# RESEARCH INTERNSHIP REPORT

## Integrating Pre-Cooling of Datacenter operated with Renewable Energies
## IRIT - Paul Sabatier University

25 March - 19 July 2019

—

Maël MADON X2016
Supervisor: Jean-Marc PIERSON
Referent: Claudia D'AMBROSIO

ÉCOLE POLYTECHNIQUE

IRIT Institut de Recherche en Informatique de Toulouse

# Déclaration d'intégrité relative au plagiat

Je soussigné MADON Maël certifie sur l'honneur :

1. Que les résultats décrits dans ce rapport sont l'aboutissement de mon travail.

2. Que je suis l'auteur de ce rapport.

3. Que je n'ai pas utilisé des sources ou résultats tiers sans clairement les citer et les référencer selon les règles bibliographiques préconisées.

Mention à recopier :

*Je déclare que ce travail ne peut être suspecté de plagiat.*

Date :                                    Signature

# Abstract

The energy consumption of the Information Technology (IT) sector is increasing at an alarming rate. In particular, as the amount of data processed every day continues to explode, it can be seen that the relative contribution to greenhouse gas emissions of data center supply in all IT categories is increasing. In this context, various studies are looking into the possibility of supplying these centers with renewable energy. The approach presented in this report targets the cooling system of a data center that is locally supplied with electricity from renewable sources. Our goal is to compensate the variability of renewable production. We are taking advantage of overproduction hours to lower the cooling system temperature set point in order to save electricity when the sources no longer produce. We formulated our problem as a Mixed Integer Linear Programming optimization and we simulated the behavior of a data center equipped with a temperature control that adapts to production curves. This study forecasts savings of around 6% on the energy bill related to cooling using what will be called the pre-cooling technique.

# Résumé

La consommation énergétique du secteur du numérique augmente à une allure inquiétante. En particulier, la masse de données traitées tous les jours ne cessant d'exploser, on observe que la part des contributions aux émissions de gaz à effet de serre due à l'alimentation des centres de données dans l'ensemble des catégories du numérique prend de plus en plus d'ampleur. Dans ce contexte, différentes études se penchent sur la possibilité d'alimenter ces centres avec des énergies renouvelables. L'approche que nous présentons dans ce rapport consiste à gérer intelligemment le système de refroidissement d'un centre de donnée alimenté localement en électricité par des sources renouvelables afin de compenser l'intermittence de la production. L'idée est de profiter des heures de surproduction pour baisser la consigne de température du système de refroidissement afin d'économiser de l'électricité au moment où les sources ne produisent plus. Nous avons formulé notre problème sous forme d'optimisation linéaire et avons simulé le comportement d'un centre de donnée muni d'un contrôle de température s'adaptant aux courbes de production. Cette étude prévoit des gains de l'ordre de 6% sur la facture énergétique lié au refroidissement en utilisant le pré-refroidissement.

# Acknowledgements

First and foremost, I would like to thank very warmly my supervisor Jean-Marc. I thank him for his advice, his trust, his humor and his friendship. He supervised me on a weekly basis while leaving me completely free to try my ideas and learn from my errors. He never let me alone in front of a difficulty and my e-mails to him always found an quick answer. Thanks to him I really discovered the research world and I enjoyed my internship .

Thank you also to all the SEPIA team full of joy and sympathy: Amal, Patricia, Paul, George, Guilherm, François. This team creates an pleasant atmosphere where no one takes themselves seriously. I enjoyed a lot participating to the Green Days with them.

I also owe a lot to the PhD students: Léo, Gustavo, Chaopeng, always with good boosting advices. Thank you also to those from TRACE team who share the same floor: Thomas, Hugues, Pascal, Armelle. I was able to find in each and every people mentioned resources to move forward.

To finish, my last thanks will go to the other interns: my friends and colleagues Florent, Emmanuel, Marie-Léontine, Emmanuel, Axel, Lucas and Quentin.
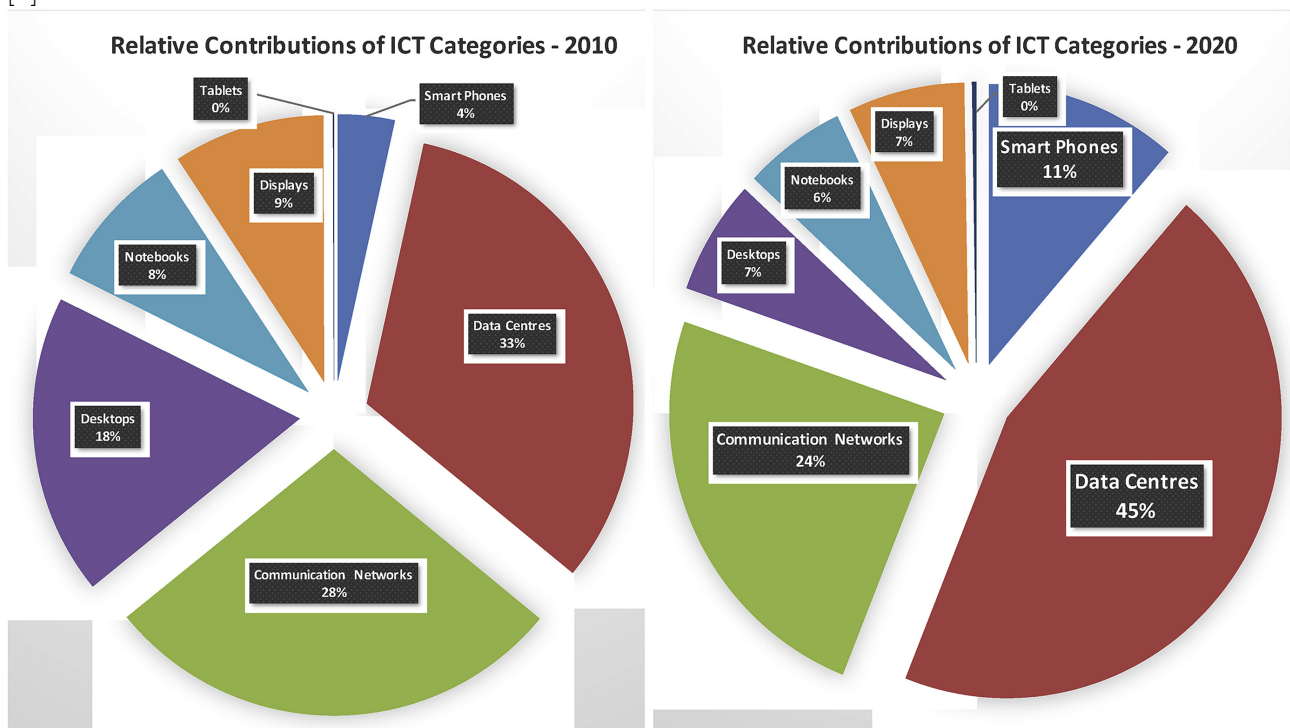
# Contents

# Introduction

**IT impact on the environment**    If awareness is growing on the impact of humans over their environment, the precise influence of each sector is not well known. The energy footprint being mostly indirect in the Information Technology (IT) sector, we tend to be unaware of its negative impact. The Shift Project, a French think tank, published in October 2018 a report on the environmental impact of digital technology [7]. Their work, based on the one of Andrae *et al.* [1] reveals that IT energy consumption in the world is increasing by about 9% per year (period from 2015 to 2020) reaching 2.7% of the world total energy consumption in 2017 and 3.3% in 2020. It has to be noticed that these figures include both the use phase and the production phase of IT devices. According to their estimate, this consumption will represent in 2020 *4% of the global greenhouse gas emissions*. This proportion and the rate of its increase leads us to consider this problem as a major concern.

Figure 1: Relative contribution of each IT category in 2010 and 2020. *Source: Belkhir et al.* [2]

**Data centers** As we mentioned before, a key aspect of IT environmental footprint is that it is mainly indirect *i.e.* that it does not only come from the electricity powering our IT devices, but also from their production or the hosting of their data. Figure 1 shows the distribution of the IT greenhouse emission footprint by IT category. It can be seen that an increasing part of this footprint relates to data centers, which reaches almost half of the IT environmental impact in 2020.

