



CARSHARING ANALYTICS LAB

Prof. Marco Mellia

AKA – a Big Data approach



SCENARIO

- Car sharing platforms are becoming popular
 - Only in Torino: Car2Go, Enjoy, ~~CarCityClub~~, BlueTorino
- They offer a “click-and-rent” paradigm
 - Using an APP, or the web, users can check which cars are available, book one, then start the rental
 - Platforms offer web API to interact with the system
 - Which cars are available, status of car, gas, etc.
 - Some API are public -- <https://github.com/car2go/openAPI>
- **Idea: use these platforms as source of information for mobility studies**
 - Lab focuses on this last aspect



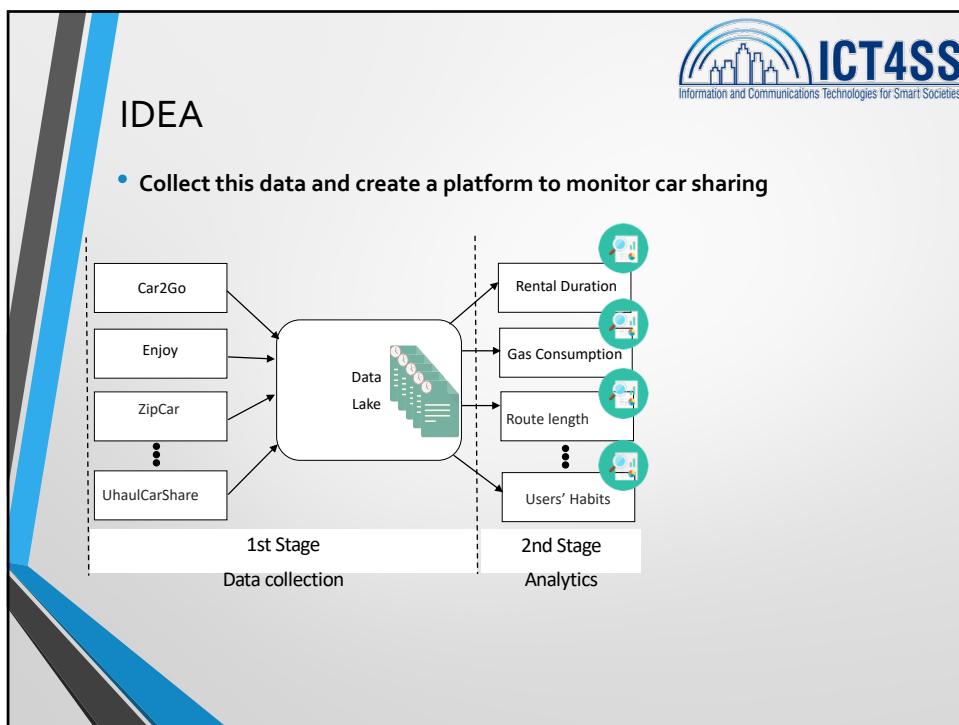
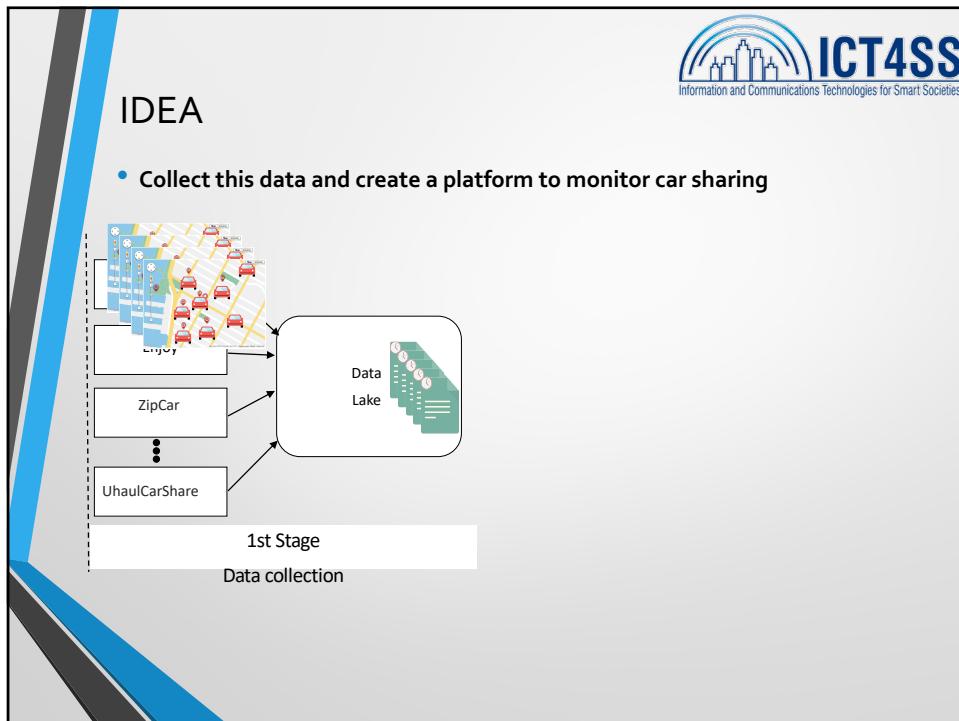
CARSHARING ANALYTICS LAB

