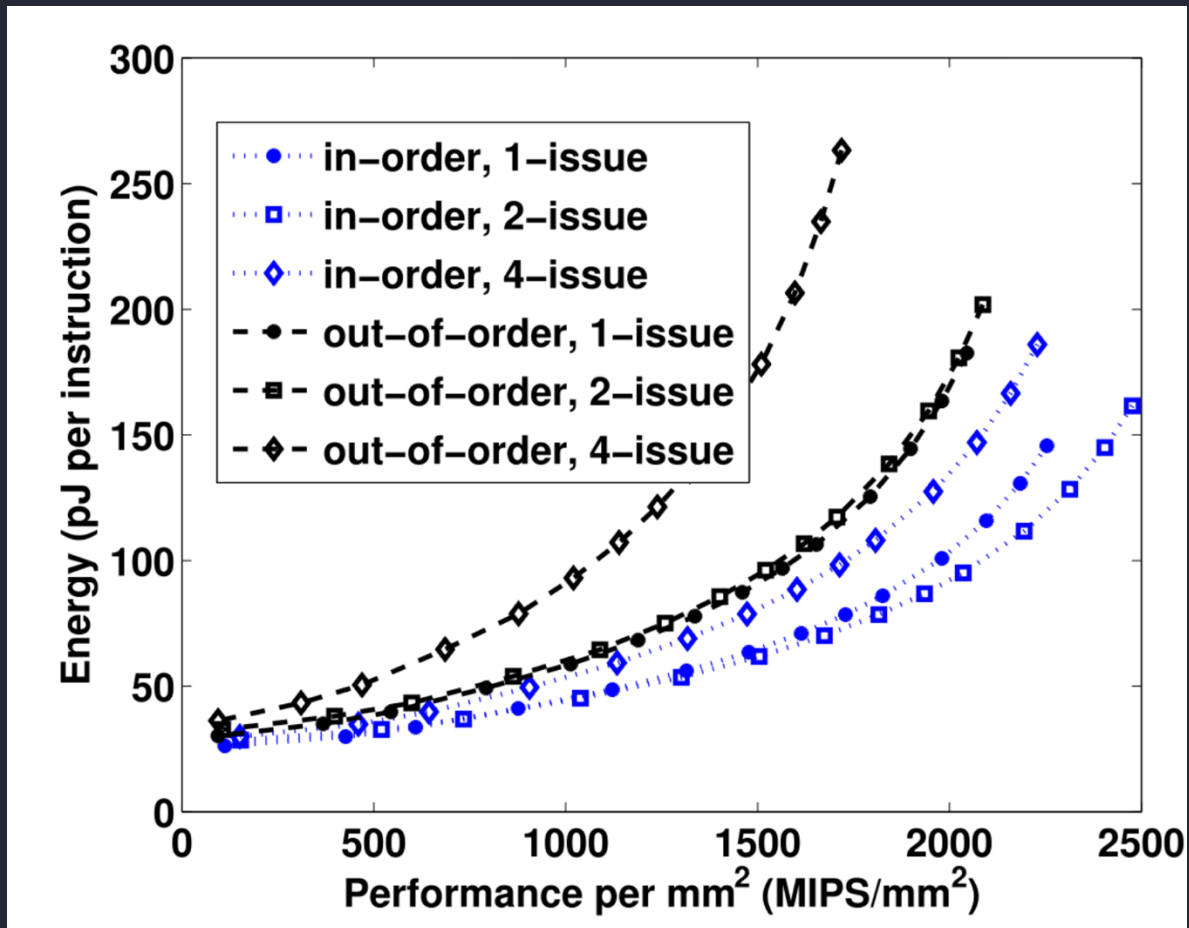
Power Conversion

Idle Power

- Design for efficiency (e.g., phones)
  - 10x better
- Tighter integration with OS/network
  - 10x better
- Optimize algo-/app-specific computation
  - 10x better

