# A First Look at Thermal Attacks in Multi-Tenant Data Centers

Zhihui Shao
UC Riverside

Mohammad A. Islam
UC Riverside

Shaolei Ren
UC Riverside

## ABSTRACT

This paper focuses on the emerging threat of thermal attacks in a multi-tenant data center. It discovers that a malicious tenant (i.e., attacker) can inject additional thermal loads beyond the shared cooling system capacity, thus resulting in overheating and possible system downtime. Importantly, the attacker can launch thermal attacks in a stealthy manner by discharging batteries inside its servers and still keeping its power drawn from the data center power distribution system under its subscribed capacity.

## 1 INTRODUCTION

Be they megawatt warehouses or micro-scale server clusters distributed at Internet edges, data centers have been so tightly integrated with our daily life that no business can run without them. Naturally, ensuring data center availability is extremely crucial.

While the importance of securing servers and networks is well recognized, recent research has also highlighted that securing data center physical infrastructures (e.g., power distribution and cooling systems) is equally, if not more, important. For example, due to the common practice of oversubscription for increasing utilization, data center power infrastructures are highly vulnerable to power attacks: when benign power demand is high, malicious power loads can be injected by running computation-intensive workloads to overload the total power capacity [1, 2]. Despite infrastructure redundancies, power attacks lead to severe performance degradations and even catastrophic data center-wide power outages. Importantly, the cost for launching power attacks only amounts to a negligible fraction of the multi-million dollar loss incurred by a victim data center [1]. In addition, malicious thermal attacks can also be launched in a similar manner to overstress the data center cooling system that is crucial for removing server heat and maintaining system uptime [3]. Specifically, because almost all server power is converted into heat, excessive server power loads injected by attackers can result in the entire data center's melt down.

In contrast with the prior research [1–3], this paper focuses on the emerging attack vector of thermal attacks in a multi-tenant data center, an important data center segment that serves almost all industry sectors (e.g., even Apple houses 25% of its servers in a multi-tenant data center). A multi-tenant data center, also called colocation, is a shared facility that provides conditioned power and cooling to physical servers owned by different organizations each viewed as a tenant, while the data center operator is only responsible for physical infrastructure support without controlling tenants' servers.

For safety and reliability, the data center operator has power meters to monitor each tenant's server power usage (also equivalently, the thermal load or cooling demand) on a per-rack or even per-server basis. This ensures that each tenant's power stays below its subscribed power capacity at all times, and non-compliance can result in warnings and/or involuntary power cuts. By doing so, the operator also keeps the total thermal load below the cooling system capacity, which is sized based on power infrastructure capacity. Thus, in today's multi-tenant data centers, the operator's power meter is essentially *dual* purposed: it monitors both power and thermal loads. Nonetheless, we discover that this practice of load monitoring creates a significant vulnerability to *behind-the-meter* thermal attacks — malicious and stealthy thermal loads that are not monitored by the operator's power meters.

Concretely, stealthy thermal loads can be injected to create cooling capacity overloads with the assistance of batteries placed inside a malicious tenant's (i.e., attacker's) servers. In recent years, distributed batteries inside physical servers have been increasingly adopted (e.g., in Google's servers) as a cost-effective backup power solution. Moreover, tenants in a multi-tenant data center can also install their own batteries inside their physical servers to lower data center leasing costs. Nonetheless, the presence of batteries inside servers also means a significant threat to data center availability: while the operator's power meter can monitor the power drawn from the data center power distribution system, its reading may *not* reflect the actual server power usage or thermal load. For example, the thermal load is actually higher than the power meter reading (which indicates the amount of power taken from the data center power distribution system) if batteries inside servers are discharged as a supplemental power source for servers, and vice versa.

An attacker can exploit the emerging architecture of distributed batteries inside servers to generate an additional thermal load (i.e., thermal attack) without being monitored by the data center operator's power meters. Specifically, when benign tenants are using a high power and producing a large thermal load, the attacker can discharge batteries inside its servers and run its servers at a power level beyond its subscribed power capacity to inject thermal attacks, overloading the data center's cooling system capacity. This not only affects benign tenants' server reliability in the long term, but also potentially makes the data center overheat and results in system shutdown. Importantly, this is done in a *stealthy* manner: the additional thermal load for attacks is supported by internal batteries and the amount of power that is actually drawn from the operator's power distribution infrastructure is still below the capacity subscribed by the attacker.

To our knowledge, while batteries are commonly used for shaving peak power, our work is the first to leverage batteries for an adversarial purpose — thermal attacks in a multi-tenant data center.

## 2 DATA CENTER COOLING SYSTEM

There are three different notions of temperature in a data center: server *inlet* temperature (i.e., temperature of cold air entering a server), server *internal* temperature (e.g., CPU temperature), and server *outlet* temperature (i.e., temperature of hot air exiting a server). Server inlet temperature is the lowest and baseline, whose increase will lead to increases in server internal and outlet temperatures. Server outlet temperature is typically elevated by 10 ~ 20°F

compared to the inlet temperature, while server internal temperature is the highest and regulated by servers' internal fans.