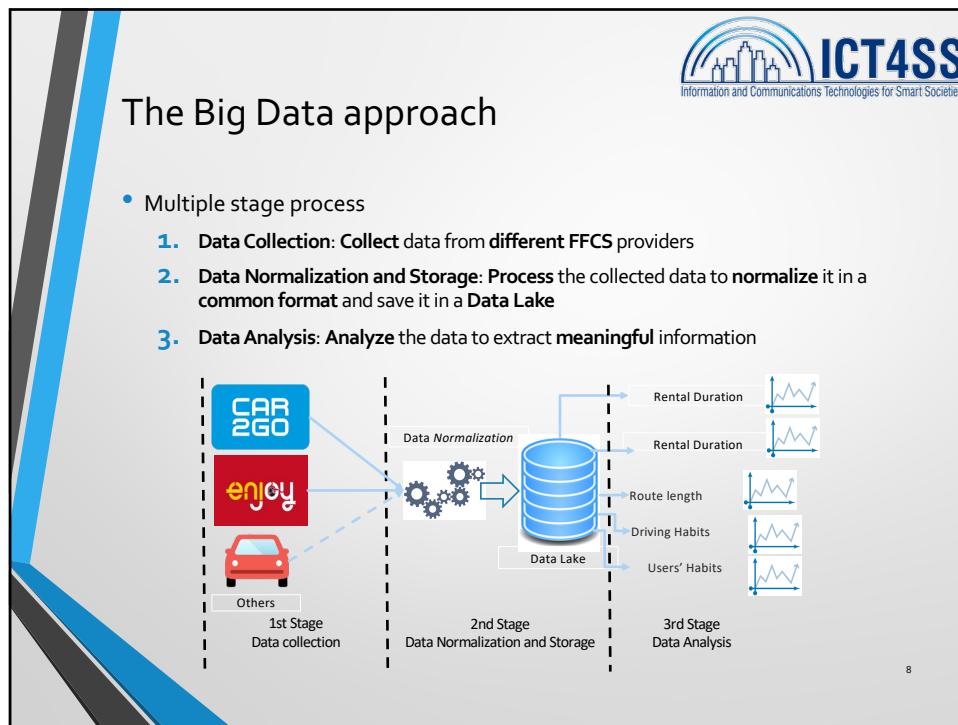
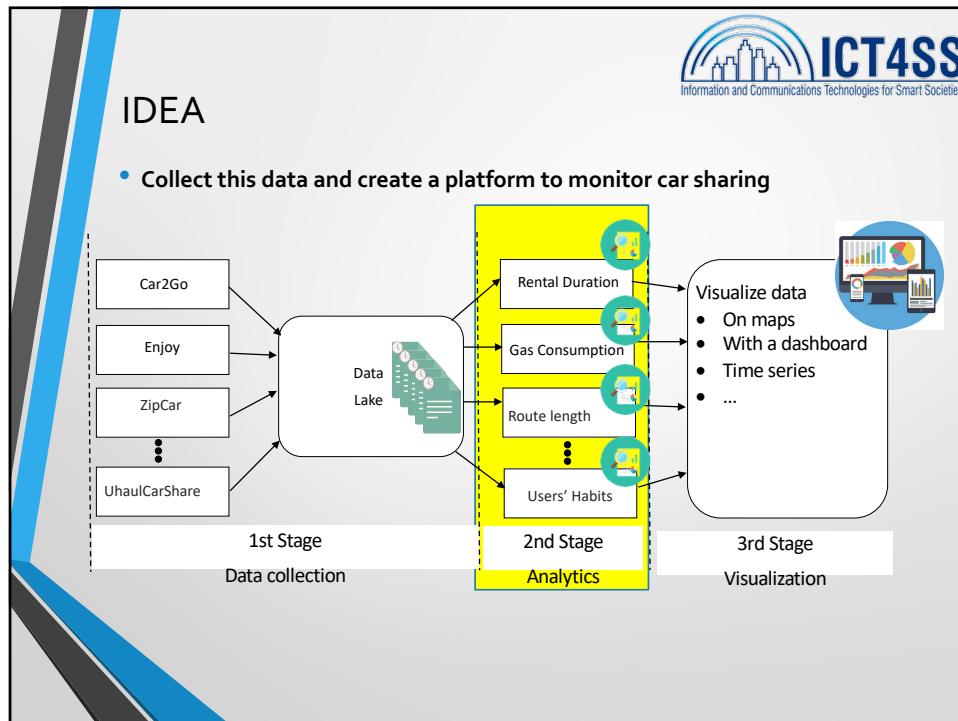