# DESIGNING FOR EFFICIENCY

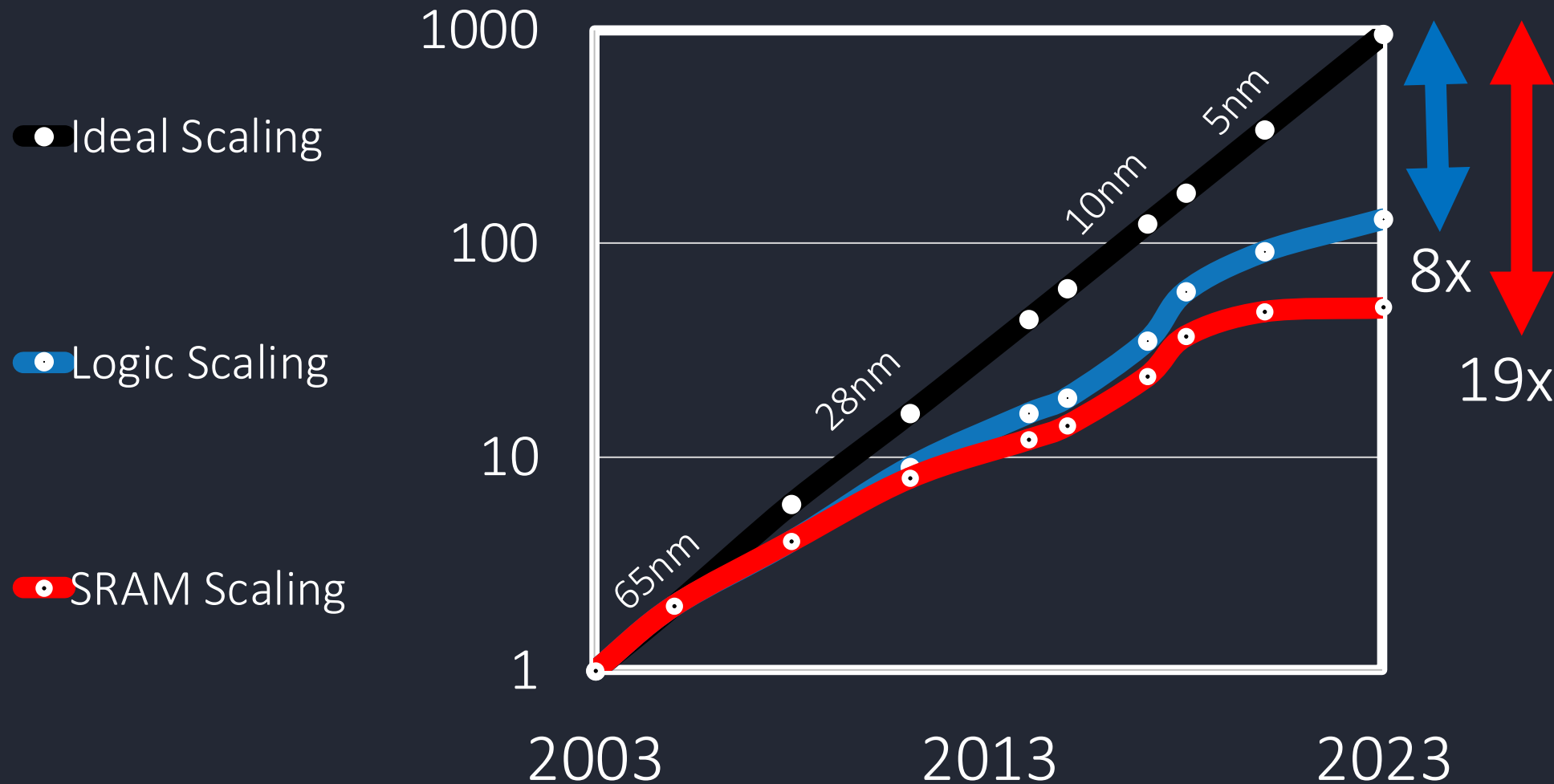
## PARETO-OPTIMAL CORE DESIGN 2010



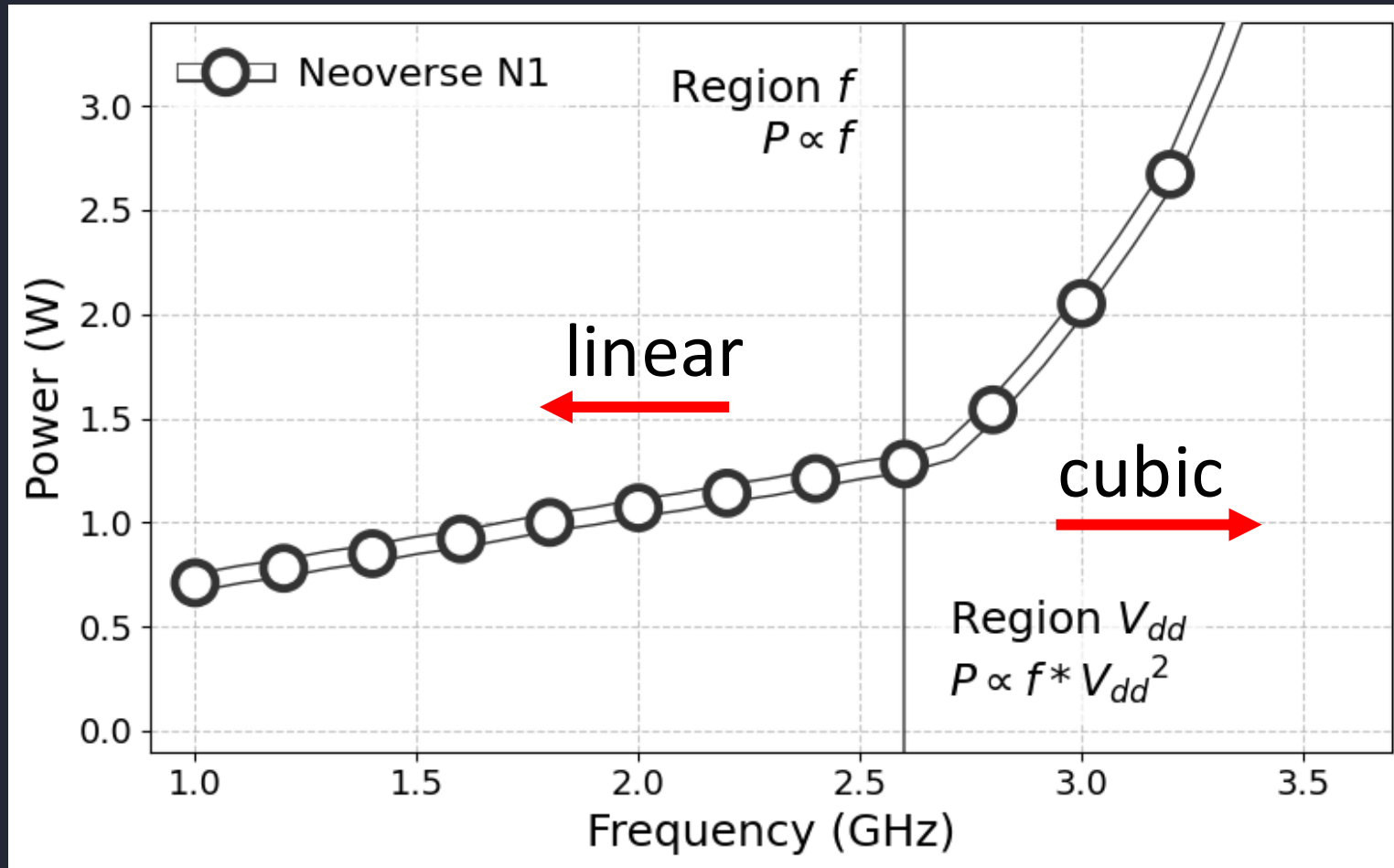
Azizi's thesis [ISCA'10]

- Per-component DSE
- Based on SPEC CPU
- No SLO
- 90nm

# SLOWDOWN IN SCALING SRAM VS. LOGIC

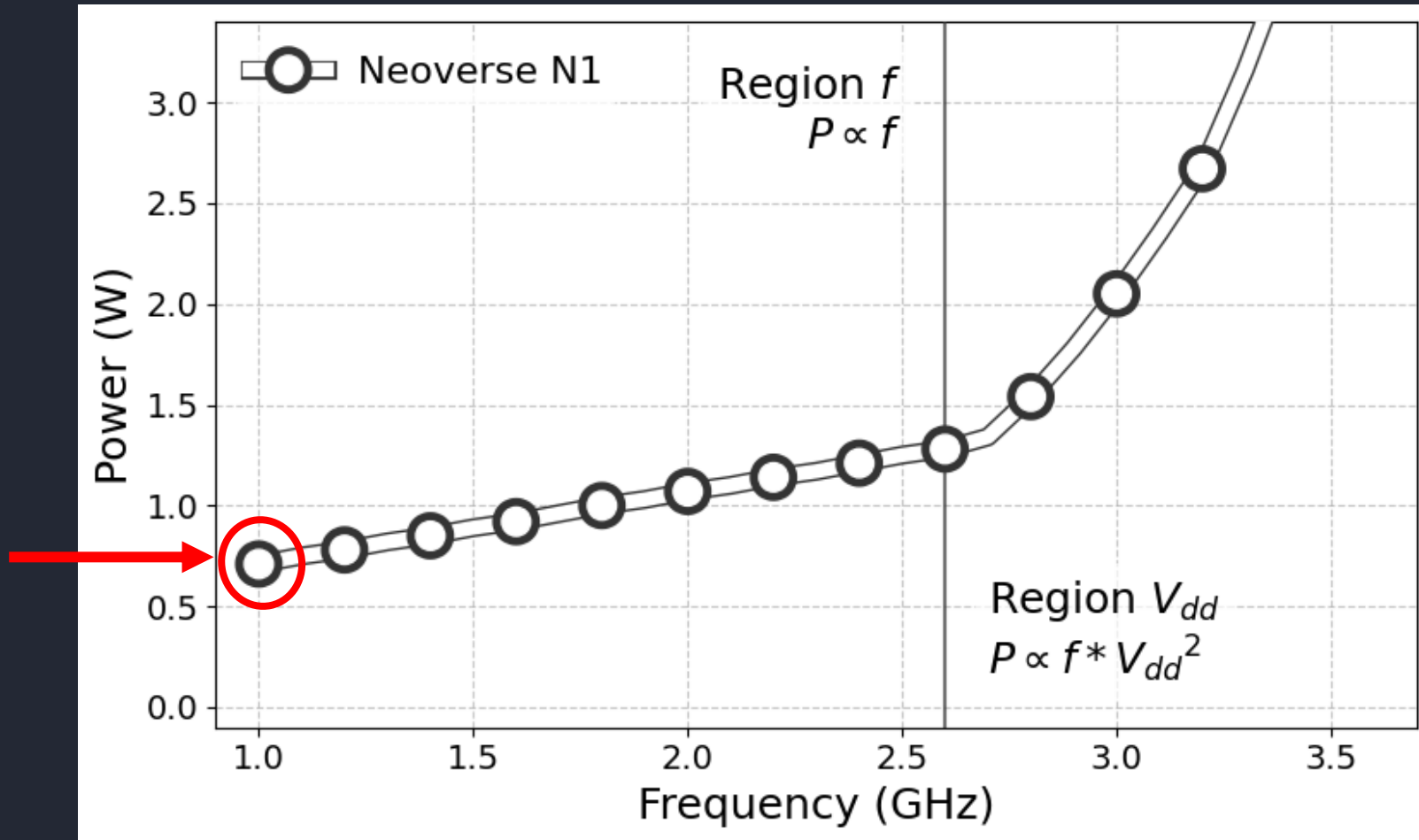


# SUPERTHRESHOLD REGION OF OPERATION CUBIC VS. LINEAR



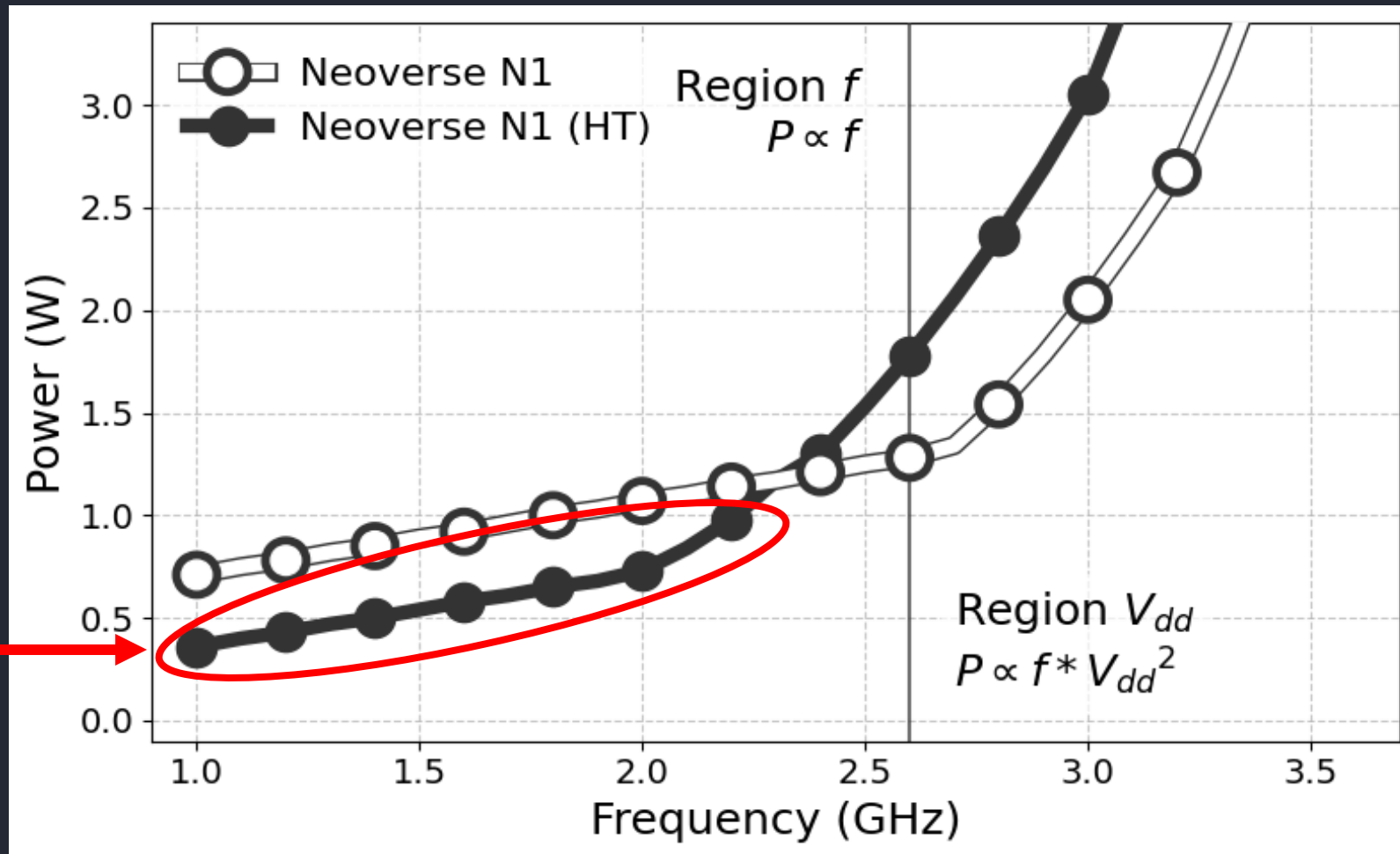
# SUPERTHRESHOLD REGION OF OPERATION CUBIC VS. LINEAR