**Internship background** It is from this last source of savings that our study claims to draw on. I am doing my third year research internship at Toulouse Computer Science Research Institute (IRIT) of Paul Sabatier University, in a team that focuses its research on energy and resource efficiency for distributed operating systems. My internship is part of the ANR Datazero project (`datazero.org`) started in 2015 in collaboration with FEMTO-ST (Franche Comté), LAPLACE (Toulouse) and financed by EATON (France). This project aims at studying a data center cut from the power grid and self-supplied in renewable energy. The electrical side is composed of photo-voltaic panels, wind turbines, batteries and fuel cells. The sources and storage units engagement is managed by the Power Decision Module (PDM). On the IT side, service placement and scheduling are done by the IT Decision Module (ITDM) to adapt to the available power. Thanks to negotiation, the PDM and ITDM are able to ensure functioning of the data center with a good quality of service.

**Pre-cooling** As well as the ITDM into Datazero is moving IT tasks to schedule them when power is available, why don't we adapt cooling loads to match renewable energy production? According to the 2019 data center survey from Uptime Institute [13], the average Power Usage Efficiency (PUE) in 2019 is 1.67. The PUE is the ratio of total amount of energy used by a data center to the energy delivered for actual IT equipment. The ideal PUE is 1 and a PUE of 2 means that half of the energy in the data center is not used directly for computing but for cooling, lighting, monitoring, etc.. In general, as one might expect, the main other energy expense in a data center is cooling. In the last decades, lot of work has been done to lessen cooling costs by choosing a better location for the data center or better cooling device, by working on the servers architecture (the well known hot aisle / cold aisle architecture) etc.. The proposed approach doesn't target a PUE reduction but study how a given installation can adapt to power supply only by modifying its cooling strategy. The American Society of Heating, Refrigerating and Air-Conditioning Engineers (ASHRAE) issued in 2016 guidelines for data center equipment [24]. The idea about pre-cooling is to use the full temperature range recommended by the ASHRAE: lessening the temperature when power is available to be able to cut the cooling device and let the temperature increase again when production is low.

**Outline** The goal of the present report is to propose a pre-cooling strategy based on weather prediction and linear solving then evaluate its ability to increase self-supply in a data center. Part 1 provides a brief overview of the state of the art of competing approaches. We will then detail in part 2 the choice we made to model a data center supplied with renewable energies from the power sources to the cooling system. We developed a simulation of the thermal behavior of this data center with which we were able to test different pre-cooling strategies. The simulation

is explained part 3. The results of the study will be given part 4 as well as some discussions about them. Finally, we will conclude part 5 by explaining the limitations of our approach and the options considered for future work.

# Part 1

# Background

## 1.1   Data centers supplied with renewable energies

In the last decade, there has been a lot of effort made to minimize greenhouse gas emissions in the IT sector. According to the 2017 Greenpeace report [5], the number of IT companies committed to a 100% renewable supply is growing. The digital giants are ranked high in Greenpeace company scoreboard with respectively 83%, 67% and 56% clean energy index for Apple, Facebook and Google. The approach is generally huge investments in green power farms or contracts with electricity providers. Each and every company can today choose "green cloud solutions" in the market to host their servers (see for example [22, 21]).

In the academic field, many work has been done to integrate co-located renewable sources in the power supply system of data centers. As always with renewable energies, the difficulty stands in the management of the variability of production. To handle this, data centers have a great potential: IT loads can be scheduled in a smart way or even be migrated to another data center located in an area where production is better. Goiri *et al.* [8] made a pioneer paper tackling the first point. They have built Parasol, a real prototype green data center powered with solar panels and grid ties as well as storage units. They developed for it GreenSwitch, a dynamic scheduler and source selector. The workloads they study were MapReduce workloads with deferable and non deferable tasks. GreenSwitch used a formulation as a Mixed Integer Linear Programming problem integrating both workload scheduling and energy management constraints. Their results are very promising as they are able to save from 40% to 100% energy from the grid, depending on the weather conditions. GreenSwitch achieves it by moving deferable IT tasks, charging the batteries with the power surplus and discharging them when production is low. In the Datazero project we introduced before, Grange *et al.* [9] and Caux *et al.* [4] proposed more decentralized approaches where scheduling is treated apart from power management. They consider an utility function given by a power decision module and test scheduling heuristics for IT tasks with due dates. Grange *et al.* achieved in their experiments a reduction of grid consumption ut to 49%.

Even more holistic approaches have been studied. Liu *et al.* [17] for example took into account renewable supply, dynamic pricing and IT workload planning as well as cooling supply in their work on data center operations. For cooling they consider free cooling, a technique that consists in blowing outside air into the building when its temperature and moisture are at an appropriate level. A concurrent approach is the one of Habibi Khalaj *et al.* [10] which is more focused on renewable production and cooling. They studied 42 places to assess their production and free cooling potential. They also give results on optimal sizing for hybrid power infrastructure. Other works study geographical allocation of tasks on data center having their own on-site or off-site renewable sources [12, 15, 14].

## 1.2   Pre-cooling in the literature

As we mentioned before, the present paper focuses only on cooling strategy and more precisely on the ability to store energy in form of cold inside the data center thermal mass. Few research has been done on this subject, as far as we could find. Zhang *et al.* [25] present in their paper TEStore, a cooling strategy exploiting thermal and energy storage to cut energy bills in data centers. They assume a varying price for electrical power during the day and study three forms of energy storage: (i) ice or water-based thermal tanks, (ii) UPS batteries, (iii) building thermal mass. When power price is low, energy can be used to charge the batteries, store cold in the thermal tanks or pre-cool the building. When power price is high, the energy stored is used to cool the infrastructure. Their results show that almost 85% of the cost saved by TEStore are introduced by exploiting the thermal tanks on the long time scale. UPS batteries can be used in the middle term and precooling potential in their model is relatively low: it represents a few percents of the total saving, with only the ability to pre-cool a few minutes before the peak price.

However, Lukawski *et al.* claim in [18] better results for pre-cooling. They are using this mechanism as a demand response mechanism for reducing coincident peak loads in the power market. The idea is to be able to cut the cooling device of the data center to save energy during a 15 minutes peak price. The model they describe, which was tested experimentally, can easily achieve it. The reasons for this difference with TEStore are the following: (i) the data center considered has a low power density (below $500W/m^2$) and (ii) it has no server air containment systems. In fact, Lukawski *et al.* highlight the close link between pre-cooling potential and floor utilization: facilities with high power density do not have sufficient thermal storage capacity to provide extended demand response time. As a result, the candidates for using a pre-cooling mechanism are small data centers such as computer rooms or data centers for telecommunications, companies or universities.

We found Zhang *et al.* and Lukawski *et al.* approaches interesting but we noticed that they don't deal with renewable energies but only financial considerations. Our work will study the synergies between renewable production and pre-cooling possibilities. Li *et al.* already covered it in a very recent paper [16] with a holistic thermal-aware workload and cooling management

for the maximization of renewable energy sources. They consider batch (deferrable) jobs and interactive (non-deferrable) jobs. They compare four methods: (i) static method where tasks are executed as soon as they arrive; (ii) load balancing distribution over time where batch-type jobs are distributed evenly over multiple time slots; (iii) best effort strategy where the batch jobs are scheduled according to the predicted power generation amount; (iv) thermal-aware workload management where not only IT but also cooling power consumption are taken into account to decide. In the last strategy, if IT doesn't consume everything, renewable surplus is used to cool the room. Pre-cooling plays here again a secondary role. The results of their comparison is that they achieve to reach more than 98% solar utilization with the last strategy with less average waiting time for batch-type jobs than for the second and third strategy. However this strategy seems quite aggressive and the goal of maximizing renewable power utilization at any cost could be discussed.

In our work we will isolate only the effects of pre-cooling with a more accurate time step, allowing ourselves to sell surplus of energy on the grid.

# Part 2

# Data center modeling

## 2.1  Description of the model components

In the chosen model, the data center is supplied with renewable energies. The data center is an isolated building containing a low density of servers. It can represent a small telecommunication data center, a computer room, internal servers of a company... The power supply system is composed of (i) a photo-voltaic source delivering at time $t$ a power $P_{pv}(t)$ (in Watt), (ii) a wind farm delivering $P_{wind}(t)$, (iii) an access to the power grid $P_{grid}(t)$.

This energy is used to power IT equipment and a cooling system, denoted by $P_{IT}(t)$ and $P_{cool}(t)$ respectively. Other energy usage such as lighting, power distribution units or auxiliary facilities is ignored in this model. Thus, if we denote by $P_{prod}(t)$ the total renewable production and $P_{conso}(t)$ the total power of the data center, we have:

$$P_{prod}(t) = P_{pv}(t) + P_{wind}(t) \tag{2.1}$$
$$P_{conso}(t) = P_{IT}(t) + P_{cool}(t) \tag{2.2}$$

Electricity from the grid is used only when necessary, *i.e.* when $P_{prod}(t) < P_{conso}(t)$.