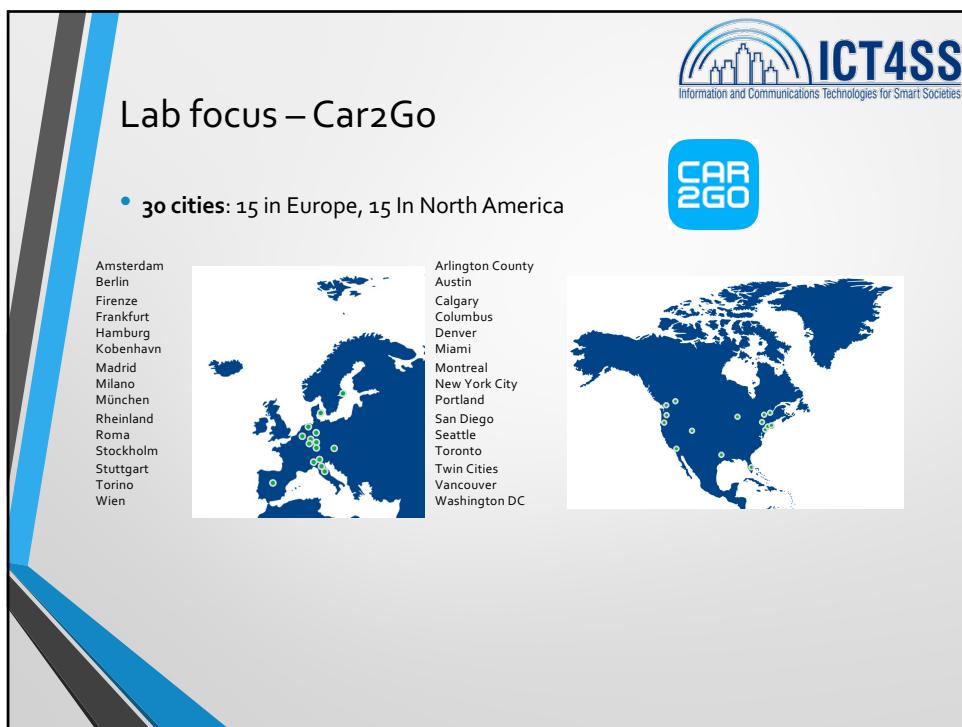
OVERALL PLATFORM

- Big Data approach
 - **Data Collection:** Collect data from Open Access Repositories
 - All car position in all Italian cities
 - Integration with maps information, public transportation, traffic status
 - **Data management and analysis:** Create a big data platform to harvest the raw data
 - Integrate data, and compute higher level metrics (car density, car flows, etc.)
 - **Design a graphical interface to visualize data**
 - Time series
 - Coloured maps
 - Dashboards







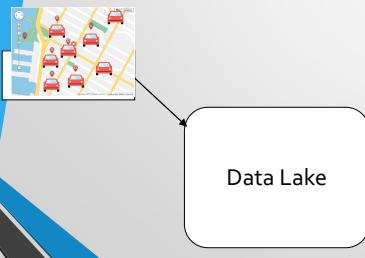


Data collection

- Harvest the data offered by car sharing platforms, and store them in a repository
- Consider the car2go platform
 - It offers public API to collect the data [need a key]

http://www.car2go.com/api/v2.1/vehicles?loc=turin&oauth_consumer_key=qetacar&format=json

