# Zero-Carbon Cloud: A Volatile Resource for High-Performance Computing

Andrew A. Chien
*University of Chicago and*
*Argonne National Laboratory*
*achien@cs.uchicago.edu*

Rich Wolski
*University of California, Santa Barbara*
*rich@cs.ucsb.edu*

Fan Yang
*University of Chicago*
*fanyang@cs.uchicago.edu*

*Abstract*—The growing deployment of renewable power generation creates a growing opportunity in "stranded power", power that is generated a close to zero-cost, but not usable by the power grid. We propose to use this stranded power to create a "zero carbon" high-performance computing resource, exploiting the batch computing model to exploit the volatile power efficiently.

*Keywords*-Data Center; Stranded Power; Green computing; sustainable computing;

## I. Introduction

Data and digital assets are rapidly becoming the fundamental building blocks for society. Commerce, government, education, science, and even social interaction are digital endeavors that are consummated and optimized with information. However, the byproducts of computing are becoming endemic in society, with its carbon footprint significant in mankind's total carbon emission (computing contributes over 2% of global carbon emissions [?]). Like with all ubiquitous societal technologies, the management of the resources consumed becomes a societal issue as well. Expanding computing in conventional fashion to meet future ambitions would further increase the damaging carbon emissions of ICT. Our ambition is to create a new source of computing which has dramatically lower carbon footprint, called "Zero-Carbon Cloud".

So how can we create this transformative breakthrough? The key ideas are to 1) the harness the excess of intermittent power created by the shift to renewables, and 2) exploit light-weight data center infrastructure to radically reduce other elements of total-cost of ownership (TCO). Our concept, "Zero-Carbon Cloud" combines these to create limitless, low-cost computing with different volatility characteristics than the ones exhibited by computing systems today. We propose ZCCloud that exploits volatile renewable power generation the non-interactive (batch) model of most high-performance computing to create a low-cost powerful HPC computing resource. Note that ZCCloud is a *pure* renewable-based computing services, a radical contrast to greening efforts [?] that purchase a balancing average of renewable power.

It is also worth noting that renewable's such as solar and wind variable output makes them challenging for integration into a reliable power grid. Such grids have been designed and engineered for controllable generation at fixed locations – and highly optimized for connectivity and cost on that basis. The shifting quantity (and consequently location) of power generation by renewables increases transmission requirements due to congestion and greater distance from generation to consumption. This poses serious technical and economic challenges for the power grid [?]. The widespread adoption of ambitious Renewable Portfolio Standards (RPS) that set goals for rapid growth in the fraction of renewable power that utilities must employ, the variability challenge is tremendous and growing. The dynamic range of such resources exceeds 50% of peak load today, and may increase to 100% within 15 years [?]. To the grid, Zero-Carbon Cloud is an example of a dispatchable load, that both creates on demand a high-value service, but and a new volatile form of computing.

New intellectual concepts for volatile cloud systems and applications, including systems, service-level agreements, prediction, and scheduling and marketplace are needed to realize this vision.

## II. Renewable Opportunity

Growing concerns about the impact on climate and environment of carbon emissions resulting from burning of fossil fuels [?], [?], [?], [?], have led to large-scale deployment of renewable sources of electricity generation. By far the fastest growing types, and those projected to address a significant fraction (¿10%) of demand are wind and solar [?], [?]. The variability and non-dispatchable nature of these renewable sources, combined with low incremental generation cost creates significant challenges for power-grid design and management [?], [?], [?], [?]. At present, when generation exceeds demand and the excess power exceeds the grid storage's limited abilities, it is simply discarded at the source – it is "stranded power." Power grids call this loss of excess power "curtailment" or "down dispatching". It is this opportunity that we propose to exploit with ZCCloud.