Static  
power  
40% of  
total  
power



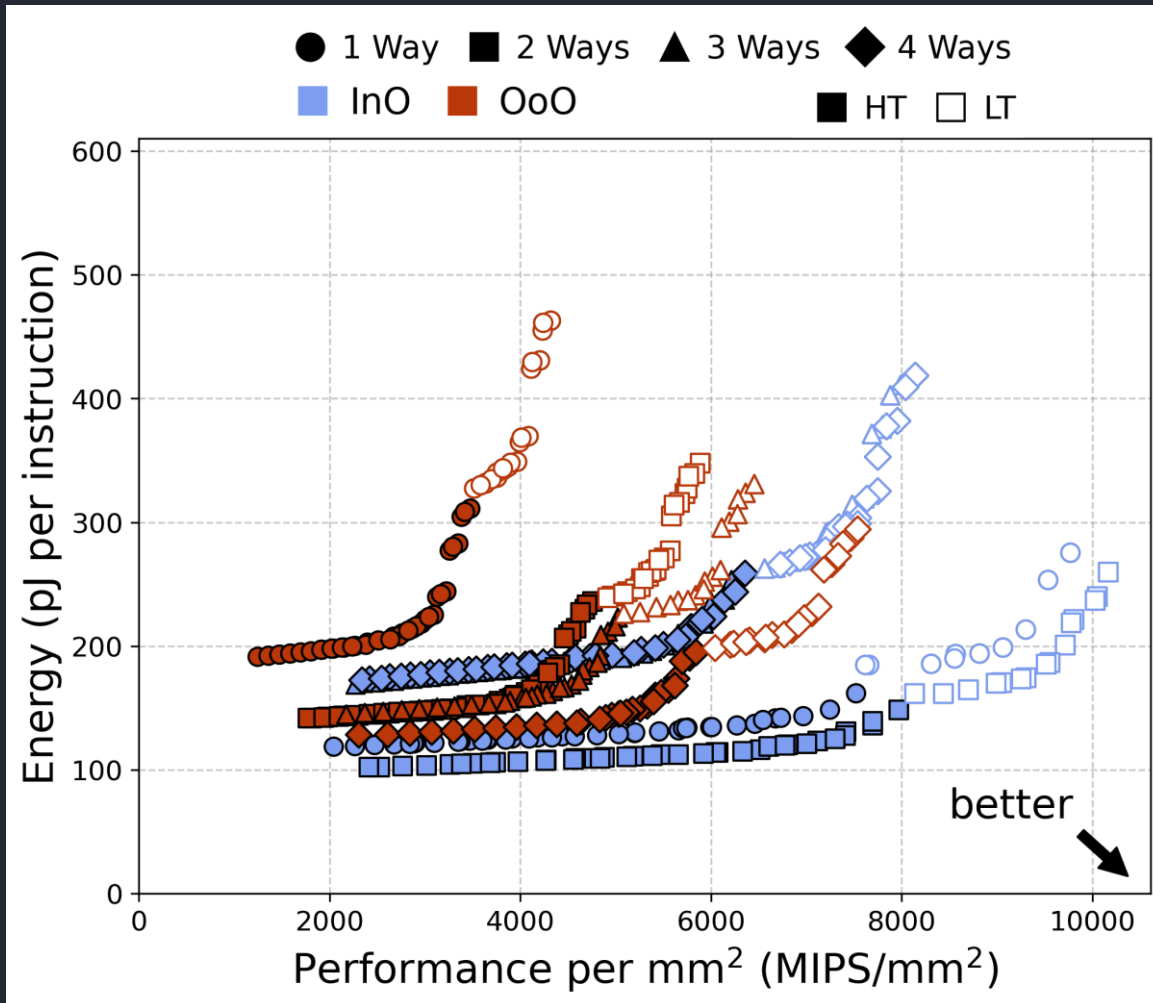
# CAN ELIMINATE LEAKAGE LINEAR REGION [VIJAYKUMAR'11]

Use high  
threshold  
voltage  
transistors  
(HT)



# DESIGNING FOR EFFICIENCY

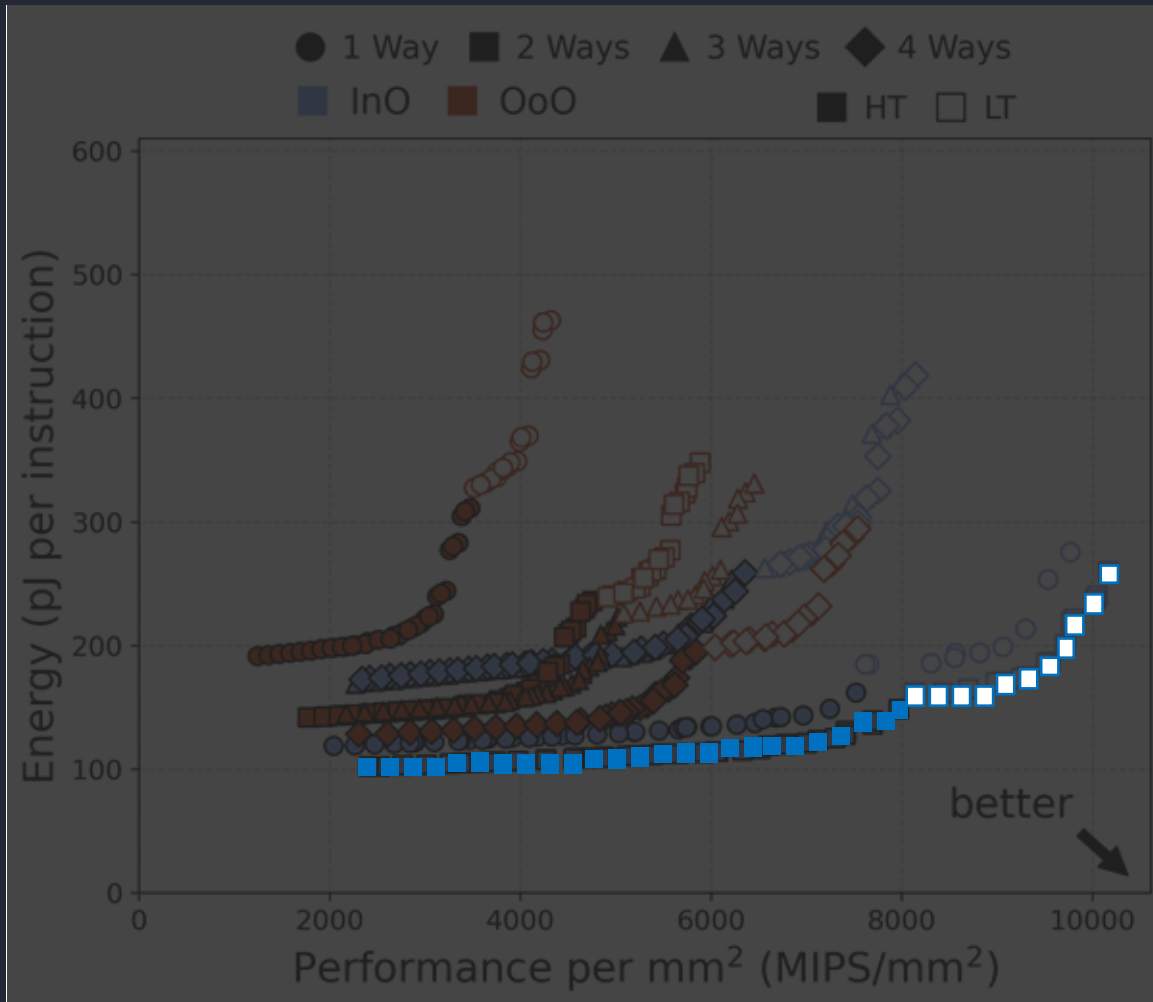
## PARETO-OPTIMAL CORE DESIGN 2025



### Redid Azizi's work

- Per-component DSE
- Server workloads
- No SLO
- 7nm
  - Supply  $V_{\min} = 0.68\text{v}$
  - Threshold LT = 0.11v
  - Threshold HT = 0.23v