• It returns a json file



```

{
  "placemarks": [
    {
      "address": "Via Argelati Filippo, 24, 20143 Milano",
      "coordinates": [
        9.1727,
        45.44989,
        0
      ],
      "engineType": "CIE",
      "exterior": "GOOD",
      "fuel": 96,
      "interior": "GOOD",
      "name": "581/FF119NT",
      "smartPhoneRequired": true,
      "vin": "WME4533421R145153"
    },
    {
      "address": "via Santander, 9, 20143 Milano",
      "coordinates": [
        9.16443,
        45.43981,
        0
      ],
      "engineType": "CIE",
      "exterior": "GOOD",
      "fuel": 96,
      "interior": "GOOD",
      "name": "285/FF457BT"
    }
  ]
}
  
```

Data processing– high level intuition

- Every **minute**, for every city, we have information about parked (available) cars.



- Every minute, we get a snapshot of available cars



Data processing– high level intuition



- When a car '**disappears**' at time T_1
 - It has been booked [!=rented]
 - We save the state "car booked at time T_1 "
 - We have temporary collection for activeBookings
 - [And we compute "**parking period = T_1 -To**"]



Data processing– high level intuition



- When the car '**reappears**' at $T=T_2$
 - The booking has ended
 - We save the state "car available at time T_2 "
 - We have temporary collection for activeParkings
 - [And we compute "**rental period= T_2 - T_1** "]



Data processing– high level intuition

- For every city, we have information about parked (available) cars



- Every minute, we get a snapshot of available cars
- We compare this with the latest status for each car
 - We have temporary collections
 - For **activeBookings**
 - For **activeParkings**

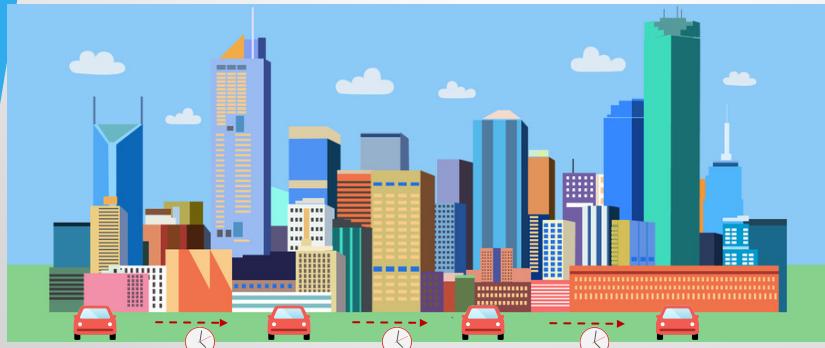
Dataset creation

- Create a dataset



- By repeating this forever, we obtain a longitudinal dataset

DATASET



The slide features a large, colorful illustration of a city skyline in the background. In the foreground, there is a row of red cars parked on a street. Each car is marked with a small circular icon containing a clock symbol and a red arrow pointing to the right, suggesting a temporal sequence or tracking of the vehicle's location.

Dataset creation

- Create a dataset



The diagram illustrates the process of dataset creation. It shows a sequence of six red cars parked along a road, each marked with a location pin. Above the cars, a circular icon with a clock symbol indicates a time interval. A red box highlights the third car with the text "The car is back". This visualizes how individual data points (car locations) are collected over time to form a longitudinal dataset.

- By repeating this forever, we obtain a longitudinal dataset
 - A collection for **PermanentBookings**
 - With durations of car booking periods
 - A collection for **PermenentParkings**
 - With duration of car parking periods

Data Lake

Dataset integration



- Integrate the data with other sources
 - Given a rental from $[x_0, y_0]$ to $[x_1, y_1]$
 - Which was the possible path the driver followed?
 - Which should be the minimum time to complete journey?
 - How long would it take to walk instead of driving?
 - How long would it take to use a bus/tram?
- Use google map to integrate the data
 - When a booking ends, query gmap for data
 - Limitations due to the number of (free) queries

Caveats



- In reality, things are more complicated
- A car may disappear, and reappear in the same place
 - GPS fix position error
 - Booking that has been cancelled
 - System issues?
- Car may disappear for (very) long time
 - Car brought to maintenance
- A lot of car may disappear at the same time
 - Maintenance?
 - System issues?
- Gmaps data may be missing
- ...
- Always double check what you get!