Numerous power grids (Independent System Operators – ISO's) around the world have stranded power, Figure 1 reports data from the Midcontinent Independent System Operator (MISO), showing total generation, total wind power, and total power "curtailed/down-dispatched" for a recent two and a half year period. Despite improved grid connectivity
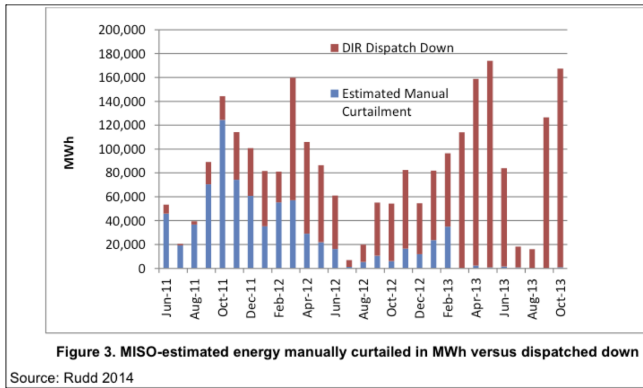
Figure 3. MISO-estimated energy manually curtailed in MWh versus dispatched down

Source: Rudd 2014

Figure 1. Stranded Power in the MISO Power Grid, June 2011-October 2013 [?]

and management has, and MISO's economic dispatching market still suffers from a few percent waste, an extraordinary amount of power – 1.6 TWh in 2014. Comparable levels prevail for ERCOT (wind) and CAISO (solar and wind), and numerous regions in Europe (Denmark, Germany, Ireland, Italy) [?]. In all of these power grids, the fraction of renewables is expected to increase by 100% or more in the next decade, creating even greater challenges to the maintenance of power-grid balance, and more stranded power [?], [?].

In the MISO region, wind power has significant penetration today with smaller states such as Iowa and Minnesota ¿20%, and larger states such as Illinois and Michigan at 5%, but all of thes states have adopted Renewable Power Standards (RPS) goals to double this percentage by 2025. [?].
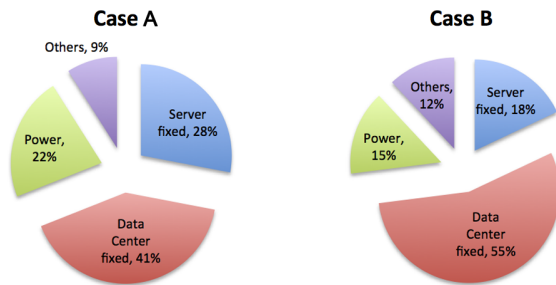


Figure 2. Total Cost of Ownership, based on [?], for low-cost server and partially filled data center scenarios. ZCCloud can substantially reduce power and data center costs, accounting for 63-70% of TCO.

## III. ZCCloud Approach

The basic approach of Zero-Carbon Cloud (ZCCloud) is to exploit recent technological advances in Information Technology (e.g. cloud computing, data center automation, system-focused analytics, etc.) to leverage "stranded" power in renewable energy settings. The result is a cloud computing
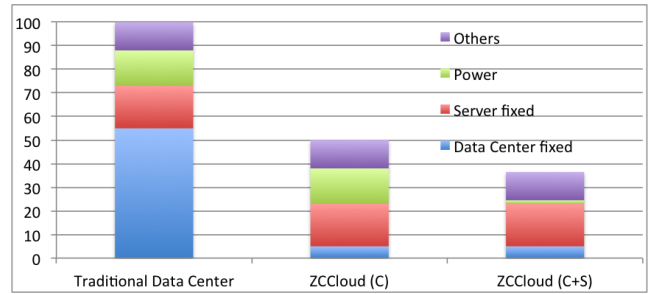


Figure 3. Zero Carbon Cloud Reduced Costs: Containers (C) and Stranded Power (S)

capability lower fixed cost (lower physical plant), lower variable costs (lower power costs), and ultimately lower overall TCO for delivered computing.

Published data suggests data centers with more than 50% physical plant and nearly 20% server costs with the sum accounting for about 75% of the TCO [?], [?]. The primary element of the remaining 25% is electricity as illustrated in Figure 2.

ZCCloud can reduce these costs in two ways.