# DESIGNING FOR EFFICIENCY NOT YOUR OFF-THE-SHELF CORE



2-way in-order

64KB L1s

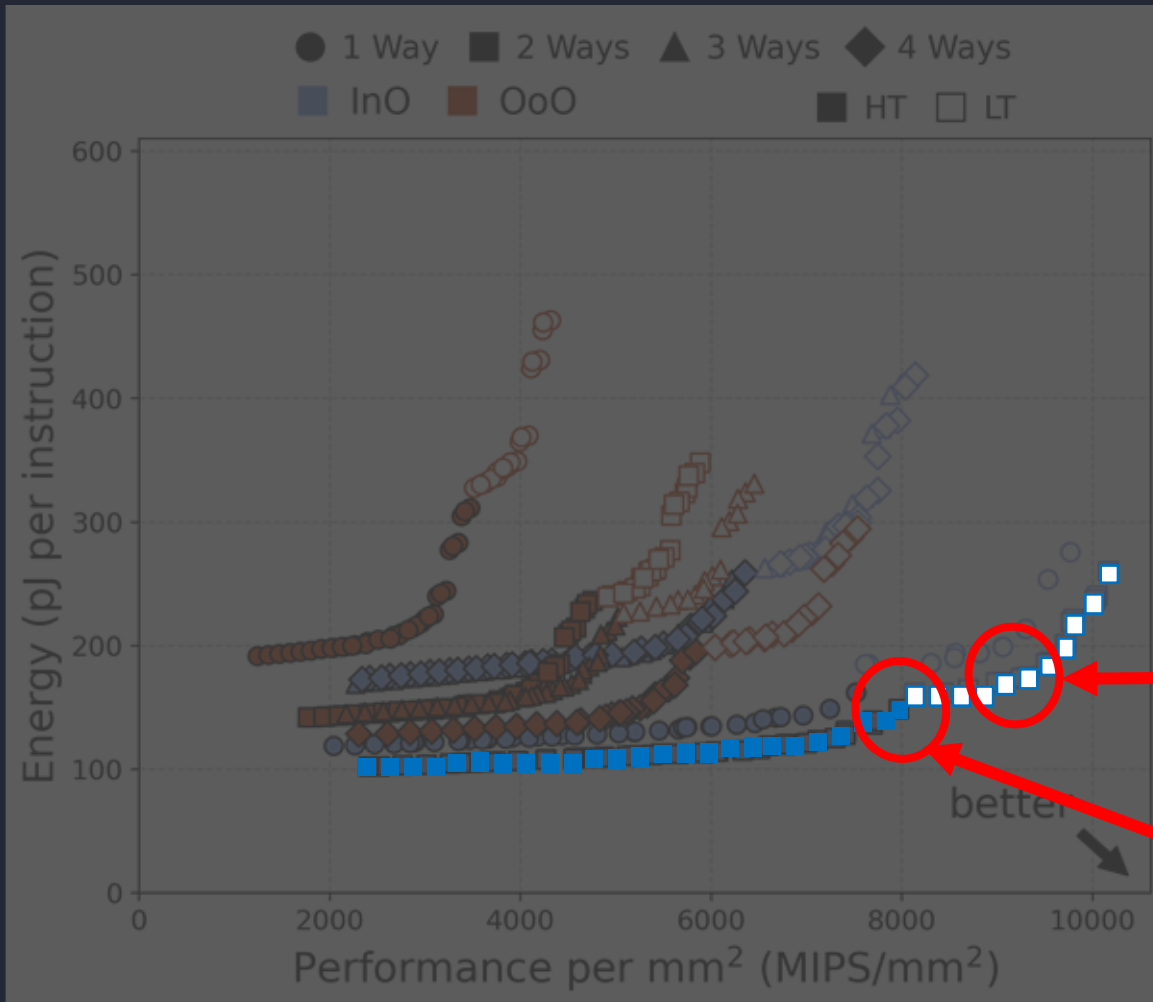
1K-4K L2 TLBs

FDIP, 8K TAGE, 2K BTB entries

SMS data prefetcher, 1K-4K

Frequency = 1.0-3.0 GHz

# DESIGNING FOR EFFICIENCY NOT YOUR OFF-THE-SHELF CORE

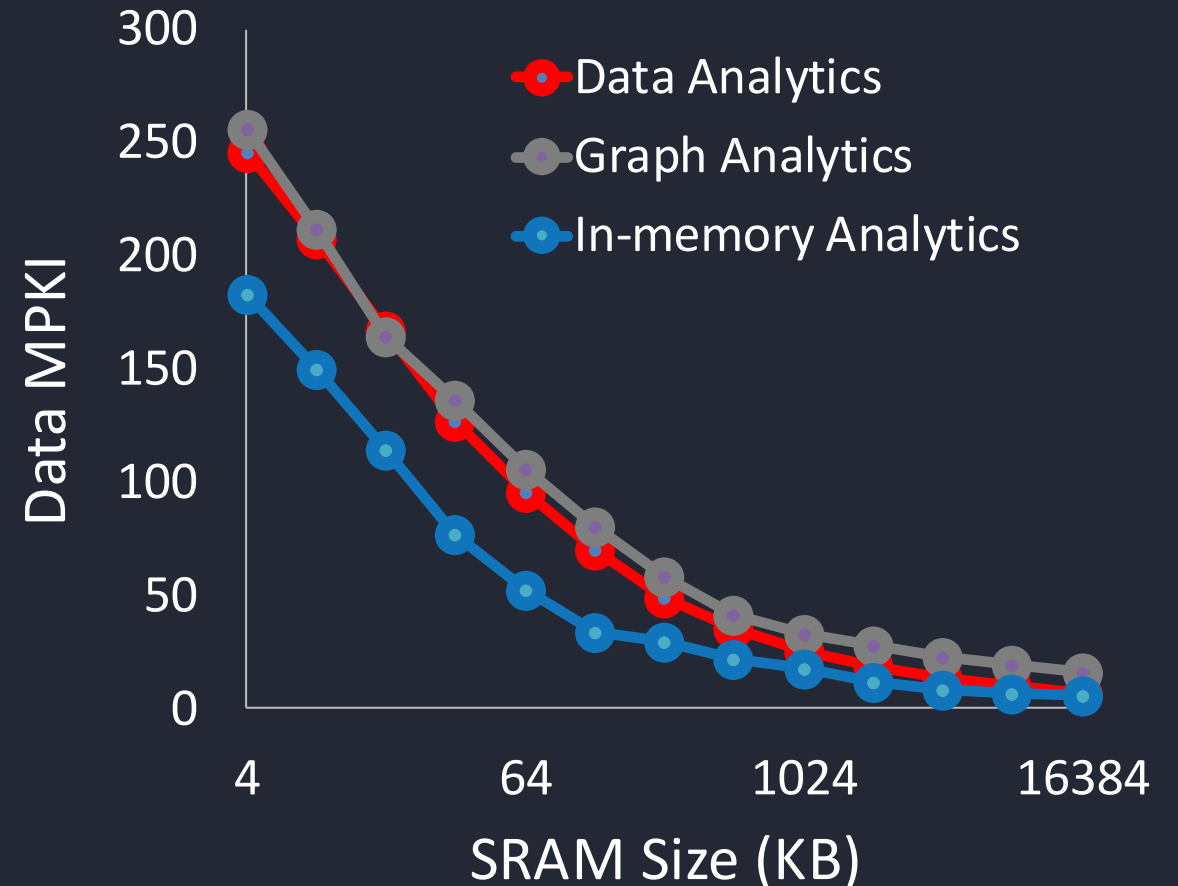
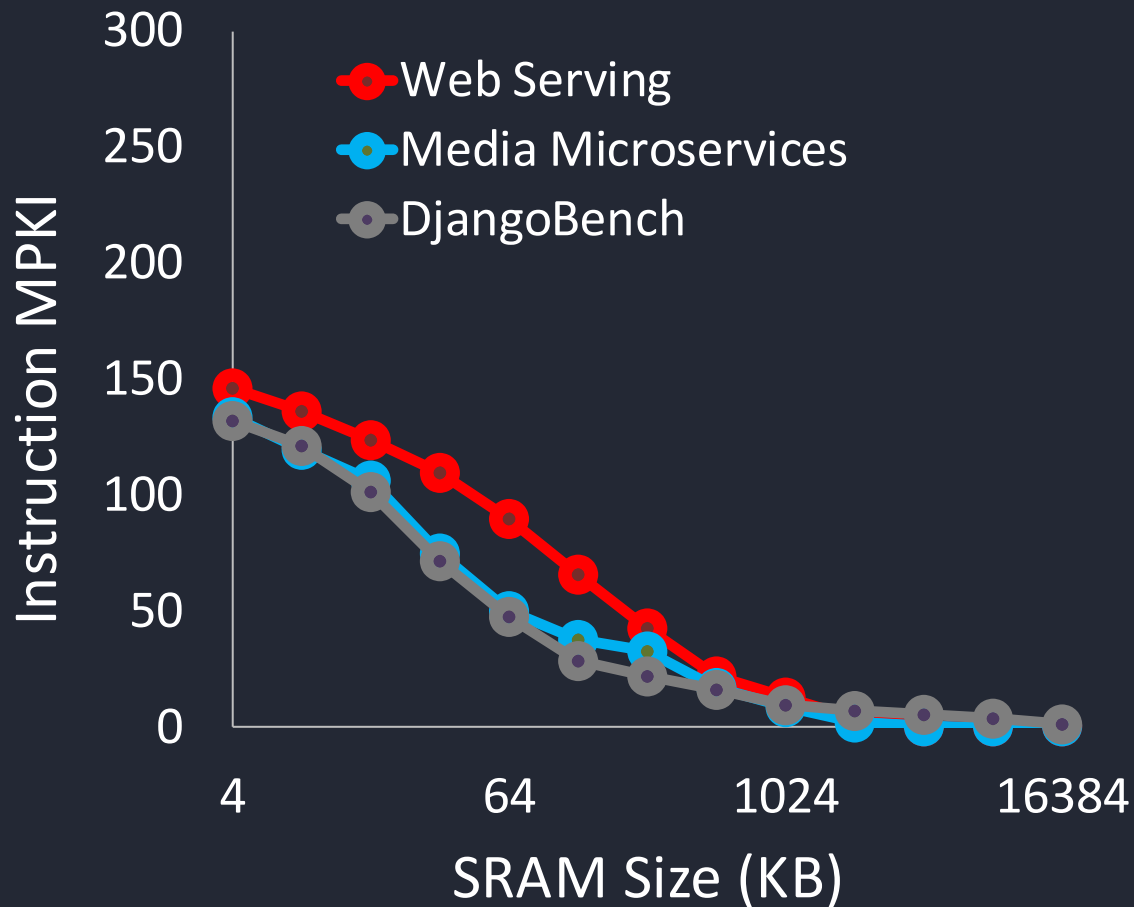