Repository for dataset



Data collection

- Which technology for our project?
- We have
 - json files
 - Heterogeneous platforms (car2go, enjoy, bluetorino, ...)
- NoSQL technology seems more adequate
 - Data can be not structured
 - Data may be missing
 - They are simpler to be used
 - They scale horizontally
- MongoDB is a free and open source document oriented DB
 - Well suited for our goals

Data Lake



MongoDB



ICT4SS
Information and Communications Technologies for Smart Societies

- It is a scalable NoSQL DB
- Uses JSON-like documents with schemas
- It has drivers for a variety of popular programming languages
 - Python, C, C++, Java, ...
- Simple query mechanisms


```
db.getCollection('ActiveBookings').find({})
db.getCollection('ActiveBookings').find({city: "torino"})
db.getCollection('ActiveBookings').find({city: "torino"}).count()
db.getCollection('ActiveBookings').find({city:
"torino"}).sort({plate: 1})
db.getCollection('ActiveBookings').distinct("city")
```
- We will learn how to use mongoDB for queries
 - The goal is not to become MongoDB masters, but rather to be practitioners

MongoDB - concepts



ICT4SS
Information and Communications Technologies for Smart Societies

- Think of **documents** as database **records**
 - Documents are just JSON objects that MongoDB stores in binary
- Think of **collections** as database tables

RDBMS (mysql, postgres)	MongoDB
Tables	Collections
Records/rows	Documents/objects
Queries return a record	Queries return a cursor



MongoDB - concepts

- Queries return "cursors" instead of a collections
 - A cursor allows you to iterate through the result set
 - A big reason for this is performance
 - Much more efficient to load results into memory
 - Especially if results are big as in big data
- The find() function returns a cursor object

```
var c = db.ActiveBookings.find( {city: "Torino"} )  
var i = 0  
while (c.hasNext() && i<10)  
{  
    var o = c.next() // this is the object  
    print(o.init_time + " " + o.city)  
    i++  
}
```

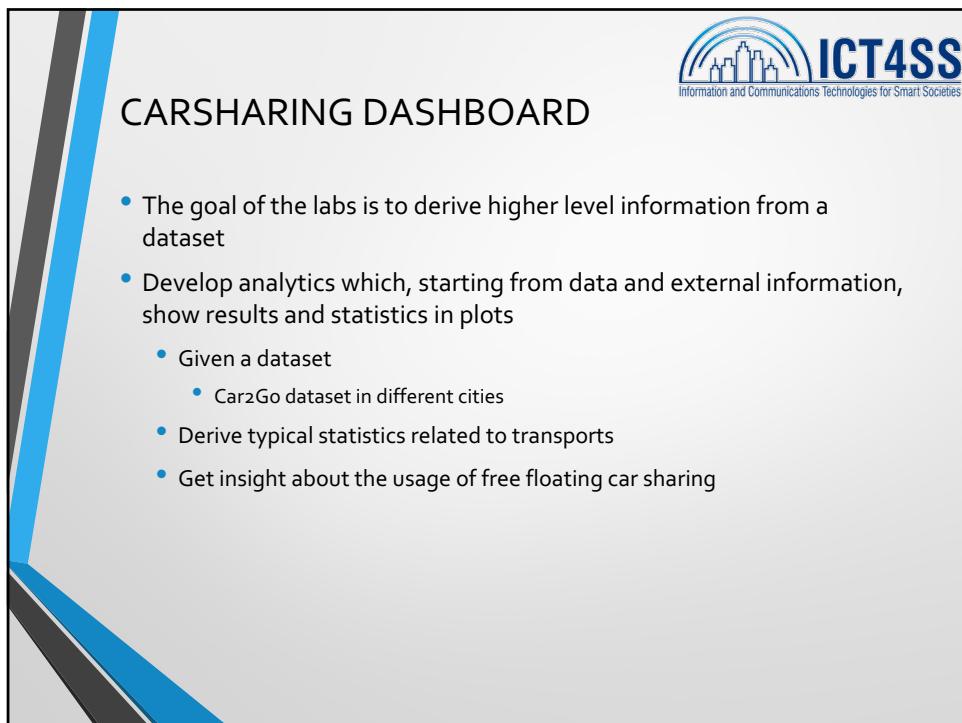


DEMO

ROBOMONGO



Analytics



- The goal of the labs is to derive higher level information from a dataset
- Develop analytics which, starting from data and external information, show results and statistics in plots
 - Given a dataset
 - Car2Go dataset in different cities
 - Derive typical statistics related to transports
 - Get insight about the usage of free floating car sharing

CARSHARING DASHBOARD

- Collect and visualize information about the car sharing service.

