

Institut  
québécois  
d'intelligence  
artificielle

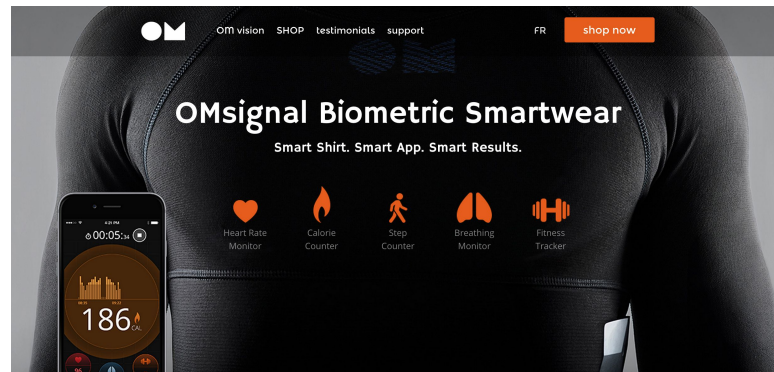


Mila

# OMsignal Project Block 1

Arsene Fansi-Tchango, PhD  
Simon Blackburn

# Company



**1,8 billion+**  
heartbeats recorded



**500 million+**  
breaths recorded



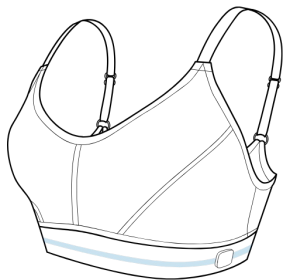
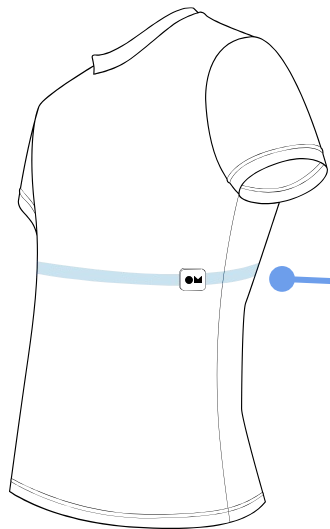
**330 thousand+**  
hours of data collected

Make personal health and wellness central to our daily lives, through the world's most advanced biosensing apparel platform.

# Technology

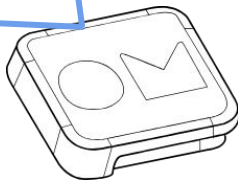
## Garments

The apparel picks up the body's signals using strategically placed ECG, Respiration, and Physical Activity sensors.



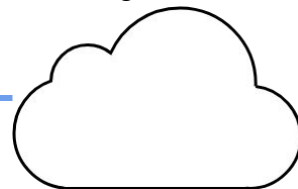
## Recording Module

The recording module transmits the signal to your smartphone, live.



## Cloud + AI

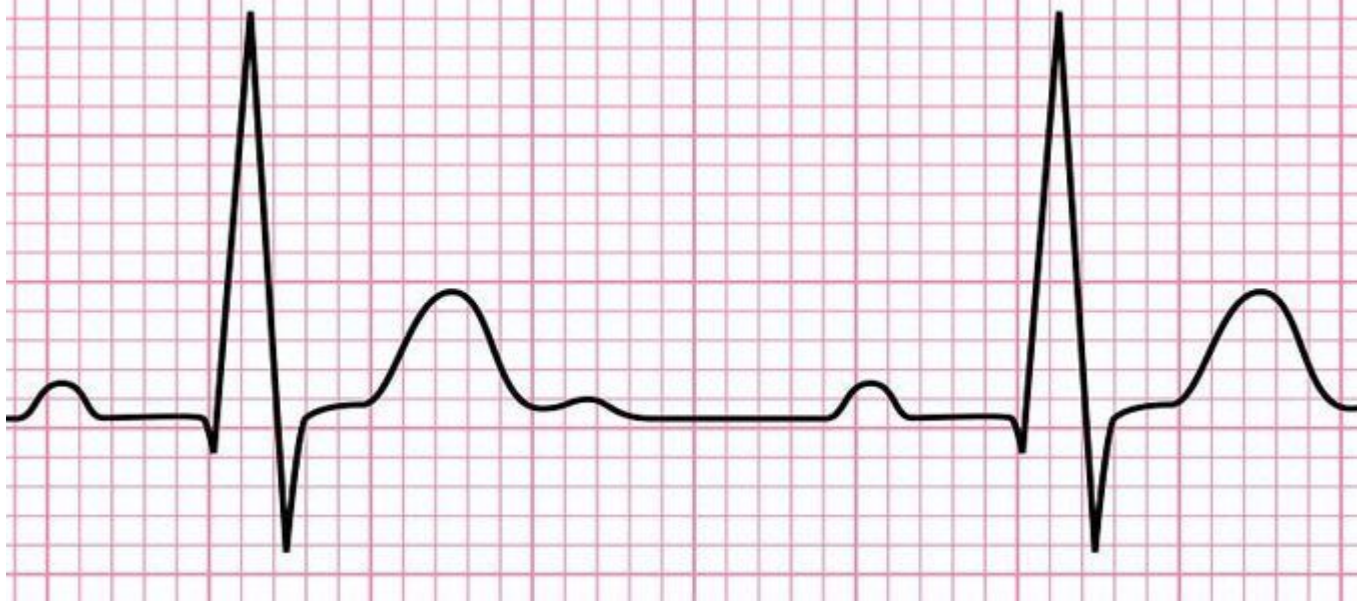
The data is sent to the Cloud to be further analysed using advanced algorithms and AI.



