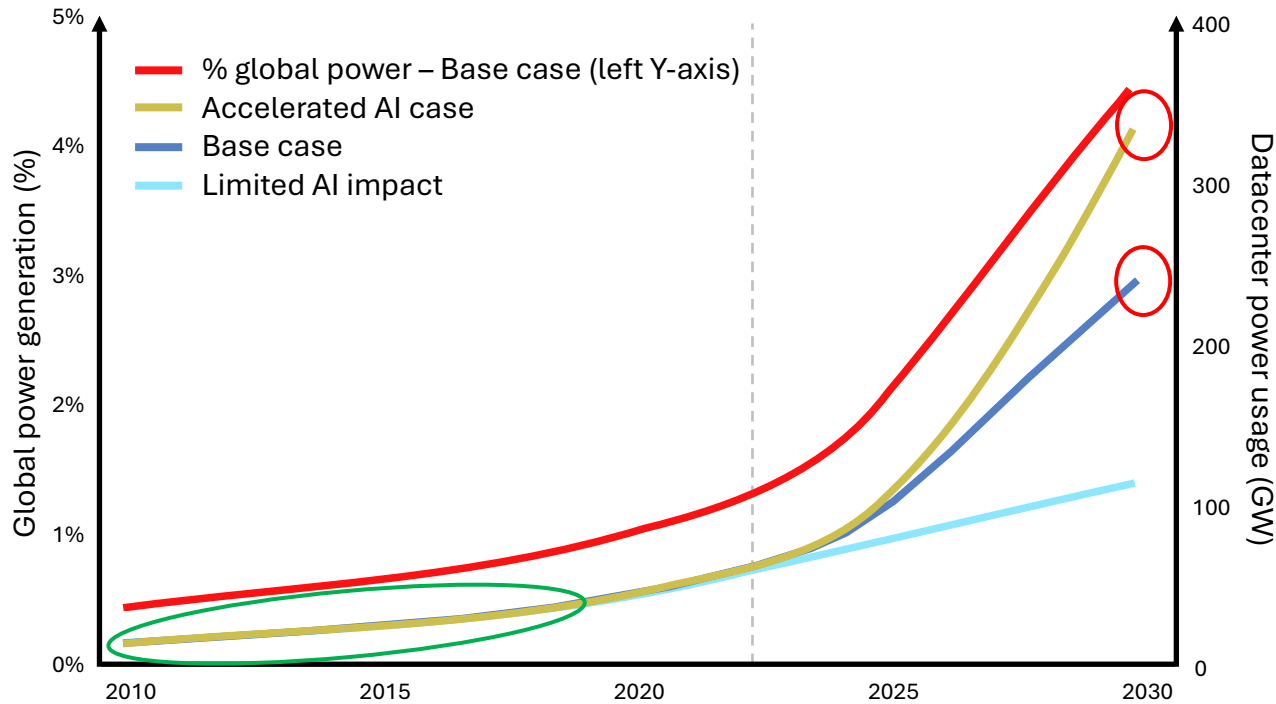# Some datacenter power and energy challenges in the AI era