We studied a geothermal heat pump as a cooling device. This system is more efficient than a conventional air-cooled computer room air conditioning (CRAC) as the heat sink in the ground stays at a more stable temperature than the outdoor air. Indeed, ground temperature remains around 13°C throughout the year, well below the ASHRAE recommended maximum temperature.

## 2.2  Power sources

In order to obtain realistic data for renewable production, we used solar irradiation and wind speed records downloaded from Solar Radiation DAta (SODA) website [23]. In compliance with the DataZero project, we used data from Belfort (France) and the year 2006. We adopted

ÉCOLE
POLYTECHNIQUE

as well the same equations to model the power output (Haddad *et al.* [11]). The equations are 2.3 and 2.4:

$$P_{pv}(t) = A_{pv} \times \eta_{pv} \times I(t) \tag{2.3}$$

where $I$ is the irradiation (in W/m$^2$), $A_{pv}$ the area of the panels (m$^2$) and $\eta_{pv}$ their efficiency.

$$P_{wind}(t) = A_{wt} \times \eta_{wt} \times \begin{cases} P_r \frac{v^3(t) - v_{ci}^3}{v_r^3 - v_{ci}^3} & \text{if} \quad v_{ci} < v(t) < v_r \\ P_r & \text{if} \quad v_r < v(t) < v_{co} \\ 0 & \text{else} \end{cases} \tag{2.4}$$

where $v(t)$ is the wind speed (m/s), $A_{wt}$ the total swept area by the blades, $\eta_{wt}$ the wind farm efficiency, $P_r$ the nominal power and $v_{ci}$, $v_r$ and $v_{co}$ wind speed thresholds.

## 2.3 Workload generation

To model our data center IT load over time, we made use of a workload generator developed by Da Costa *et al.*. The generator comes from the analysis of the Google Cluster Workload Traces and its implementation is described in [6]. It is a Python script which takes as input a parameter $n$ as well as a parameter $d$ representing the average number of tasks processed by the data center each hour. It outputs a trace of $n$ points similar to a standard Google trace, each point indicating the number of tasks in process at time $t$. We chose a low value for $d$, to imitate a small data center with a low task arrival. We then assumed that each task requires the same amount of power to run and thus normalized the output to have an average power of $\overline{P_{IT}}$, a constant chosen to match the sizing of our model. This IT load model is obviously approximate but sufficient for our study.

## 2.4 Indoor temperature

To simulate the thermal behavior of the data center building, we adopted a discrete time model used in similar works [19, 18]. Time between 0 and $t_{end}$ is discretized into $K$ periods of time step $h$. The system we study is the whole building represented as a homogeneous thermal mass. The heat transfers we consider are the following:

- the IT thermal load $Q_{IT}(t)$

- the heat dissipated by the cooling device $Q_{cool}(t)$

- the heat transfers through the walls of the building

The equation modeling the thermal behavior comes from the application of the second principle of the thermodynamic on our system. It has been tested experimentally by Lukawski *et al.* in [18]:

$$T(t + h) = a.T(t) + (1 - a)[T_{ext}(t) + R(Q_{IT}(t) - Q_{cool}(t)] \tag{2.5}$$

where $T_{ext}$ is the ambient temperature, $R$ the thermal resistance (in °C/W) and $a$ a dimensionless parameter capturing the thermal inertia of the building.

$$a = exp\left(-\frac{h}{C.R}\right) \tag{2.6}$$

with $C$ the thermal capacitance (in J/°C).

According to [20], heat dissipation from IT can be approximated by its power consumption, as the power transmitted by computing or other information technology equipment through the data lines is negligible:

$$Q_{IT}(t) = P_{IT}(t) \tag{2.7}$$

## 2.5   Cooling system

Once more, the cooling model has been adopted from Lukawski *et al.* and validated in [26]. The data center acts as a thermostatically controlled load (TCL): the heat pump maintains the temperature within a dead band by switching on and off. The regulation is done thanks to three available stages for the pump: off, nominal rate or double rate. Heat dissipated by the cooling device can be written as follows:

$$Q_{cool}(t) = r(t) \times Q_{cool,nom} \times (a_C \cdot T(t) + b_C) \tag{2.8}$$

where $Q_{cool,nom}$, $a_C$ and $b_C$ are parameters depending on the heat pump model and $r(t)$ is the variable representing the current stage of the pump ($r(t) \in \{0, 1, 2\}$).

The thermostatic control is expressed by:

$$r(t+h) = \begin{cases} 0 & \text{if} & T(t) \leq T_{min} \\ 2 & \text{if} & T(t) \geq T_{max} \\ 1 & \text{if} & T_{min} < T(t) < T_{med} & and & r(t) = 2 \\ 1 & \text{if} & T_{med} < T(t) < T_{max} & and & r(t) = 0 \\ r(t) & \text{else} \end{cases} \tag{2.9}$$

Finally, power consumption of the cooling system can be computed using the coefficient of performance (COP) with the equations 2.10 and 2.11:

$$P_{cool}(t) = \frac{Q_{cool}(t)}{COP(t)} \tag{2.10}$$

$$COP(t) = COP_{nom} \times (a_{COP} \cdot T(t) + b_{COP}) \tag{2.11}$$

where $COP_{nom}$, $a_{COP}$ and $b_{COP}$ are here again specific to the pump.

Table 2.1: Values of the model's parameters

| Symbol | Parameter | Value |
|---|---|---|
| $\overline{P_{IT}}$ | Average IT power | 14040 [W] |
| *Power sources* | | |
| $A_{pv} \times \eta_{pv}$ | Solar constants | 50 [m$^2$] |
| $A_{wt} \times \eta_{wt}$ | Wind constants | 30 [m$^2$] |
| $v_{ci}$, $v_r$, $v_{co}$ | Wind speed thresholds | 4, 10, 30 [ms$^{-1}$] |
| $P_r$ | Wind turbine nominal power | 1800 [W] |
| *Thermal model* | | |
| $R$ | Thermal resistance | $4.67 \cdot 10^{-3}$ [°CW$^{-1}$] |
| $C$ | Thermal capacitance | $15.76 \cdot 10^6$ [J°C$^{-1}$] |
| $h$ | Time step | 300 [s] |
| $Q_{cool,nom}$ | Heat pump nominal cooling capacity at stage 1 | 15767 [W] |
| $a_C$, $b_C$ | Heat pump parameters | 0.024, 0.361 |
| $COP_{nom}$ | Heat pump coef. of perf. at nominal conditions | 3 |
| $a_{COP}$, $b_{COP}$ | COP temperature dependency parameters | 0.022, 0.406 |
| $T_{min}$, $T_{med}$, $T_{max}$ | Thermostatic control dead-band | 25.6, 26.4, 27.2 [°C] |
| *Linear solving* | | |
| $p_{buy}$ | Purchase price of electricity on the grid | 0.15 [€/kWh] |
| $p_{sell}$ | Resale price of electricity on the grid | 0.06 [€/kWh] |
| $T_{lo}$, $T_{up}$ | ASHRAE recommended temperature | 18, 27 [°C] |