# Operational Challenges

- **Easy** to collect unlabeled data
  - Huge amount of data captured under different conditions
    - running / walking / sitting / sleeping, etc...
    - different levels of signal to noise ratio
- **Hard** to label this data for supervised learning
  - Experts (e.g., medical doctors) are expensive
  - Time demanding
    - E.g., walk through all the samples of a signal

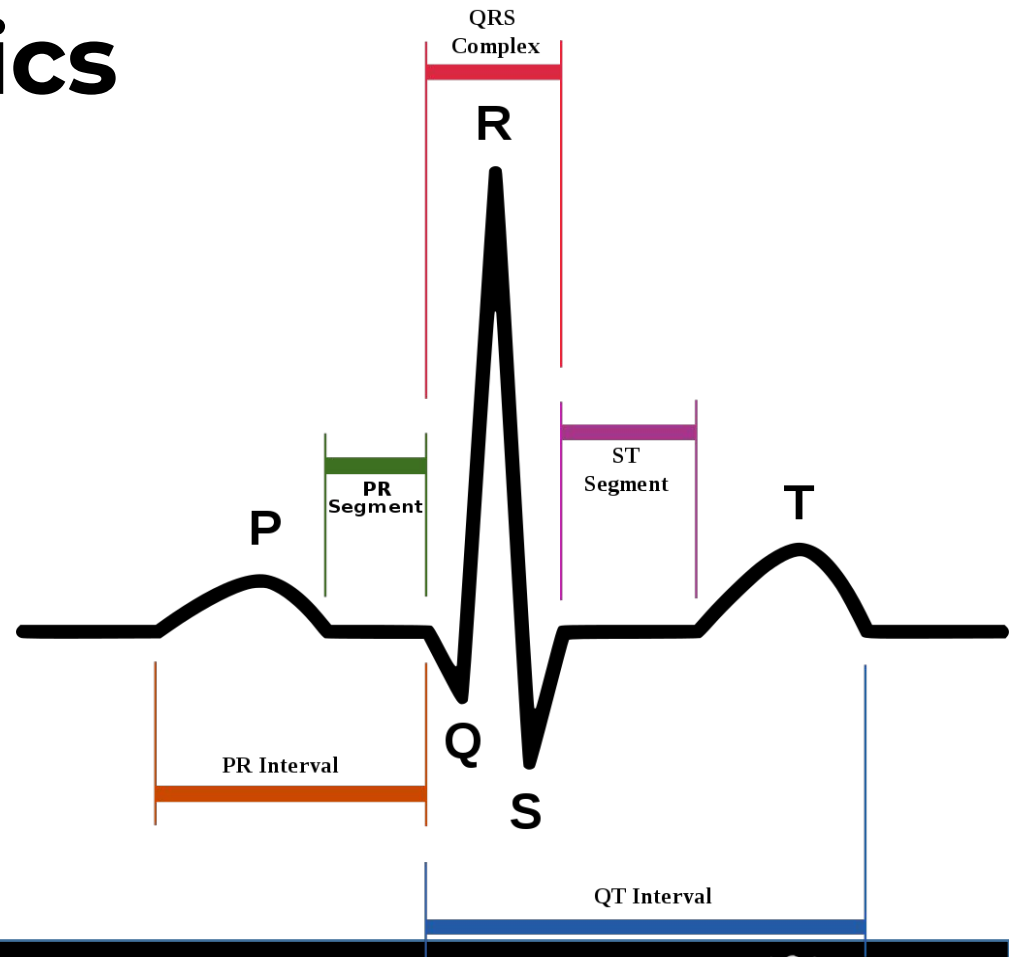
# ECG Example (1 lead)



From <http://www.onlinebiologynotes.com/electrocardiogram-ecg-working-principle-normal-ecg-wave-application-of-ecg/>