## Ricardo Bianchini

Microsoft Azure

# AI datacenter electricity demand
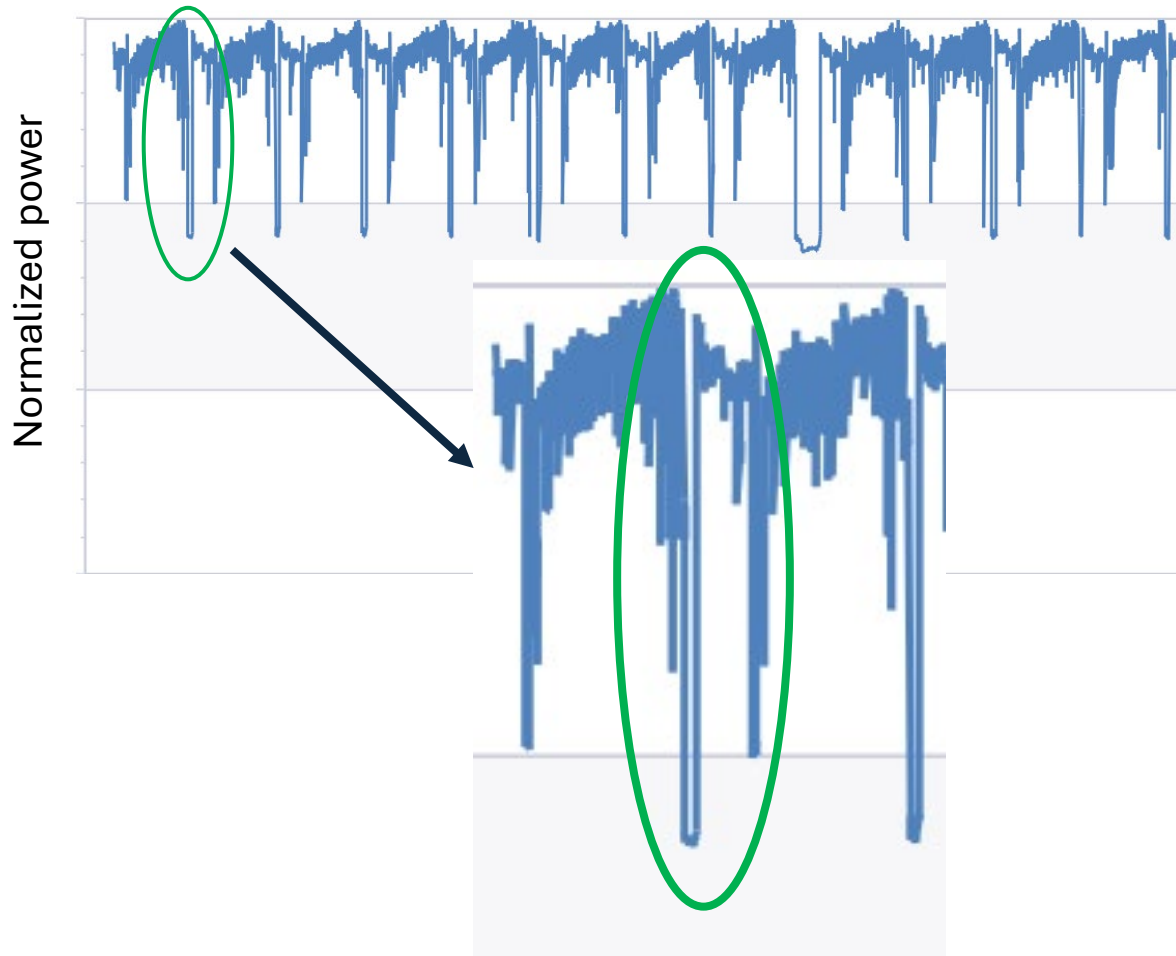
## Global Datacenter Power Usage



Legend:
- **% global power – Base case (left Y-axis)** (red)
- **Accelerated AI case** (olive/yellow)
- **Base case** (blue)
- **Limited AI impact** (light blue)

Left Y-axis: Global power generation (%) — 0% to 5%
Right Y-axis: Datacenter power usage (GW) — 0 to 400
X-axis: 2010, 2015, 2020, 2025, 2030

Base case suggests 171 GW in 2030 (3x today!)
Assuming the "accelerated AI" case, 242 GW!
Efficiency kept demand increasing slowly
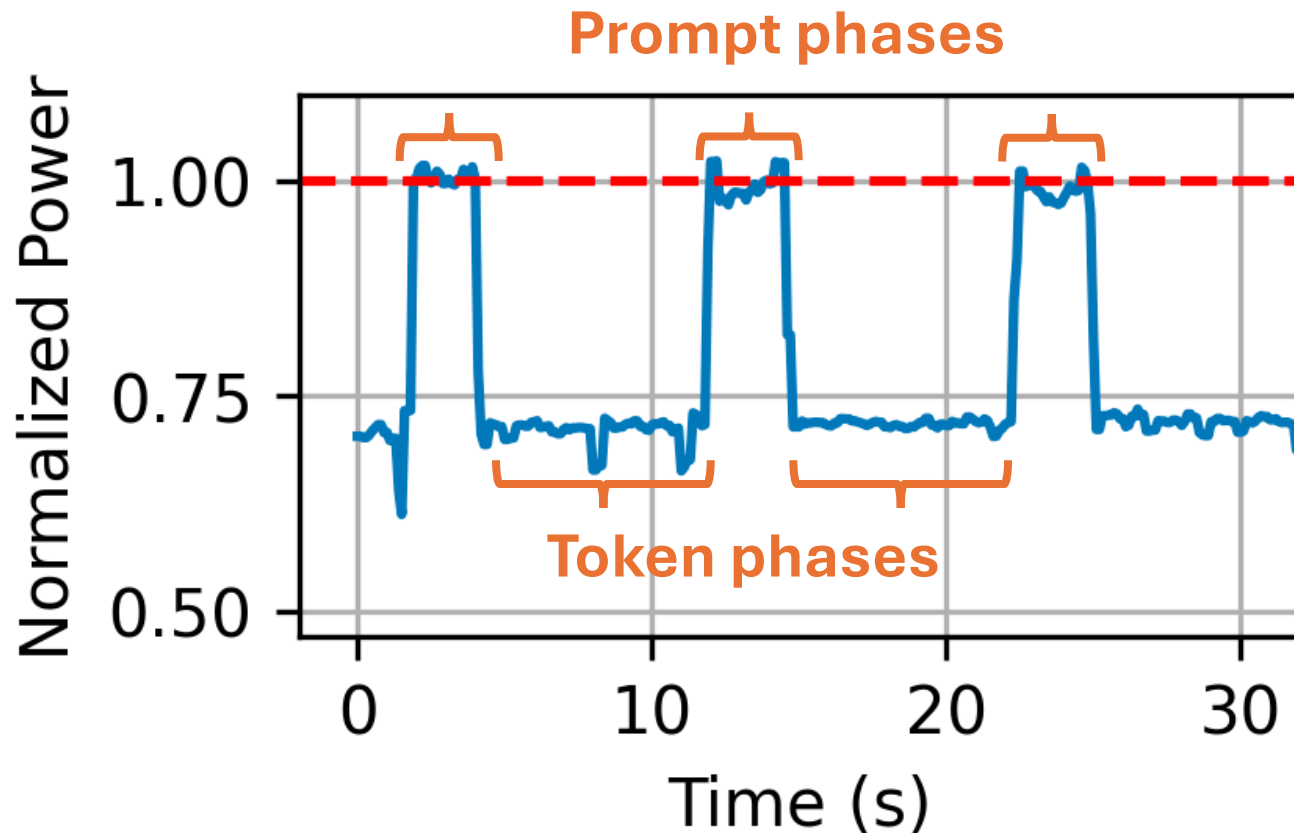
# Software: Training



Batch job with high power utilization
Massive synchronous computation

Challenges:
- Power swings can de-stabilize the grid
- Decreasing effective utilization due to failures

# Software: Inference

**Prompt phases**



**Token phases**

3x BLOOM-176B inference requests on 8 GPUs

Interactive service with lower utilization
Phases behave very differently

Challenges:
- Optimizing for such different behaviors
- Maximizing power oversubscription