1) Using containerized server facilties, sited at renewable generation sites, ZCCloud eliminates the need for purpose-built buildings to house the infrastructure and power distribution.

2) Exploiting stranded power, the cost of power can be reduced well below even the wholesale prices paid by large data centers , perhaps ten-fold.

We outline a case study of traditional data center in Figure 3, along with a projection for achievable ZCCloud cost. Together these improvements suggest a system that could provide computing 2-fold cheaper by exploiting a decentralized architecture, and as much as 3-fold cheaper if stranded power is exploited.

ZCCloud provides an energy sink for curtailed renewable power that is capable of producing useful computational work. That is, the curtailment is simply converted to computation and storage capabilities (albeit with different volatility and availability cycles than traditional systems) rather than being discarded as extra, unconsumable power. Once the excess power becomes computation, however, modern data center efficiency technologies can be used to maximize the utilization of this power. For example, heat reclamation techniques commonly used today [?] in many top-end data centers to maximize energy efficiency become relevant. That is, the curtailed power that is converted to computing and storage can be made to do so with increasing efficiencies using the technological advances that are improving datacenter efficiencies today and in the future.

## IV. Design and Scaling ZCCloud

Superficially, wind or solar power may seem to be unusable for computing due to their intermittent availability. In

Figure 4. Scaling from **Small**: 2-container, 2.2 Petaflops, 0.59MW, **Medium:** 4-container, 4.4 PF both can be power by a single turbine. **Extreme**: 42-container, 45.5 PF, 12.4MW, 5% of turbines in a wind farm.
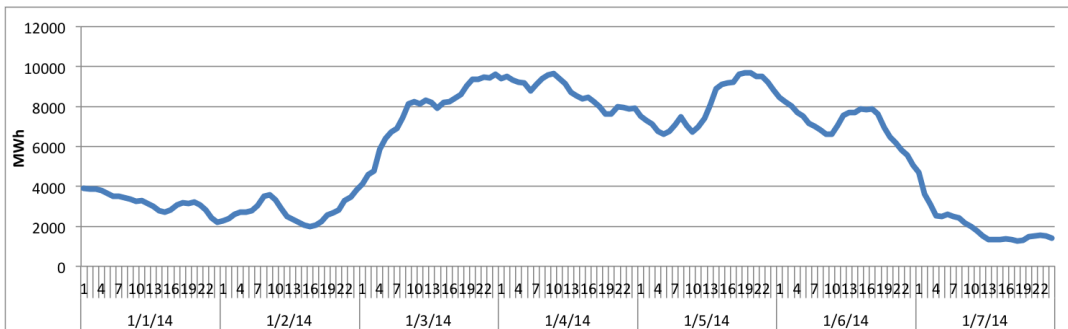


Figure 5. MISO Wind Generation for One Week.

fact the rates of available wind power change on long time scales (several days or weeks) or as short as a few hours. Solar power is more predictable, varying in a similar daily cycle approximately matched to increases in traditional demand (daytime higher, nighttime lower). We detail a sketch ZCcloud system design, and then discuss intermittence.

**Design of ZCcloud Building Blocks:** ZCCloud uses convention building blocks that achieve high densities of computing per rack and per container. The computing nodes are connected with low-latency, high-speed 10-gigabit Ethernet switches. We assume the containers have a Power Usage Effectiveness (PUE) of 2.0, which means the non-compute facilities, e.g. cooling system etc. has the same power consumption as compute nodes. This conservative and published numbers for commercial container-based products and hyperscalers such as Google are now below 1.25 and as low as 1.19. To enable real-time response even when stranded power is unavailable, ZCCloud also deploy an always-on frontend server for each container, the power consumption of which is nontrivial comparing to the total container power. The resulting power and computing density is summarized in Table I.