2.6 GHz

2.3 GHz

# THE SINGLE-CORE PERFORMANCE SAGA

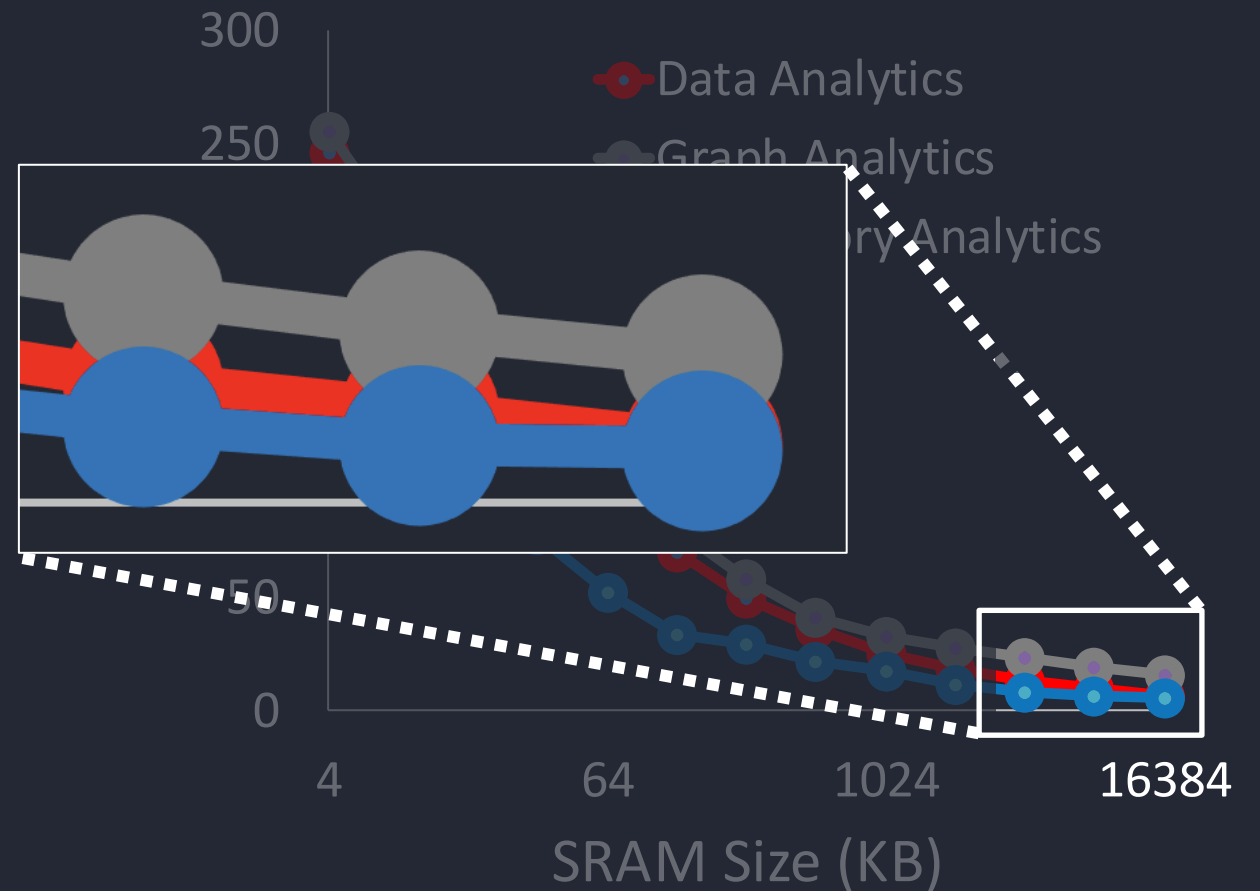
## HOW MUCH CACHE CAPACITY?