# Part 3

# Simulating the pre-cooling

We developed in Python a simulation of the data center described in part 2. This entirely modular program allowed us to try different approaches with different data in order to assess the potential and limits of the pre-cooling.
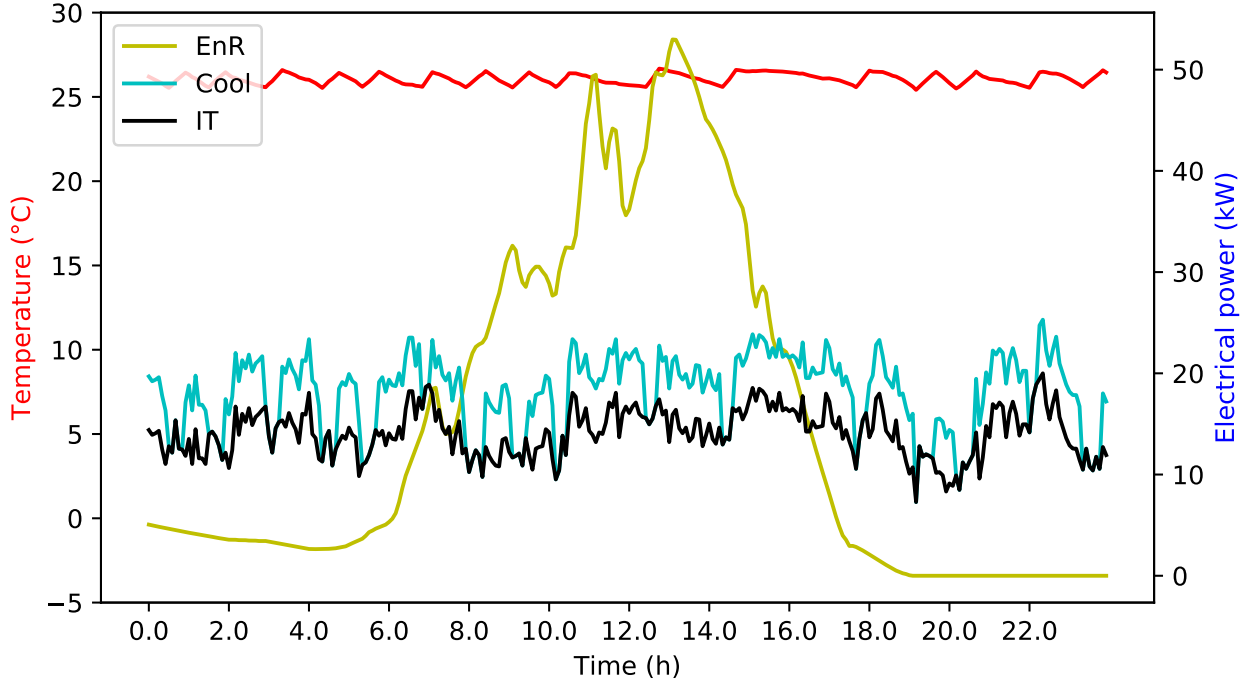
## 3.1  Naive heuristic

The program takes a time window as entry along with the IT load and weather data (ambient temperature, solar irradiation, wind speed) over this period. It then simulates the thermal behavior of the data center over time thanks to the equations previously described. By default, the algorithm has nothing to decide: indoor temperature and heat pump stage are determined at each iteration based on the values of the model variables at the previous iteration (equations 2.5 and 2.9). The program is then able to output the evolution over time of the two computed variables, *i.e.* room temperature and cooling device power consumption. Figure 3.1 shows a standard output of the program, where the cooling strategy is the basic TCL.

On this figure we can notice that the second stage of the heat pump is never used. The temperature can be controlled only by switching between stage 0 and stage 1. Also, looking at the shape of the room temperature curve can help us visualize the thermal inertia of the system. Depending on the ambient temperature, it takes about one or two hours to make the temperature drop by one degree with cooling device at stage 1 and about a quarter to a tenth of this time to regain it with cooling device off. It means that our building has a rather large thermal inertia and enables us to store energy in the thermal mass, which is the basis of pre-cooling. We will see in part 4.3 that this thermal mass storage capacity is closely linked to the density of IT equipment per square meter, relatively low in our model.

The pre-cooling potential stands in the period where production is greater than consumption. We could cool at that time our building more importantly to store electricity surplus in form of cold. Therefore, figure 3.2 displays the behavior of our system with a simple heuristic for the cooling control: as long as renewable production is in surplus, we try to make the best use of this energy by increasing the heat pump stage. This cooling strategy can be expressed as:

Figure 3.1: Standard output of the program for one day. Temperature in red is thermostatically controlled. We see renewable production in yellow, here principally solar. IT load is in black and cooling power is plotted above in blue.
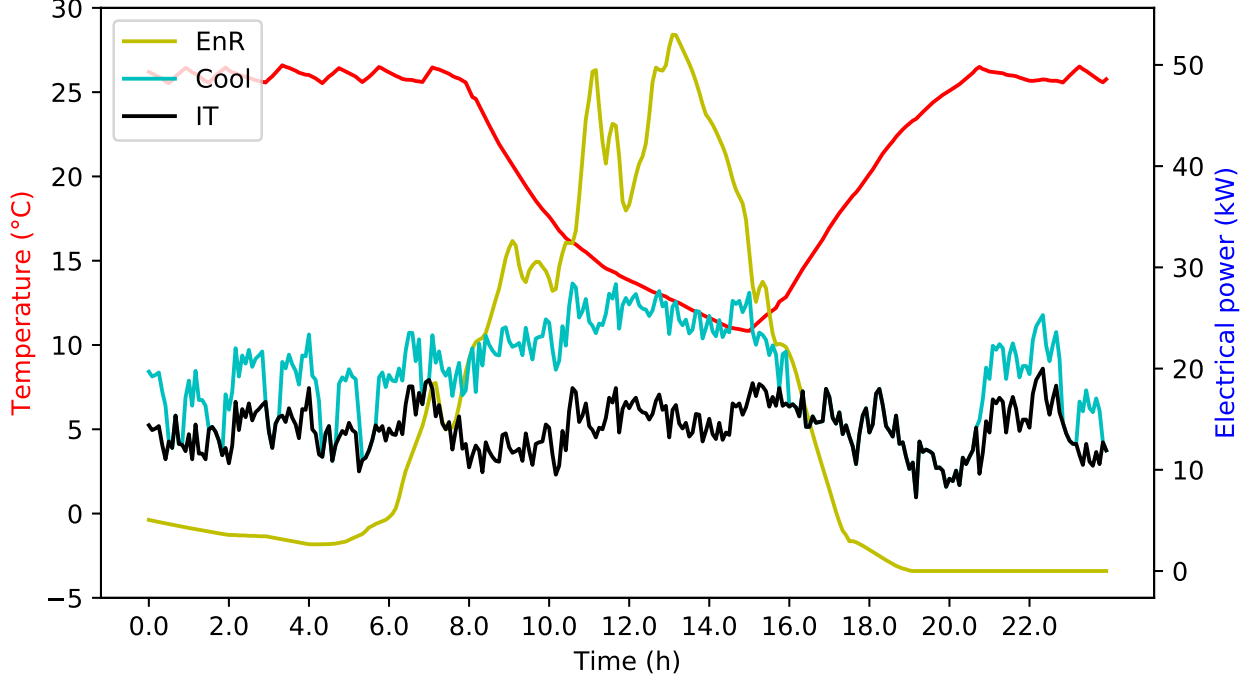


$$r(t) = \begin{cases} r_{TCL}(t) & if \quad P_{prod}(t) < P_{IT}(t) + P_{cool,nom} \\ 1 & if \quad P_{IT}(t) + P_{cool,nom} \leq P_{prod}(t) < P_{IT}(t) + 2P_{cool,nom} \\ 2 & if \quad P_{prod}(t) \geq P_{IT}(t) + 2P_{cool,nom} \end{cases} \quad (3.1)$$

with $r_{TCL}$ being the standard cooling strategy presented in equation 2.9 and $P_{cool,nom}$ the power of one stage of the heat pump at average temperature.

We see in figure 3.2 that we were able to cut off the heat pump during a few hours at the end of the day thus saving electricity from the grid. It has costed on the down side an overconsumption of the cooling system during the day. Overall, we have traded a loss in the total energy used (the area under the consumption profile, in kWh) to the benefit of a better self-supply. However, it is easy to realize that this simple pre-cooling approach is not optimal. Firstly it can make the room temperature fall under the recommended temperature as it is the case in this example. Secondly, cooling at the maximum of our capacity is likely to be a too aggressive management. That is why the purpose of what follows is to find the best time to start pre-cooling or more generally the best pre-cooling policy.

Figure 3.2: Thermal behavior and power consumption of the data center with naive pre-cooling heuristic.



## 3.2 Linear solving

As we saw before, the problem of optimizing the data center cooling policy in order to maximize the direct use of renewable production has only one decision variable: the heat pump stage $r(t)$ at each time step. Modifying its value at a given time will impact the system during the rest of the simulation. Consequently, the best choice can only be made by taking into account all in one the evolution of the exogenous variables over the time window and every possible decision for $r$ at each time step. As it is a non trivial algorithmic problem, we decided to use a linear solver. This section will first explain the choice for the objective function then detail the constraints of our model. We used Gurobi linear solver for our experiments.