While various cooling methods (e.g., computer room air conditioner, chiller, and "free" outside air cooling) can be employed depending on the size and climate condition of data centers, they all ensure that the server inlet temperature stays below 81°F as recommended by ASHRAE. As almost all server power is converted into heat, cooling system capacity is commonly measured in kilowatt (kW) and sized based on the provisioned power infrastructure capacity. If the thermal load is below the cooling capacity, all server inlet temperatures can be conditioned below 81°F. Under this situation, the server internal temperature can be well regulated below the safety threshold by servers' cooling fans, even when servers are running at the maximum power. Nonetheless, when excessive server heat is generated beyond the available cooling capacity, the temperature of cold air supplied by the cooling system will quickly increase, leading to an elevation in the server inlet temperature and possible data center overheating.

## 3 THREAT MODEL AND RESULTS
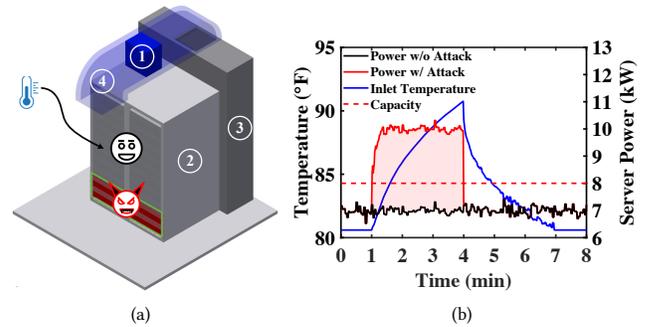
### 3.1 Threat Model

We consider a malicious tenant (attacker) housing its physical servers in a multi-tenant data center. Except for its intention to overload the shared cooling capacity and cause data center overheating, the attacker behaves normally and follows all the constraints on server power usage set by the data center operator.

The attacker subscribes a certain amount of power capacity from the data center operator, but its servers can consume more power than the subscribed capacity. The attacker has batteries inside its servers and can discharge them as a supplemental power source in order to run servers at the maximum power for stealthy thermal attacks without drawing more power from the operator's power distribution system than its subscribed capacity. This is not a restrictive assumption, because distributed batteries inside servers have been increasingly adopted in practice (e.g., in Google's servers). Moreover, recent research has suggested that tenants in a multi-tenant data center also install their own batteries inside their physical racks/servers to lower their data center leasing cost.

Finally, instead of launching random thermal attacks that may drain batteries at wrong times, we assume that the attacker can time its attacks when benign tenants have a high thermal load, which is more likely to overload the shared cooling capacity. This can be achieved by estimating the aggregate thermal load through physical side channels (e.g., acoustic signal [1]).

### 3.2 Preliminary Results

As illustrated in Fig. 1(a), we consider a multi-tenant edge data center, an increasingly more popular type of data center hosting latency-critical workloads such as assisted driving. The data center houses two server racks with a total power demand of 8kW, and the cooling system can supply cold air at 81°F for a total thermal load of up to 8kW. While an attacker only subscribes 2kW of capacity, its servers can use up to 5kW (i.e., generating 5kW thermal load), out of which 3kW is supplied by batteries stored inside the servers. For the best cooling efficiency, the data center also has hot isle containment to prevent hot air from mixing with cold air, and hence



Figure 1: (a) Layout of a multi-tenant edge data center. ① Air conditioner. ② Server racks. ③ Return air duct. ④ Supply air duct. (b) Power and inlet temperature trace during an attack.

the cold air temperature supplied by the cooling system is almost identical to the temperature at all server inlets. We perform widely-used computational fluid dynamics (CFD) analysis to simulate the temperature dynamics inside the data center.

Fig. 1(b) illustrates the server inlet temperature trace during a thermal attack. When the total server power (including the attacker's) and hence thermal load are below the capacity, the server inlet temperature can be well-conditioned at 81°F. Nonetheless, after a 3-minute thermal attack that overloads the cooling capacity, the server inlet temperature rises to nearly 92°F, which is even higher than the allowed temperature limit set by ASHRAE for enterprise servers. This can significantly damage server reliability and even result in data center overheating and downtime incidents, highlighting the danger of thermal attacks. Even when a thermal attack successfully results in a cooling capacity overload, it may not always drive up the server inlet temperature to as high as 92°F due to limited battery capacities and/or drops in benign tenants' thermal loads. Nonetheless, as shown in Fig. 1(b), overloading the cooling capacity for even less than 1 minute can result in a noticeable server inlet temperature increase. Such frequent temperature variations over a long term can greatly impact servers' lifetime.

Note that, although the data center operator has temperature sensors and can detect server inlet temperature increases, it cannot precisely locate the source of excessive heat (i.e., attacker). This is because the cold air supplied by the overloaded cooling system is already hotter than normal and has the same temperature as all server inlets due to hot isle containment. Furthermore, like as any benign tenants, the attacker's power drawn from the shared power distribution system is still below its subscribed capacity.

*In our future work, we will study how to utilize limited battery capacity for thermal attacks and also develop defense mechanisms.*

## REFERENCES

[1] M. A. Islam, L. Yang, K. Ranganath, and S. Ren, "Why some like it loud: Timing power attacks in multi-tenant data centers using an acoustic side channel," in *SIGMETRICS*, 2018.
[2] Z. Xu, H. Wang, Z. Xu, and X. Wang, "Power attack: An increasing threat to data centers," in *NDSS*, 2014.
[3] X. Gao, Z. Xu, H. Wang, L. Li, and X. Wang, "Reduced cooling redundancy: A new security vulnerability in a hot data center," in *NDSS*, 2018.