# THE SINGLE-CORE PERFORMANCE SAGA

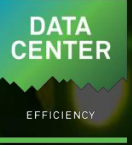
## MARGINAL GAIN WITH MORE SRAM

Diminishing returns  
from adding SRAM

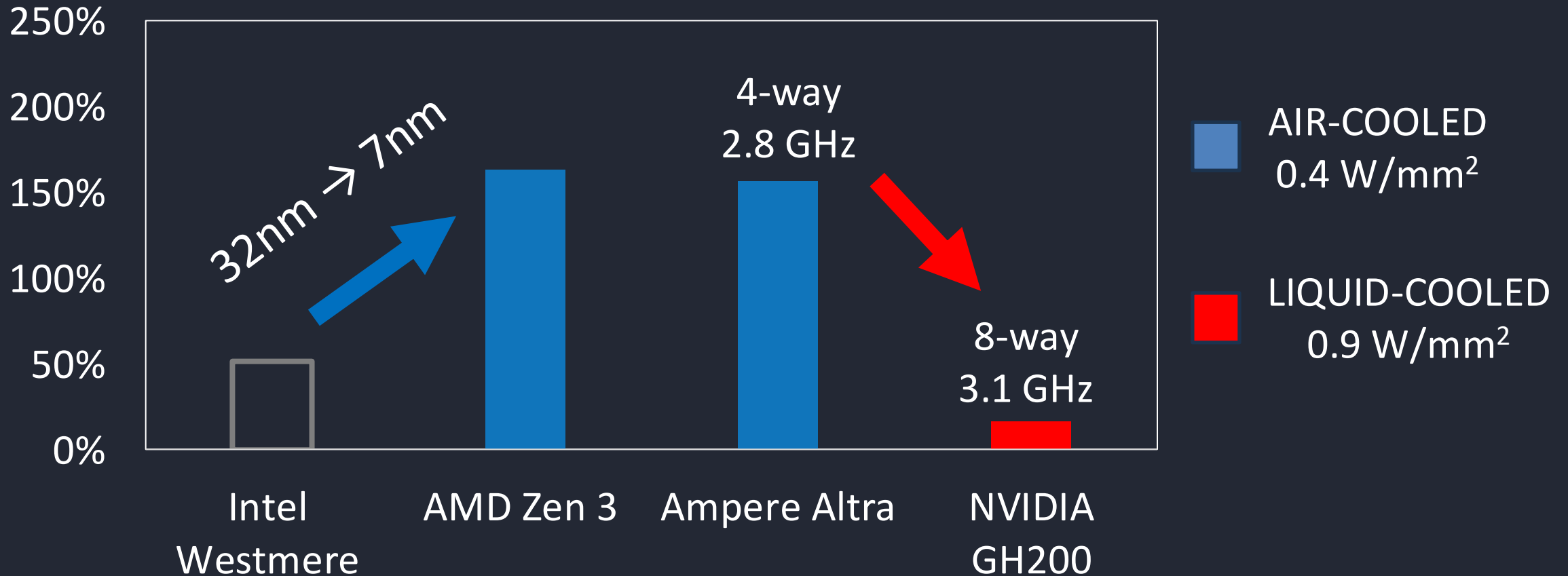


# CORES EAT POWER

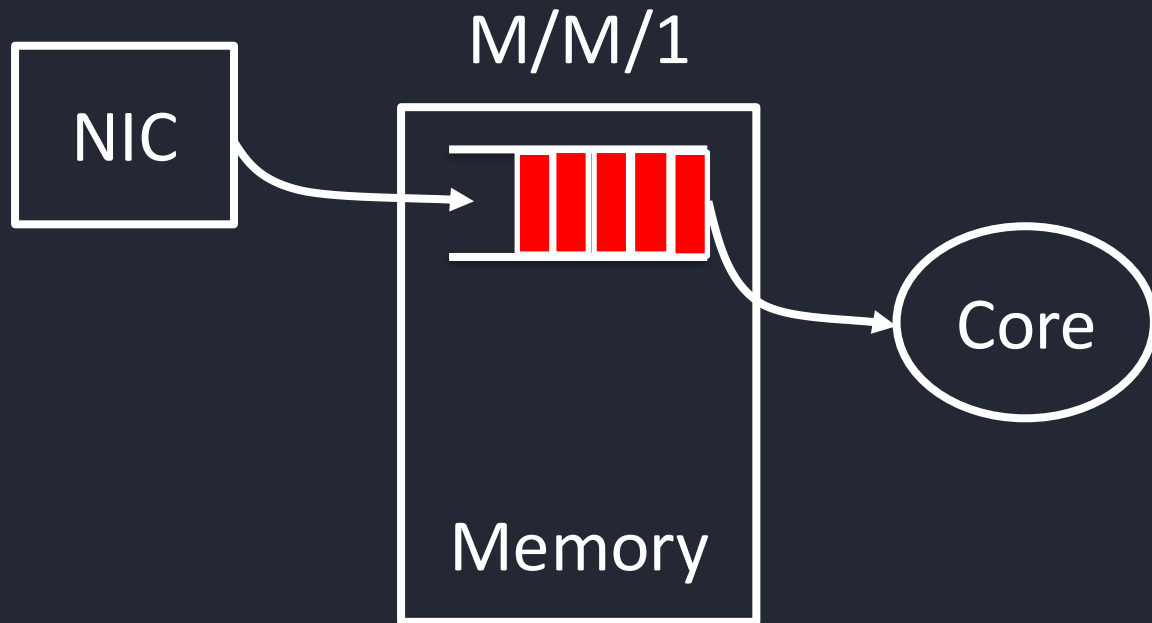
# SRAM ADDED TO MEET POWER DENSITY



Per-core additional area in SRAM (L2/LLC)



# SINGLE-CORE PERFORMANCE & SLO PERFORMANCE VS. TAIL LATENCY

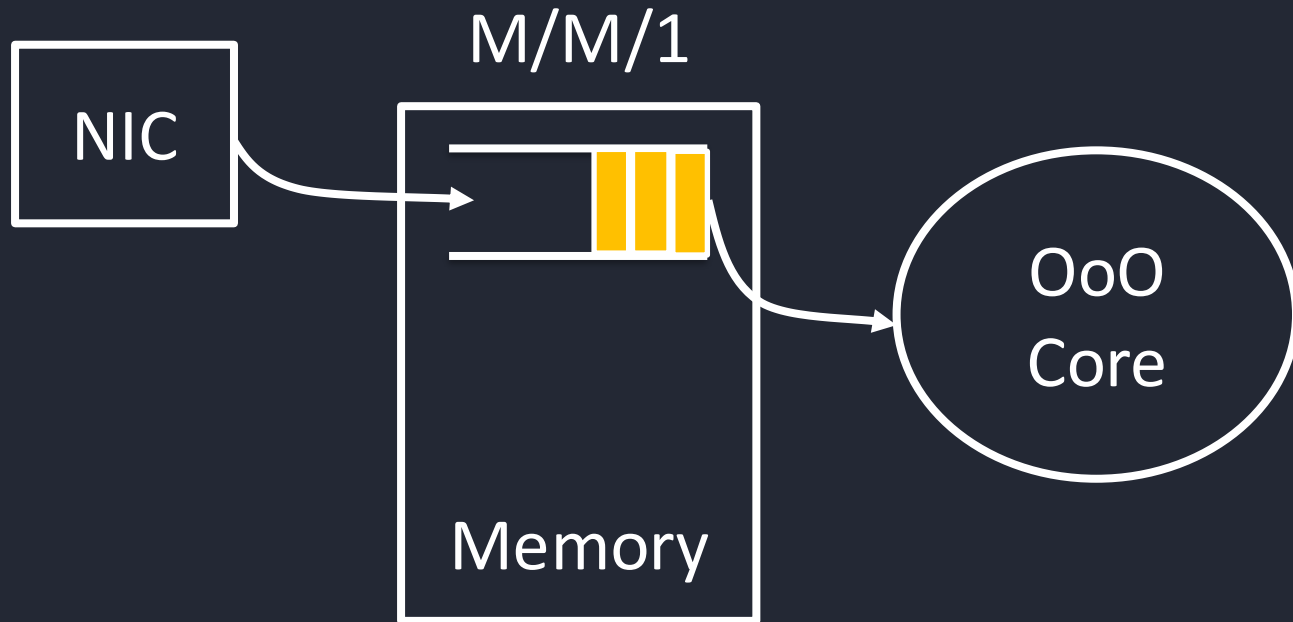


Tail latency  $\propto$  queueing  
[Delimitrou, CACM'18]