### 3.2.1 Objective function

At the beginning we tried a simple objective: minimizing the grid consumption. Grid power can be expressed as the power that is needed by the data center but can not be supplied by the renewable sources.

$$P_{grid}(t) = max\left(0, P_{conso}(t) - P_{prod}(t)\right) \tag{3.2}$$

In our discrete model comes the objective function as follows:

$$\text{minimize} \sum_{k=0}^{K-1} P_{grid}(kh) \tag{3.3}$$

The drawback of this method is that it leads to irrational choices when energies are available. In fact, for the problem as it stands, the only quantity we minimize being grid consumption it is as if the energy coming from renewable sources was totally free. As a result, this objective function is not significantly better than the naive approach as it does not correct the bias we mentioned above. We needed to capture the fact that we also wanted to minimize total electricity consumption. Gurobi allows us to do this by setting multiple objectives with priority between them. We tried it by programming it to find the most energy-efficient solution within a certain tolerance around the solution of the first objective. Quickly though, we changed the approach for one that seemed more appropriate to us.

The approach we chose for the objective function is to introduce a price for the energy self-supplied and the one bought on the grid. This pricing can represent either the intensity of each of these kilowatt-hours in greenhouse gas emissions or their purchase and resale price on the grid. For simplicity, we will talk in the following about purchase and resale prices. Let $p_{buy}$ and $p_{sell}$ be those prices. We assume by common sense that $p_{buy} > p_{sell}$. An electron we use is an electron we cannot sell so it is as if we buy our own production at price $p_{sell}$.

Thus, the objective function comes as:

$$\text{minimize} \sum_{k=0}^{K-1} \left( p_{buy} \cdot P_{grid}(kh) + p_{sell} \cdot P_{self}(kh) \right) \tag{3.4}$$

where we have introduced the self-supplied power $P_{self}(t)$:

$$P_{self}(t) = min\left(P_{conso}(t), P_{prod}(t)\right) \tag{3.5}$$

With such an objective, solving the problem will really find the good trade-off between pre-cooling and overconsumption while allowing us to precisely tune this trade-off with the two prices.

### 3.2.2   Constraints from the physical model

Now that we have formalized what we aim at optimizing, let us see how our model can be translated into constraints. We kept the variables and the equations presented in the modeling part (part 2). In order to take into account the problem as a whole, we solve a system containing for each time step $k$ the equations and inequalities simulating our model. Therefore we will note in the following our variables with $k$ indices rather than as time functions. With this notation, we obtain for each $k$ the following constraints to capture the thermal behavior and cooling power consumption:

$$T_k = \begin{cases} T_{init} & \text{if } k = 0 \\ a.T_{k-1} + (1-a)[T_{ext,k-1} + R(Q_{IT,k-1} - Q_{cool,k-1})] & \text{else} \end{cases} \tag{3.6}$$

$$Q_{cool,k} = r_k \times Q_{cool,nom} \times (a_C \cdot T_k + b_C) \tag{3.7}$$

$$P_{cool,k} = \frac{Q_{cool,k}}{COP_{nom} \times (a_{COP} \cdot T_k + b_{COP})} \tag{3.8}$$

To this we add two equations defining the variables involved in the objective function:

$$P_{grid,k} = max\left(0, P_{cool,k} + P_{IT,k} - P_{prod,k}\right) \tag{3.9}$$

$$P_{self,k} = min\left(P_{cool,k} + P_{IT,k}, P_{prod,k}\right) \tag{3.10}$$

Among the equations presented above, 3.8 and 3.7 are non linear. We linearized the first by neglecting the heat pump $COP$ dependency on $T_k$. Indeed, thanks to the choice of a geothermal source, the slope of this dependency is relatively low. With our numerical parameters, the $COP$ varies from $COP(18°C) = 2.41$ to $COP(27°C) = 3$. After the numerical resolution, this dependency is taken back into account to correct the error introduced in the results and figures we plot. Equation 3.7, however, requires more work. We will detail it in the next section along with the constraints on $r_k$.

### 3.2.3   Constraints on the cooling device

As we saw before, the purpose of using a linear solver is to make it find the best strategy for the cooling device, *i.e.* find the best $r_k$ for all $k$. Thus, the first constraints we add are upper and lower bounds for the room temperature. We will use ASHRAE recommendations for all classes that can be found in their 2015 thermal guidelines [24]: $T_{lo} = 18°C$ and $T_{up} = 27°C$. Hence the set of constraints for all $k$:

$$T_{lo} \leq T_k \leq T_{up} \tag{3.11}$$

Regarding $r_k$, we will define two operating modes: (i) the "free mode" where the solver is free to set $r_k$ evolution as it wishes and (ii) the "TCL mode" where we force the solver to make it work as a thermostatic control. We will translate both modes into constraints later on. It might seem more appropriate to set the solver in free mode for all the duration of the simulation. However we did not for two reasons. First, despite its parameterization, the Gurobi solver was not able to find the optimal solution in a reasonable time for a 24-hour time window problem with a 5-minute time step. It can find an almost optimal solution in a few seconds but struggles to demonstrate its optimality even in more than 20 hours of calculation on a 12 cores i7-8700 machine. This phenomenon is probably due to the impossibility in our problem of cutting branches at the root in the decision tree on the variable $r_k$, for the interdependence reasons mentioned page 19. The second reason why we didn't use only free mode is that the quasi-optimal solutions found by the solver seemed physically unrealistic: in order to save as

much energy as possible, the cooling system ended out switching on and off at each step so that the temperature remains as close as possible to the maximum temperature. As a result, we chose to trigger free mode only if the production is sufficient:

$$
\begin{array}{lll}
\text{if } P_{prod,k} > P_{IT,k} & : & \text{free mode} \\
\text{else} & : & \text{TCL mode}
\end{array}
\tag{3.12}
$$

In the following we detail the constraints for both modes. First of all, $r_k$ being a trinary variable we will represent it without loss of generality as $r_k^1 + r_k^2$ with $r_k^1$ and $r_k^2$ two binary variables and $r_k^2 \leq r_k^1$. As such, we can linearize equation 3.7 using a classical trick in linear optimization. To linearize expression $Q_k^1 = r_k^1 \times T_k$ with $r_k^1$ a binary variable and $T_k$ a continuous variable bounded below by zero and above by $M = T_{up}$, it suffices to add the following inequalities in the model:

$$
\left\{
\begin{array}{lll}
Q_k^1 & \leq & M \times r_k^1 \\
Q_k^1 & \geq & 0 \\
Q_k^1 & \leq & T_k \\
Q_k^1 & \geq & T_k - (1 - r_k^1)M
\end{array}
\right.
$$

By using this trick for $r_k^1$ and $r_k^2$ we obtain the set of linear constraints corresponding to equation 3.7. Free mode then consists of these inequalities plus the lower and upper temperature constraint.

Regarding TCL mode, we will add more constraints to the free mode previously described. The new constraints correspond to the thermostatic control formalized in equation 2.9. We explain in the following the trick we used to linearize. We first introduce $inf_k$ and $sup_k$ two binary variables. We manage to have $sup_k = 1$ if and only if $T_k \geq T_{med}$ by adding this two constraints:

$$
\left\{
\begin{array}{lll}
sup_k & \geq & \frac{T_k - T_{med}}{c} \\
sup_k & \leq & 1 + \frac{T_k - T_{med}}{c}
\end{array}
\right.
$$

where $c$ is chosen big enough to avoid the right terms exceeding 1. In the same way, we force $inf_k = 1$ if and only if $T_k \leq T_{min}$. Thermostatic control for $r_k^1$ can then be expressed as follows:

$$
\left\{
\begin{array}{lll}
r_k^1 & \geq & sup_k \\
r_k^1 & \leq & 1 - inf_k \\
r_k^1 - r_{k-1}^1 & \leq & sup_k + inf_k \\
r_{k-1}^1 - r_k^1 & \leq & sup_k + inf_k
\end{array}
\right.
\tag{3.13}
$$

By using the same trick for $r_k^1$ between $T_{med}$ and $T_{max}$, we obtain the full set of constraints for the TCL mode.