**Scaling ZCCloud:** Significant computing facilities are of modest scale compared to modern wind farms. Our **small** scale system of 2.2 petaflops and **medium** scale system of 4.4 petaflops require below $1.2MW$ and thus could be easily placed beneath a single modern $2MW$ turbine – the dominant size being deployed in commercial farms today [**?**] (see Figure 4) Our **Large** system is still modest in size. A 42-container system with 45.5 petaflops capacity would require 2 containers at the bases of 21 wind turbines – a small fraction of a modern commercial wind farm. For example, the Twin Groves farm [**?**] includes 240 turbines, each $1.65MW$ for peak generating capacity of $398MW$. Even a 3% curtailment can power significant computing capacity (our **Large** system of 45.5 PF is only $12.4MW$). From our 2.2 petaflop ZCCloud building block, replication to achieve 100-petaflop systems is straightforward. One such would cover only 18.9% of a large wind farm such as Twin Groves, and there are dozens of such facilities in the MISO region, so total capacity exceeding exaflops is possible.

ZCCloud offers intermittent and variable capacity based on the availability of stranded power. Volatility tied to wind power will support continuous availability from hours (overnight) to days, due to the change in weather patterns (see Figure 5). Volatility tied to solar stranded power appears to be likely tied to regular daily and weekly cycles, but may involve shorter periods. Commercial uses and markets for volatile computing resources exist. Large-scale cloud provides, like Amazon AWS, offer "spot instances" – rentals at a bid "spot price" that are terminated when the market (or perhaps the provider) decides that they should be reclaimed.

| | Description | Performance | Power |
|---|---|---|---|
| Node | Dual sockets, Intel E7-8890 v3 CPUs | 4,838 GFLOPS | 0.66 kW |
| Rack | 14 nodes per rack | 67.732 TFLOPS | 9.24 kW |
| Container | 16 racks per container | 1.083 PFLOPS | 295.68 kW |

Table I
ZCCLOUD COMPUTING CONTAINER SKETCH DESIGN

In short, they can be revoked at any time, yet are deemed useful by a large user community [**?**]. Unlike these other volatile cloud computng rentals, however, ZCCloud volatility results from the fluctuations in available power and not market or other commercial forces. Thus it is possible to offer more reliable minimum guarantees of service (compared to current spot-market offerings) in the form of Service Level Agreements (SLAs). In short, usage will resemble Amazon's spot-instance facility but adding a guarantee of minimum time to "eviction" and subsequent spot-instance termination. Extensions such as Amazon Spot Fleet's, combining sets of these systems (or their virtualization in instances) also make sense.

We plan to enhance the capabilities of ZCCloud through the construction and operation of a series of software and hardware prototypes. These prototypes will demonstrate the economic benefits of ZCCloud, create a volatile computing resource - demonstration to application users -, and enable advanced research to improve service-level agreements (SLA's) to increase the value of the delivered computing services.

## V. INITIAL EVALUATION

To assess the utility of a Zero-Carbon Cloud HPC computing resource, we have performed a series of simulation using over 12 months of job traces from the Argonne Leadership Computing Facility's Mira system [**?**], considering a number of different stranded power scenarios. Using a system with two times the hardware resources, but only intermittent power of 8 hours/day (33% duty factor), our preliminary results show:

1) >30% of the jobs experience comparable or better turnaround time
2) the largest jobs experience improved turnaround time
3) Simple rules can identify which jobs will benefit, making ZCCloud useful as a complementary resource to traditional HPC platforms.

In short, we are encouraged that ZCCloud is a promising approach to create a new class of HPC computing resources – supporting new capabilities and promsing for cost-effectiveness.

## VI. SUMMARY AND DISCUSSION

We have described the Zero-Carbon cloud concept, describing the basic elements of exploiting stranded power and light-weight physical infrastructure. Our initial design and evaluation suggests the approach is promising – with the potential to dramatically reduce computing costs, and create a complementary capability to traditional approaches. We look forward to exploring the ZCCloud concept more deeply in the future, including open challenges and potential directions, including 1) rigorous simulation studies showing the benefits and exploring the huge configuration space of ZCCloud systems, 2) a detailed design and demonstration of the ZCCloud system, 3) exploring the addition of limited energy storage, 4) exploring opportunities of geographic distribution, including job migration, and many more.