Linear  $\uparrow$  perf =  
Exponential  $\downarrow$  queueing

# PERFORMANCE VS. TAIL LATENCY CORE COMPLEXITY

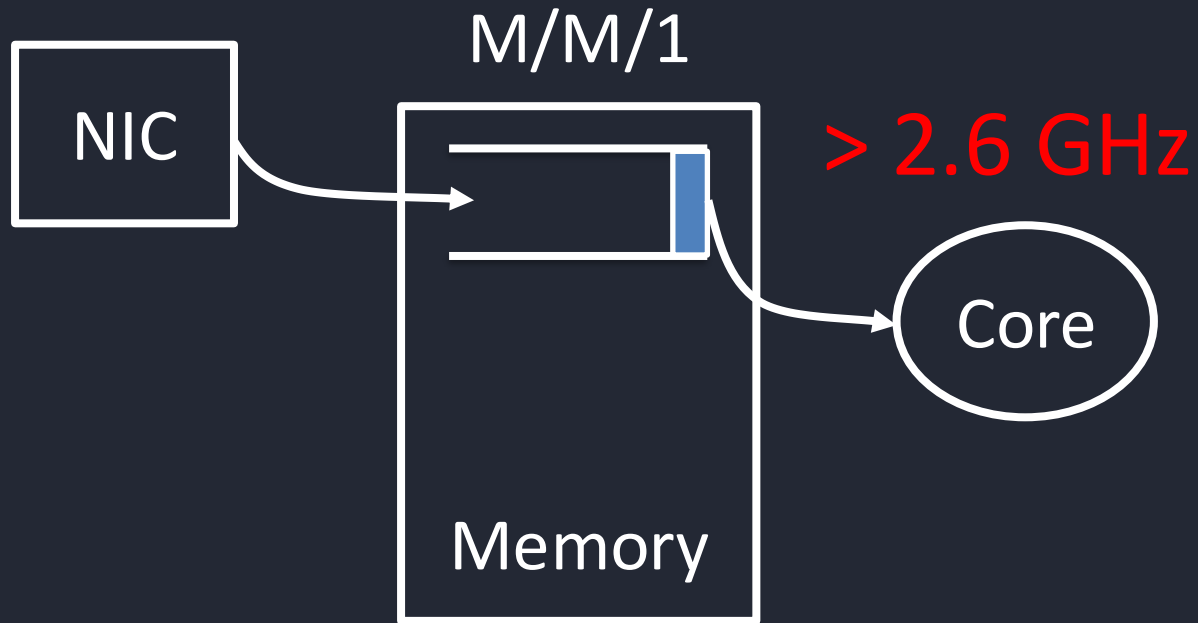
But,



with core complexity  
Pollack's Rule  
 $\text{area} \propto \text{perf}^2$

Big impact on area 

# PERFORMANCE VS. TAIL LATENCY FREQUENCY

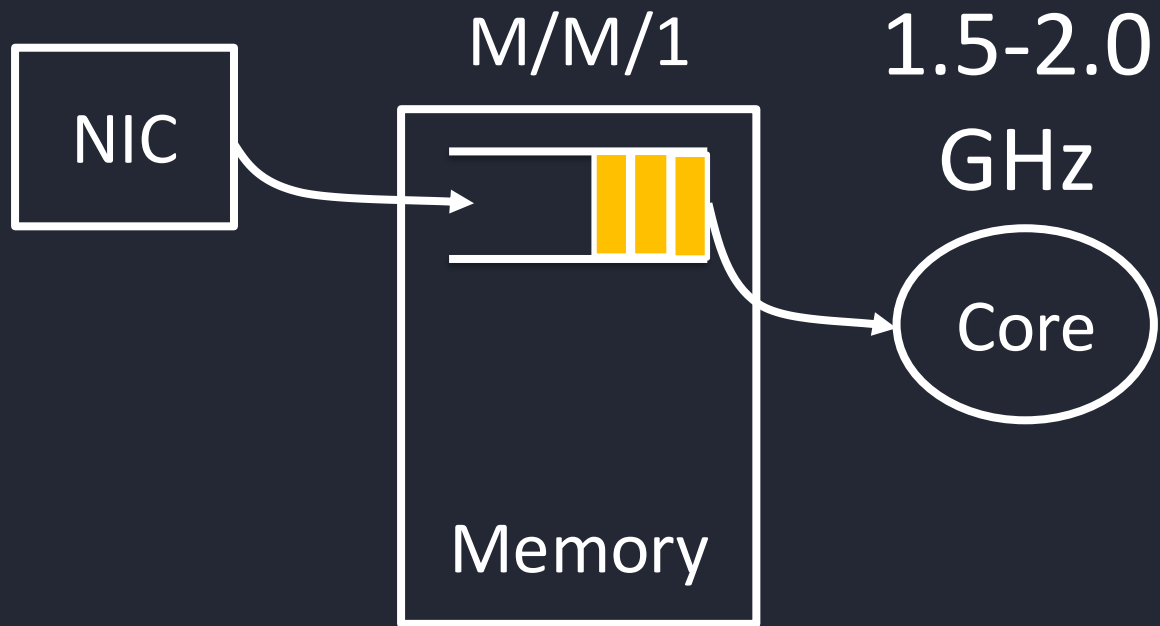


Region  $V_{dd}$   
power  $\propto$  cubic(perf)

Huge impact on power



# PERFORMANCE VS. TAIL LATENCY FREQUENCY

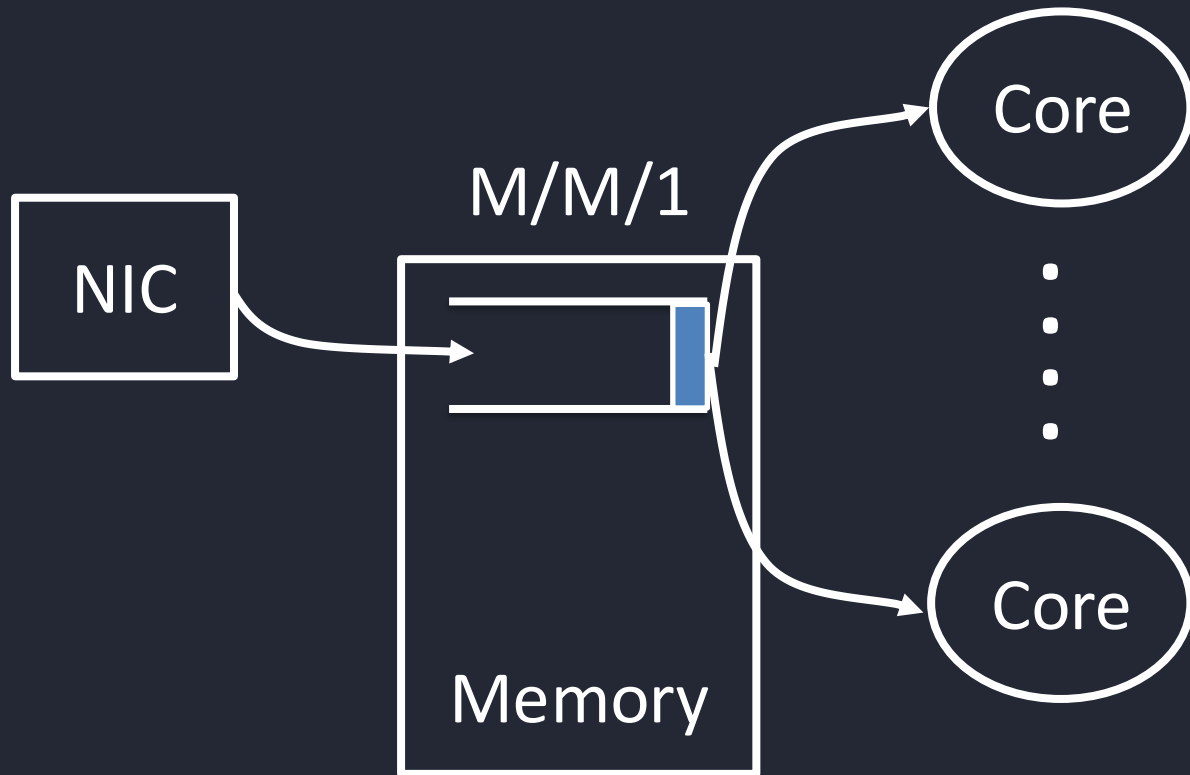


Region  $f$

Linear

$\uparrow$  perf =  $\uparrow$  power

# PERFORMANCE VS. TAIL LATENCY CORE COUNT



Single-queue servers  
[ASPLOS'19]

Linear  $\uparrow$  core count =  
Exponential  $\downarrow$  queueing

Optimal chiplets with low-  
frequency wimpy cores!

# DESIGNING FOR EFFICIENCY BACK TO THE FUTURE

DATA  
CENTER

EFFICIENCY

Compute

Thermal

Power Conversion

Idle Power



48 in-order ARM cores @ 2.5 GHz  
10x throughput/W over Intel Xeon  
16MB of L2 (no LLC)  
Blueprinted at EPFL, 2014

# EFFICIENCY GAINS WITHOUT BREAKING THE STACK

## Linux Support

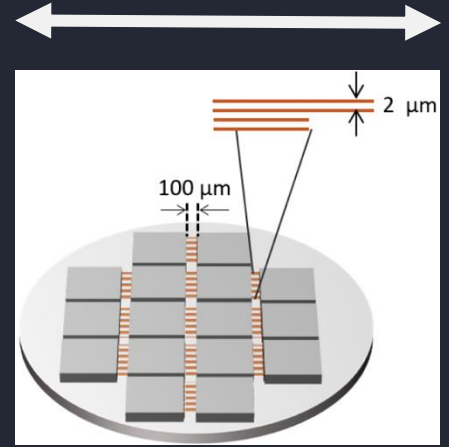
- Single-address space OS services [ISCA'25]
- User-level interrupts [Aydogmus, ASPLOS'25]
- Massive threading [Humphries, HotOS'21]
- CPU & I/O VM contracts [ISCA'21]

## Memory & NIC

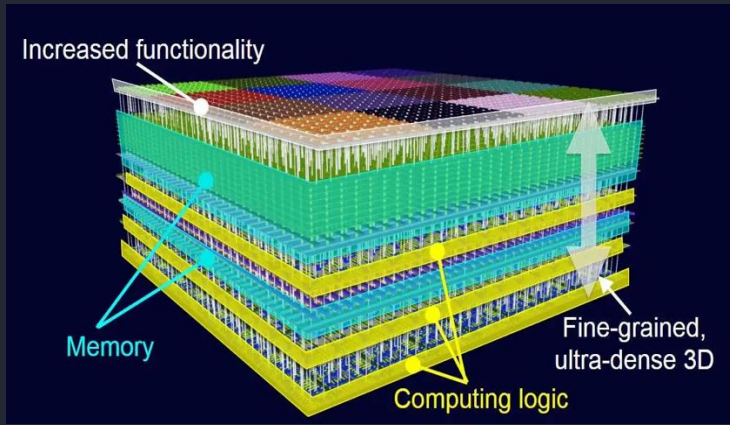
- Memory pooling [ASPLOS'15]
- Disaggregated memory [Lim, ISCA'09]
- Coherent NICs [ISCA'20]
- Single-queue dispatch [ASPLOS'19]
- RPC processors [MICRO'21]

# DESIGNING FOR EFFICIENCY FUTURE TECHNOLOGIES

Device/Chip Level



Horizontal  
lower cost  
not much denser



Vertical  
lower pJ/bit,  
higher BW/mm<sup>2</sup>

Rack Level

Disaggregated Servers



# OPTIMAL DESIGN & OPERATION

DATA  
CENTER

EFFICIENCY

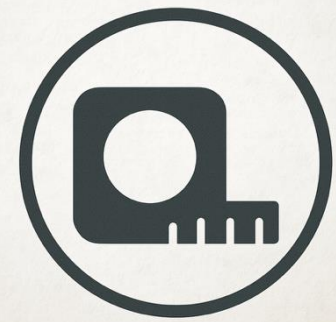
Metrics



Design



Best Practices



# THE MOST COMPREHENSIVE MEASUREMENT AVAILABLE

## DCE

Data Center Efficiency

DC INFRASTRUCTURE

IT INFRASTRUCTURE

CARBON & WATER

Power Usage Effectiveness,  
enhanced with heat recycling and  
on-premise renewables