# Part 4

# Results and discussion

## 4.1 Decomposition of an output

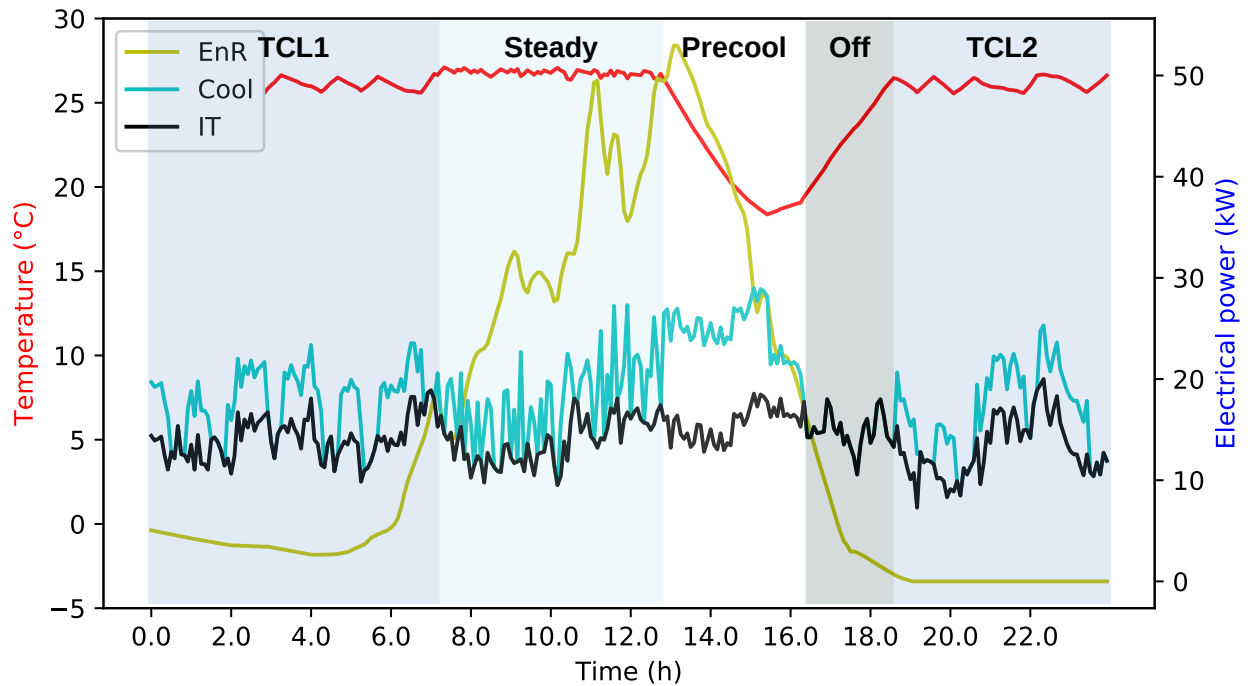Figure 4.1: Cooling strategy calculated by linear solving



Figure 4.1 shows the result of the linear solving for a weather data recorded in Belfort on May 21, 2006. This day corresponds to the day of median photo-voltaic production in May of this year. We used the numerical parameters presented in Table 2.1. For pedagogical and comparison purposes, we cut the graph into 5 phases. The first and the last two phases stand when the renewable production is below the IT power consumption. According to the

constraints previously described, the cooling device is here forced to act as a thermostatic control. In phase 2 and phase 3, the optimization occurs as the linear solver is freer to set the pump stages. We clearly see the pre-cooling happening.

Table 4.1: Decomposition of cost and energy expenses for cooling. The control execution is the standard thermostatic cooling plotted in Fig. 3.1 whereas the optimized execution correspond to Fig. 4.1. Energy is in kWh and cost is in €. The ratios in the last two columns are the relative participation of each phase into the total energy (resp. financial) gap.

| Phase | Control execution | | Linear solving | | Normalized deviation | |
|---|---|---|---|---|---|---|
| | Energy | Cost | Energy | Cost | Energy | Cost |
| TCL 1 | 25.14 | 3.77 | 25.14 | 3.77 | 0.00 % | 0.00 % |
| Steady temp | 22.06 | 1.32 | 21.20 | 1.27 | - 0.92 % | - 0.49 % |
| Pre-cooling | 16.77 | 1.01 | 32.32 | 1.94 | + 16.62 % | + 8.85 % |
| Cooling off | 10.59 | 1.59 | 0.00 | 0.00 | - 11.32 % | - 15.08 % |
| TCL 2 | 18.97 | 2.85 | 18.97 | 2.85 | 0.00 % | 0.00 % |
| **TOTAL** | **93.54** | **10.54** | **97.63** | **9.83** | **+ 4.38 %** | **- 6.72 %** |

To be more precise on the analysis of this output, quantified results for each phase of this execution are given in Table 4.1. These results are compared with the standard output where nothing is done to maximize self-supply. We will call this execution the "control execution". As expected, we can notice that the optimized data center behaves exactly like the control execution in phases TCL1 and TCL2. In fact, room temperature is within the dead band and we did everything to force this behavior and having only one solution in underproduction period. The pre-cooling phase starts around 1:00PM. We clearly see the cooling device going into stage 2 then stage 1 to remain under the yellow curve and room temperature evolving accordingly. In this example, it generates an overconsumption of the cooling device of 15.54 kWh compared to the control execution. It represents 16.62 % of the total energy consumed by the contol execution. The point is that the cost of the energy at this time is lower because it is self-supplied. That is why in financial terms, this phase only generates an extra cost of €0.93, which participates up to 8.85 % in the overall cost difference. This extra cost is more than compensated by the next phase where the cooling is switched off while it costed €1.59 at market price for the control execution.

One can notice that the phase we called "steady" can be seen as non physically realistic. Indeed, as we are in production period, the solution output by the linear solver switches freely between the heat pump stages to stick as close as possible to the maximum temperature. It is a behavior we managed to avoid by forcing TCL in underproduction period but that we cannot avoid here. However, it brings only 0.49 % overall cost reduction, to be compared with the 6.72 % reached.

## 4.2   Discussion about the results

Finally, we reached in this example 6.72 % cost reduction compared to the control execution. As it is mentioned page 20 when choosing the objective function, this cost reduction can not only be interpreted as profit for the data center manager, but also as a diminution in greenhouse gases emission. It all depends on the choice of the prices $p_{buy}$ and $p_{sell}$ and what they represent. This cost reduction has been made to the detriment of an overall energy increase of 4.38 %. What we won is that we consumed more when the energy was available thus coping with variability, one of the big challenge that renewable energies have to face.

This result can be compared with other cost reduction. For example, we computed how much energy was saved by simply setting the temperature one or two degree higher in our data center. We run the standard iterative simulation at different temperatures and reported the results in Table 4.2. We can see that pre-cooling potential in our model, even forgetting about the non realistic "steady" phase, leads to a cost reduction comparable to the one we would have had by augmenting temperature by one degree.
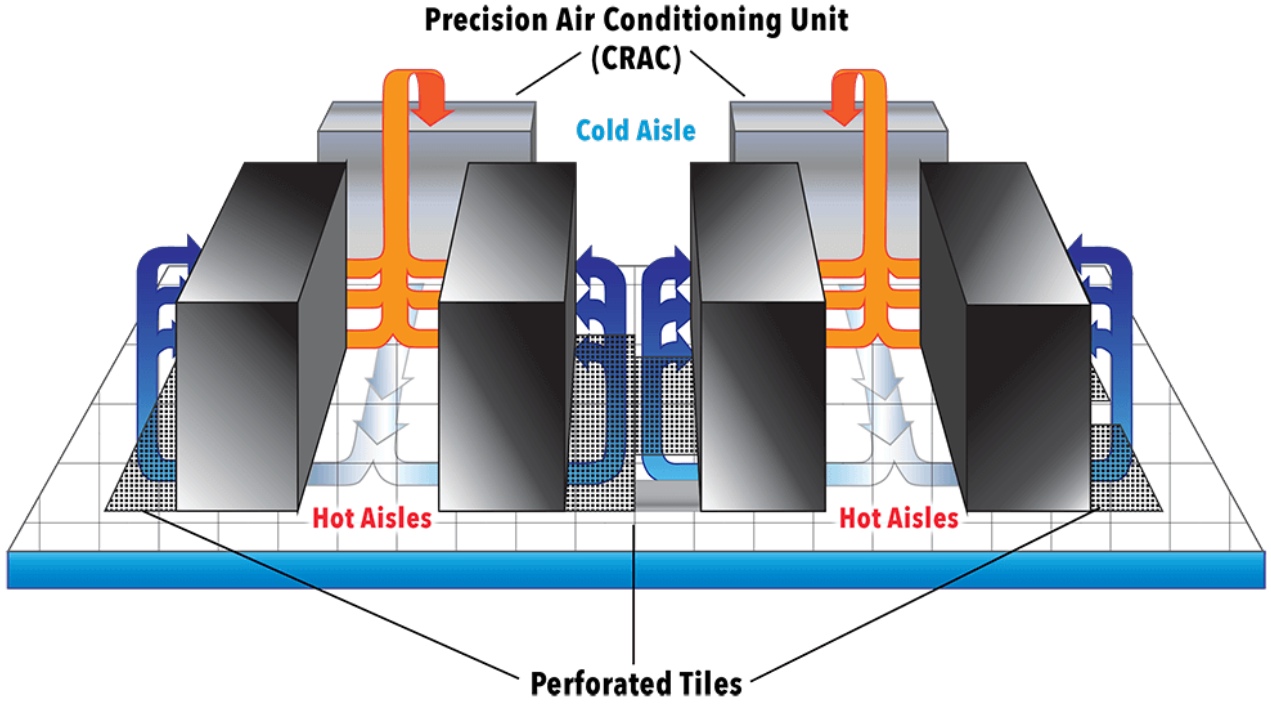
Table 4.2: Increase in temperature set point

|  | **Control** | **+ 1 °C** | **+ 2 °C** |
|---|---|---|---|
| Total energy need for cooling (kWh) | 93.54 | 89.28 | 84.12 |
| Relative difference to control execution | - | - 4.55 % | - 10.06 % |