CarSharing Statistics

Number of booked cars

Booking duration → (Guessed) time of booking
Parking Duration → (Guessed) time of driving

Locations

CARSHARING DASHBOARD

- Collect and visualize information about the car sharing service.

Number of booked cars

Line Graph: N° noleggi (Number of rentals) vs Ora (Hour). Legend includes: lunedì 18 gennaio 2016, martedì 22 dicembre 2015, mercoledì 13 gennaio 2016, giovedì 7 gennaio 2016, venerdì 15 gennaio 2016, sabato 9 gennaio 2016, domenica 17 gennaio 2016, and Media feriale (Festive average).

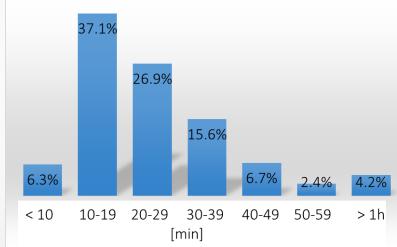
Bar Chart: % Verco (Percentage of Verco) vs Data (Date). The chart shows daily fluctuations in the percentage of Verco usage.

CARSHARING DASHBOARD

• Collect and visualize information about the car sharing service.

CarSharing Statistics





Booking duration ^(*)	Fraction [%]
< 10	6.3%
10-19	37.1%
20-29	26.9%
30-39	15.6%
40-49	6.7%
50-59	2.4%
> 1h	4.2%

Rental Type^()**

Rental Type ^(**)	Total	Fraction [%]
Chain	3297	8.8%
Single	31940	84.9%
Single*	1517	4.0%
tour	849	2.3%

^(*) for a proper definition of duration...

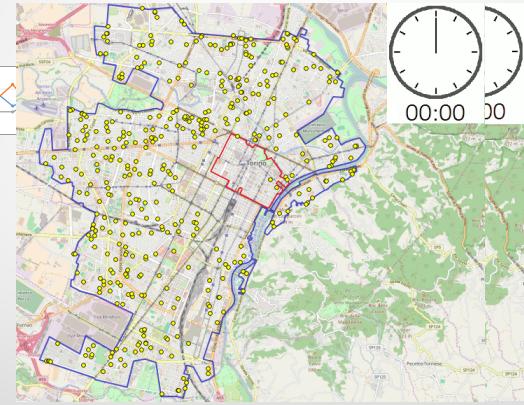
^(**) for a proper definition of type

CARSHARING DASHBOARD

• Collect and visualize information about the car sharing service.

Locations







CARSHARING DASHBOARD

- Collect and visualize information about the car sharing service.



Google APIs
[Google Maps and Google Direction API]

- Develop a system which, starting from data, shows a dashboard offering statistics
 - Correlate rentals
 - with Public Transport availability
 - Walking distance
 - Time of day
 - ...



Lab organization



Lab organization

- Students work in groups of 3 people
 - 2 is too small, 4 is too big
- Better having a mixed background
 - As when you get your job
- We give you access to a MongoDB database
- You extract the data, process it, obtain statistics, and plot them
- Useful tools
 - Data extraction: A bit of MongoDB
 - Data postprocessing: A bit of python | Java | awk | bash | yourchoice
 - Data plotting: Matlab | Matplotlib | Gnuplot | excel | yourchoice
 - Report editing: Word | Openoffice | LaTeX | yourchoice



Lab Evaluation

- Each group writes a 5pp report + appendix
 - Describing what you have done
 - The results you got
 - Source code goes in appendix
- The group report must be uploaded on the didattica website
 - Filename: **groupXX.pdf**, one file per group, latest copy is evaluated
 - Deadline: 23:59 of the same day of the exam reservation deadline



Lab Evaluation

- The report is evaluated
 - On a range between [ins:30]
 - It's about 25% of final vote
 - All students get the same vote for the report
 - No report => no oral exam
 - Insufficient report => no oral exam
 - Copied report => no oral exam + disciplinary committee
- During the oral exam
 - One question on your report to discuss what you did
 - And check that you did it
 - Copied report => no oral exam + disciplinary committee
 - No answer => no pass



Lab Organization

- Formation of lab groups
 - Go to the gDrive document and enter your group
<https://qoo.qi/SqdTDY>
 - No preference => I choose
 - One or more PC per group



QUESTIONS

- Ask whether there are doubts or something not really clear.