Utilization of IT (servers, storage,  
network), efficient technology &  
operating temperature

Operational emissions and water  
consumption, including recycling  
and water stress

[navigator.sdea.ch](https://navigator.sdea.ch)

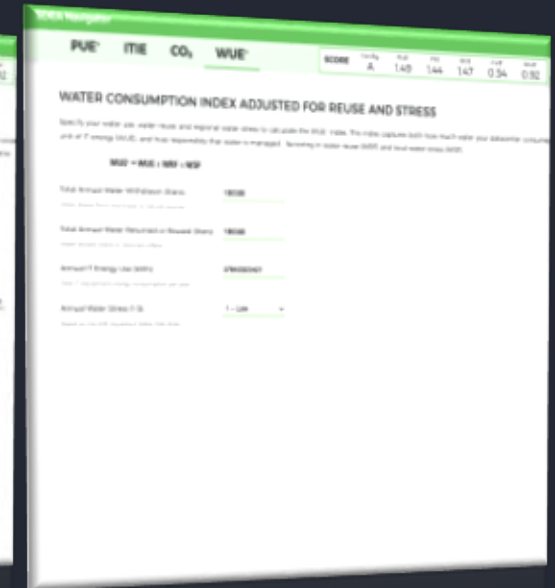
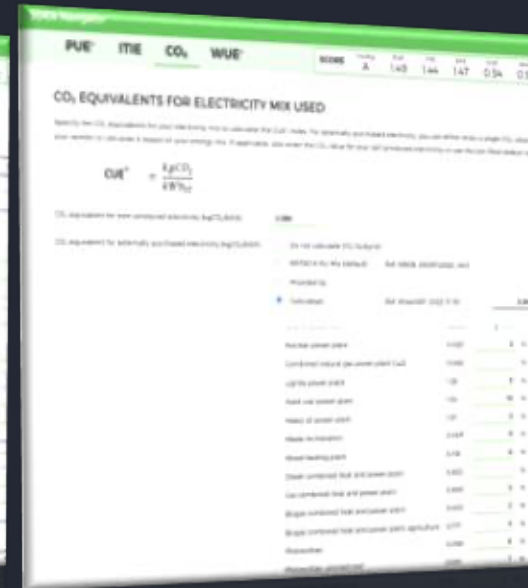
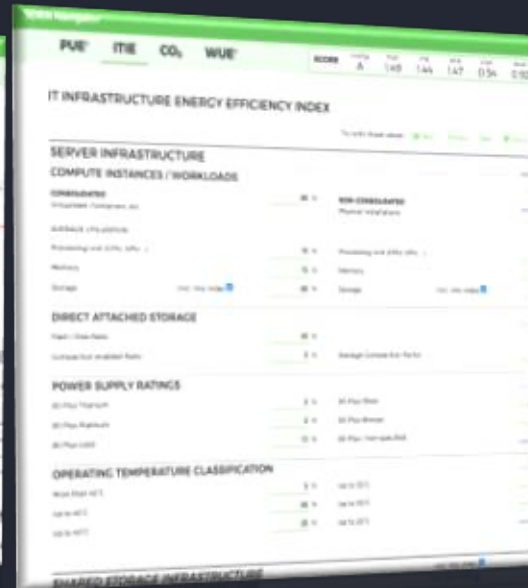
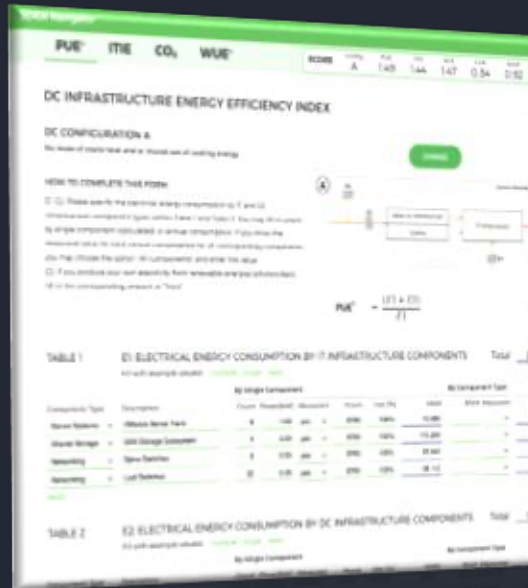
# TO MEASURE IS TO KNOW RESOURCE EFFICIENCY

FACILITY

IT

CARBON

WATER

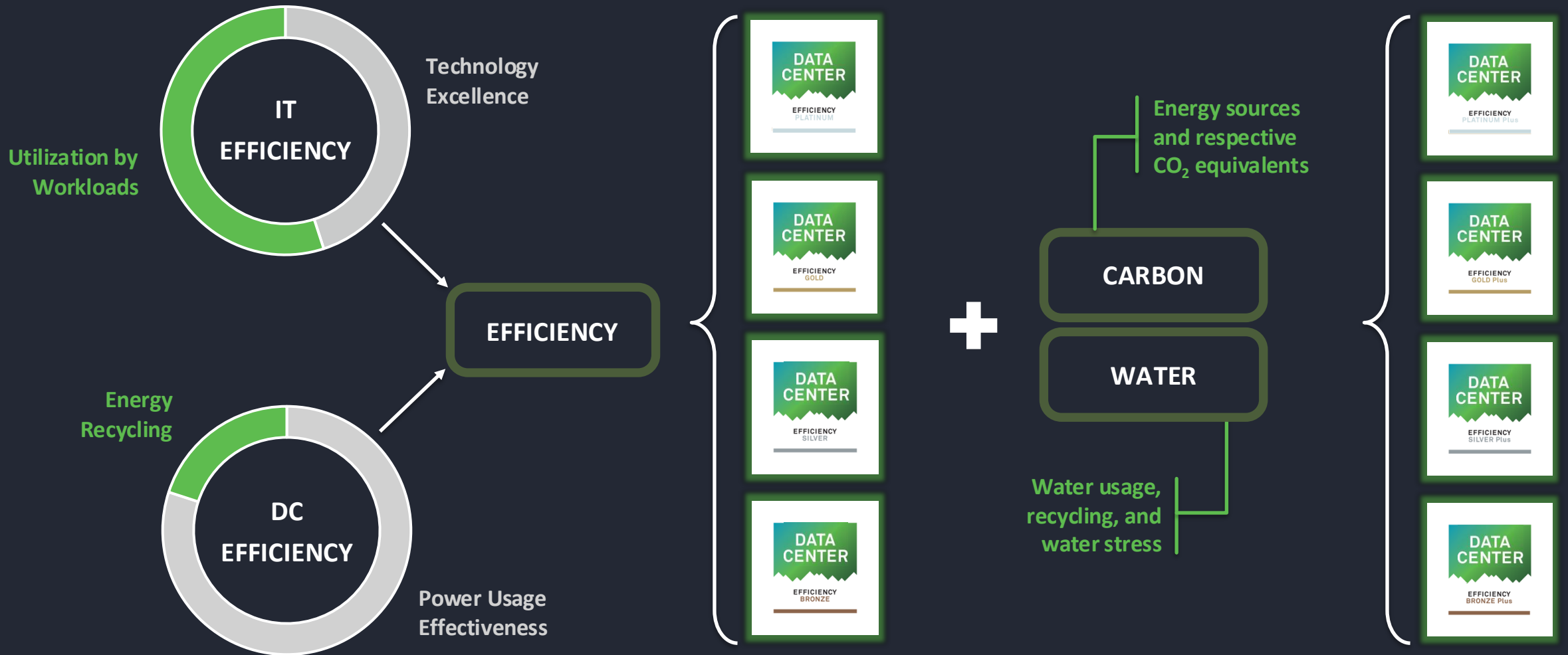


# QUALIFICATION

## COMBINED EFFICIENCY, EMISSIONS, AND WATER INDEX

DATA CENTER

EFFICIENCY



# CONCLUSIONS

## Useful compute under real constraints

- Metrics
- Design
- Best practices

## Stop chasing after single-core performance

- Opportunities w/ integration, specialization, approximation

Efficiency measured not guessed!



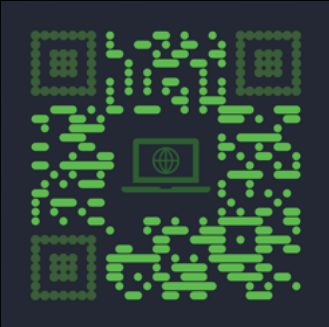
EPFL



Website



LinkedIn



Navigator Tour



Email