## 4.3   Limitation of our aproach

**Thermal simulation**   Pre-cooling mechanism is not appropriate for all data centers. We remind that the thermal model and numerical parameters used in the simulation are adopted from a previous work, where it had been tested both experimentally and with TRNSYS software [26]. One of the assumption is strong: we consider the whole building with its walls, air mass and servers as a homogeneous thermal mass. According to the validation, the model is in line with experimental measurements, but it does not allow us to study in detail the temperature gradient inside the room. In particular, our cooling mechanism may cause hot spots to appear close to the servers. A Computational Fluid Dynamic study might help to figure it out. Also, the assumption we made is not suitable for data centers using the containment technique. This technique is widely used today for improving cooling efficiency, and consists in separating physically cold supply from hot air recuperation in different corridors (see Figure 4.2). With this architecture, we can no longer afford to model the data center as a homogeneous thermal mass. The cold air being blown directly on what we want to cool, there is no more thermal mass we can pre-cool.

**IT density**   As already mentioned a few times in this report, pre-cooling requires the building to have a sufficient thermal inertia and a low IT power density by square meter in order to have

Figure 4.2: Hot aisle / cold aisle layout design



an interest. The thermal inertia in our model is the parameter $a = exp\left(-\frac{h}{C.R}\right)$. The thermal capacitance $C$ measures the energy required by the thermal mass to increase its temperature by one degree. The thermal resistance $R$ measures how the walls resist to heat flow. According to [3], $C$ ranges from 0.015 to 0.065kWh/°C/m² while $1/R$ ranges from 0.001 (for a very efficient build- ing) to 0.003 kW/°C/m². In the same unities, the experimentally found values we used are: $C = 0.23$ and $1/R = 0.0019$. We will study the influence of these parameters in future work.

It can be noticed that the equation we use to describe the evolution of temperature does not depends on the surface:

$$T(t + h) = a.T(t) + (1 - a)[T_{ext}(t) + R(Q_{IT}(t) - Q_{cool}(t)]$$

$R$ and $C$ can thus be considered per unit area. The value of $a$ won't change because they appear in it as a product. The only real exogeneous parameter that will then influence the cooling strategy will be the density of heat generated by IT per square meter. This influence will also be studied in more details in future work.

# Part 5

# Future work

## 5.1   Research article

My internship not being finished yet, we still have a few things to work on until the last day. In particular, this study will be the subject of a research article that we hope to submit in September. We want to make at least the experiments described in this section.

**Results for all the year**   We plan to run our simulation on different days, representative of the different seasons and weather conditions. Figure 5.1 plots a view of the yearly data we have to our disposal.

**Quantifying the limits**   We have mentioned part 4.3 the limits of pre-cooling as a load shift mechanism. We want to be more specific and study the impact of thermal $R$ and $C$ parameters as well as the impact of IT density on the results. We want to be able to answer this question: from how many watts per square meter does pre-cooling lose its interest?

**IT workload**   The figures shown in this report are obtained with simulated Google IT traces. We can notice that they are relatively stable and do not take into account day/night variations. We want to complete this study by running in parallel other simulations with real traces. Alibaba published last year an 8-day record from 4000 of its servers (available on `github.com/alibaba/clusterdata`). We have plotted on Figure 5.2 the shape that IT power consumption would have over 48 hours. These very different data may give us other results and especially on synergies between pre-cooling and IT load variations.

**Pricing**   An influence that is not studied here is the influence of pricing over our results. We took $p_{sell} = $ €0.6 and $p_{buy} = $ €0.15 as it was in adequacy with French power market. We would like to try other value pairs representing different electricity contracts if prices are used to represent market prices or different energy mixes if they are used to represent greenhouse gas emission intensity.

Figure 5.1: By month, minimal, median and maximal renewable production curves for Belfort in 2006
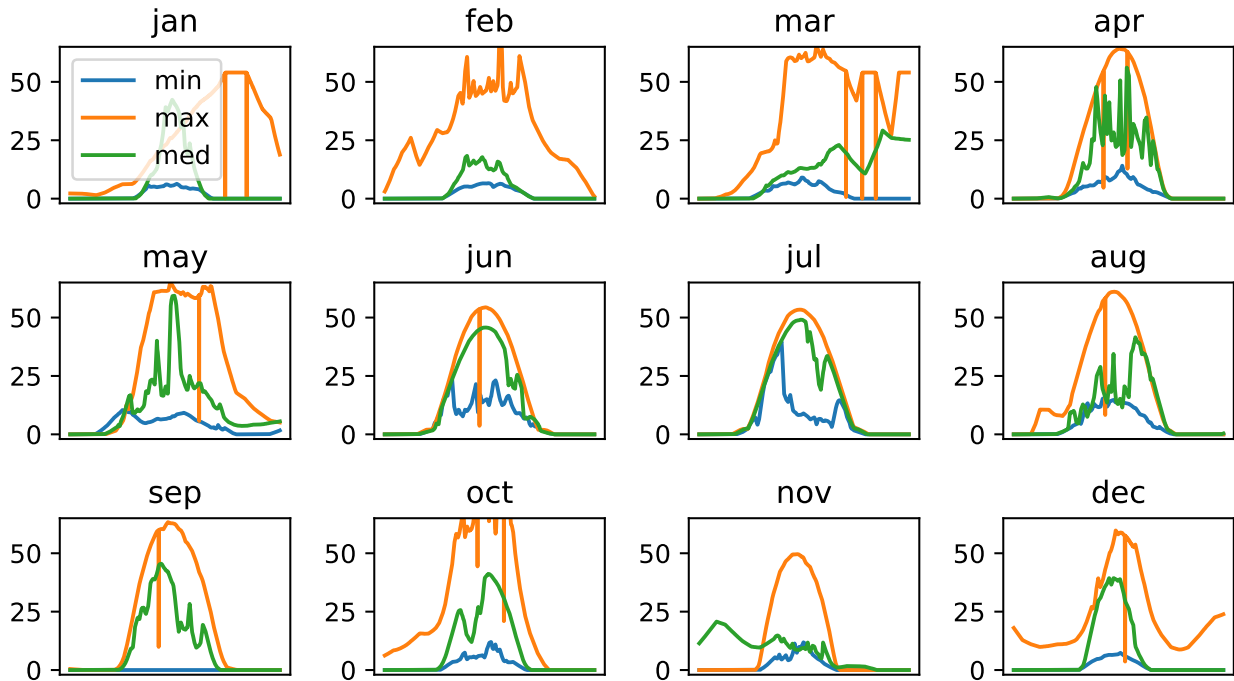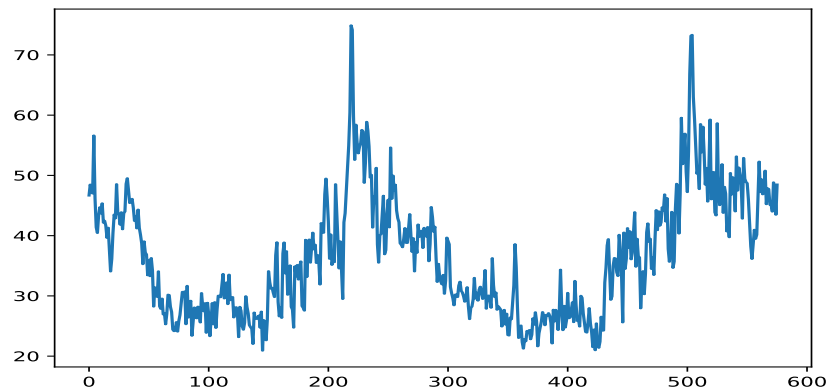


Figure 5.2: CPU usage (in percent of total capacity) for the first 48 hours of Alibaba 2018 cluster trace

## 5.2 Questions left open

There are also questions that have not been treated in this internship, but which deserve interest. Firstly, one may notice that we took into account in our model one cooling device only: a geothermal heat pump. Equations for the linear solver may be adapted to simulate other systems such as the often used Computer Room Air Conditioning system. Working on free cooling techniques as it is done in [17, 10] might be a good idea, as it has interesting synergies with renewable energy production. Secondly, as mentioned before, a stronger and more accurate thermal model like a Computational Fluid Dynamic study would help assessing the potentials and risks of lowering temperature and switching off the cooling device. ASHRAE thermal guidelines for data centers may be used more widely because we have only considered in our work the recommended temperature for all classes of equipment. Lastly, integrating this work into the more general framework of ANR Datazero project might have an important added value. Indeed, for the project as it is now, cooling is taken as a constant load. This could be refined by including decisions for cooling either on the power side or during the negotiation process. while optimizing source commitment in order to suggest to the IT side a reachable power profile,

All the suggested complementary work mentioned above could be the subject of an other master thesis or even a PhD.

# Conclusion

The present report represents the work of a 4-month research internship. It addresses an issue that is not often addressed in the literature: the ability to pre-cool a data center to face variability of renewable energy production. I am quite satisfied of our findings, the pace at which we have progressed and all that it has taught me about myself and the research in general.

If the results need yet to be refined before the end of this internship, we demonstrated that there was a potential of synergy between cooling strategy and renewable energy production curves. We created a model capable of capturing the thermal and power behavior of a data center. We then developed a simulation based on linear resolution that allowed us to try different data and parameters. We plan to make a research article with the last results before the end of the internship.

# Bibliography

[1]     Anders Andrae and Tomas Edler. "On Global Electricity Usage of Communication Technology: Trends to 2030". en. In: *Challenges* 6.1 (Apr. 2015), pp. 117–157. ISSN: 2078-1547. DOI: 10.3390/challe6010117.

[2]     Lotfi Belkhir and Ahmed Elmeligi. "Assessing ICT Global Emissions Footprint: Trends to 2040 & Recommendations". en. In: *Journal of Cleaner Production* 177 (Mar. 2018), pp. 448–463. ISSN: 09596526. DOI: 10.1016/j.jclepro.2017.12.239.

[3]     Duncan S. Callaway. "Tapping the Energy Storage Potential in Electric Loads to Deliver Load Following and Regulation, with Application to Wind Energy". en. In: *Energy Conversion and Management* 50.5 (May 2009), pp. 1389–1400. ISSN: 01968904. DOI: 10.1016/j.enconman.2008.12.012.

[4]     Stephane Caux et al. "IT Optimization for Datacenters Under Renewable Power Constraint". en. In: *Euro-Par 2018: Parallel Processing*. Ed. by Marco Aldinucci, Luca Padovani, and Massimo Torquati. Vol. 11014. Cham: Springer International Publishing, 2018, pp. 339–351. ISBN: 978-3-319-96982-4 978-3-319-96983-1. DOI: 10.1007/978-3-319-96983-1_24.

[5]     Gary Cook. *CLICKING CLEAN: WHO IS WINNING THE RACE TO BUILD A GREEN INTERNET?* en. Tech. rep. Greenpeace, Jan. 2017, p. 102.

[6]     Georges Da Costa, Leo Grange, and Ines De Courchelle. "Modeling and Generating Large-Scale Google-like Workload". en. In: *2016 Seventh International Green and Sustainable Computing Conference (IGSC)*. Hangzhou, China: IEEE, 2016, pp. 1–7. ISBN: 978-1-5090-5117-5. DOI: 10.1109/IGCC.2016.7892623.

[7]     Hugues Ferreboeuf, Maxime Efoui-Hess, and Zeynep Kahraman. *Rapport Lean ICT : Pour une sobriété numérique*. fr. Tech. rep. The Shift Project, Oct. 2018.

[8]     Inigo Goiri et al. "Parasol and GreenSwitch: Managing Datacenters Powered by Renewable Energy". en. In: (Mar. 2013), p. 13.

[9]     Léo Grange, Georges Da Costa, and Patricia Stolf. "Green IT Scheduling for Data Center Powered with Renewable Energy". In: *Future Generation Computer Systems* 86 (Sept. 2018), pp. 99–120. ISSN: 0167-739X. DOI: 10.1016/j.future.2018.03.049.

[10]    Ali Habibi Khalaj, Khalid Abdulla, and Saman K. Halgamuge. "Towards the Stand-Alone Operation of Data Centers with Free Cooling and Optimally Sized Hybrid Renewable Power Generation and Energy Storage". In: *Renewable and Sustainable Energy Reviews* 93 (Oct. 2018), pp. 451–472. ISSN: 1364-0321. DOI: 10.1016/j.rser.2018.05.006.

[11]  Maroua Haddad, Marie-Cécile Pera, and Christophe Varnier. "Stand-Alone Renewable Power System Scheduling for a Green Data-Center Using Integer Linear Programming Version 1". en. In: *HAL archives ouvertes [Research Report]* (2019), p. 39. DOI: hal-02081951.

[12]  D. Hatzopoulos et al. "Dynamic Virtual Machine Allocation in Cloud Server Facility Systems with Renewable Energy Sources". In: *2013 IEEE International Conference on Communications (ICC)*. June 2013, pp. 4217–4221. DOI: 10.1109/ICC.2013.6655225.

[13]  Uptime Institute. *Annual Data Center Survey Results*. en. 2019.

[14]  Atefeh Khosravi, Adel Nadjaran Toosi, and Rajkumar Buyya. "Online Virtual Machine Migration for Renewable Energy Usage Maximization in Geographically Distributed Cloud Data Centers". In: *Concurrency and Computation: Practice and Experience* 29.18 (Sept. 2017), e4125. ISSN: 1532-0626. DOI: 10.1002/cpe.4125.

[15]  Demetrio Laganà et al. "Reducing the Operational Cost of Cloud Data Centers through Renewable Energy". en. In: *Algorithms* 11.10 (Oct. 2018), p. 145. DOI: 10.3390/a11100145.

[16]  Yuling Li et al. "Thermal-Aware Hybrid Workload Management in a Green Datacenter towards Renewable Energy Utilization". en. In: *Energies* 12.8 (Apr. 2019), p. 1494. ISSN: 1996-1073. DOI: 10.3390/en12081494.

[17]  Zhenhua Liu et al. "Renewable and Cooling Aware Workload Management for Sustainable Data Centers". en. In: (June 2012), p. 12.

[18]  Maciej Z Lukawski et al. "Demand Response for Reducing Coincident Peak Loads in Data Centers". en. In: (Jan. 2019), p. 11.

[19]  Jean-Marc Pierson et al. "MILP Formulations for Spatio-Temporal Thermal-Aware Scheduling in Cloud and HPC Datacenters". en. In: *Cluster Computing* (Apr. 2019). ISSN: 1573-7543. DOI: 10.1007/s10586-019-02931-3.

[20]  Neil Rasmussen. "Calculating Total Cooling Requirements for Data Centers". en. In: *White Paper #25* American Power Conversion (APC) (2007), p. 8.

[21]  *Secure Microsoft Hosted Exchange Cloud Service*. en-US. https://ecocloud360.com/.

[22]  *Server Colocation | Green House Data*. https://www.greenhousedata.com/colocation.

[23]  *Solar Radiation DAta (SODA)*. http://www.soda-pro.com/.

[24]  ASHRAE TC9.9. *Data Center Power Equipment Thermal Guidelines and Best Practices*. en. 2016.

[25]  Yanwei Zhang, Yefu Wang, and Xiaorui Wang. "TEStore: Exploiting Thermal and Energy Storage to Cut the Electricity Bill for Datacenter Cooling". en. In: (2012), p. 9.

[26]  David P. Zurmuhl et al. "Hybrid Geothermal Heat Pumps for Cooling Telecommunications Data Centers". en. In: *Energy and Buildings* 188-189 (Apr. 2019), pp. 120–128. ISSN: 03787788. DOI: 10.1016/j.enbuild.2019.01.042.