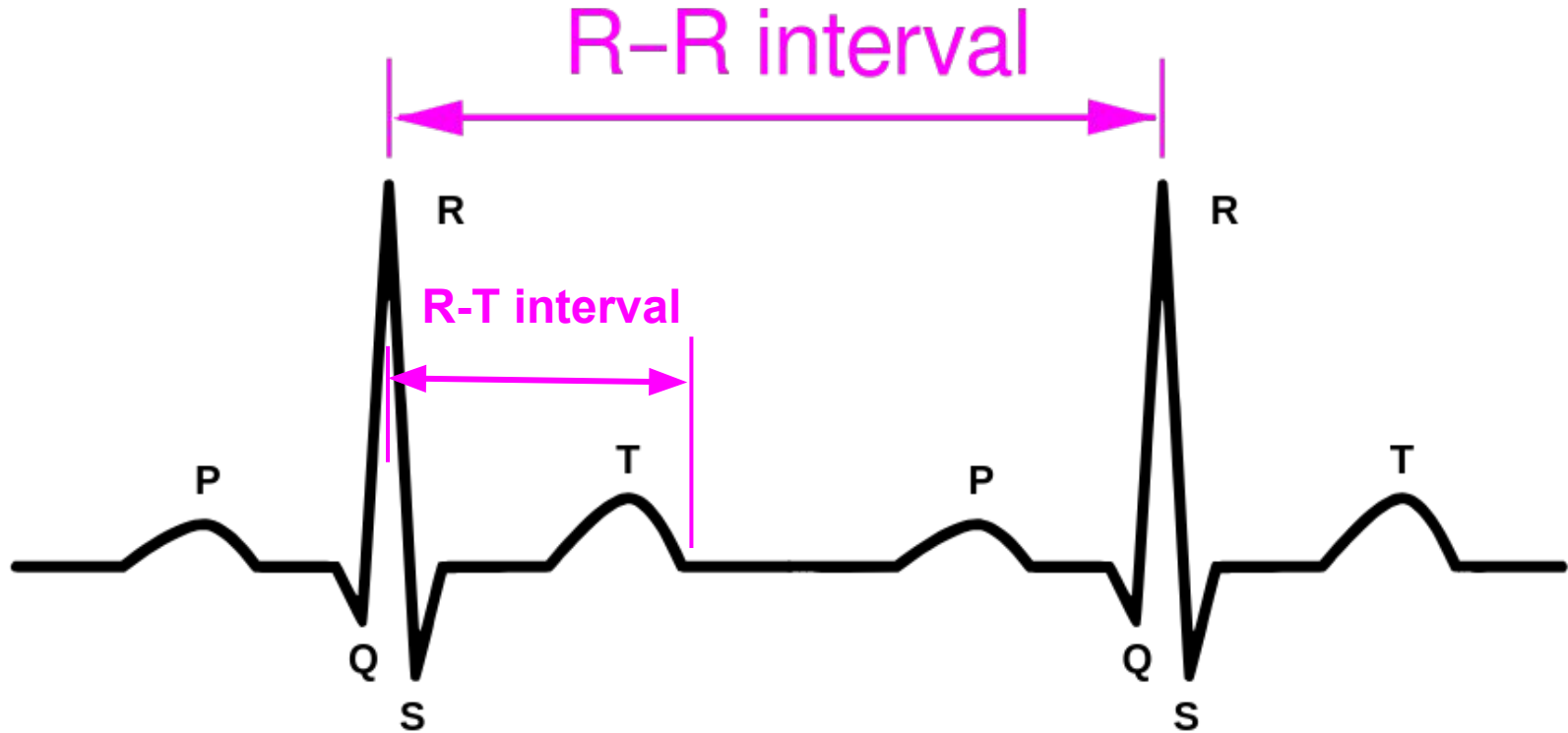
# ECG Characteristics

- **Fiducial points:** P, Q, R, S, T
- **P-Wave:**
  - Indicates atrial depolarization (systole)
- **QRS wave:**
  - Represents the ventricular depolarization (systole)
- **T- wave:**
  - Indicates ventricular repolarization (diastole)
- **P-R interval:**
  - Represents the time required for an impulse to travel through the atria
- **S-T segment:**
  - Represents the time when ventricular fibres are fully depolarized



From <https://en.wikipedia.org/wiki/Electrocardiography>

# ECG Characteristics



# OMsignal Project

- **Overall goal:** develop an **unsupervised/semi-supervised** representation learning approach that produces representations useful for tasks that have little labeled data:
  - Identification of the user
  - Fiducial point distributional information
    - Mean of the PR-Interval (real value)
    - Mean of the RT-Interval (real value)
    - Standard deviation of the RR-Interval (real value)



# Data

## OMsignal MyHeart project:

- **Private** data
- **32** Participants
- ECG signals are divided into windows of 30 seconds each at 125 Hz (**3750** samples per window)
- Labeled data:
  - **15** windows for each participant are labeled
  - Among them, **5** windows are used as test data
  - The remaining **10** are provided as train/validation data
- Unlabeled data:
  - **657233** windows

# Data Formats

- The data is provided as **numpy.memmap** files:
  - MILA\_TrainLabeledData.dat: **160 x 3754**
  - MILA\_ValidationLabeledData.dat: **160 x 3754**
  - MILA\_UnlabeledData.dat: **657233 x 3750**
  - MILA\_TestLabeledData.dat: **160 x 3754 (Blind)**
- The samples in each of train, validation and test originate from 3 different days.

# Data Formats

- **Labels:** 4 last columns of labeled datasets
  - **PR\_Mean:** 3751th column
  - **RT\_Mean:** 3752th column
  - **RR\_StdDev:** 3753th column
  - **UserID:** 3754th column
- Code to read/write numpy.memmap files into/from numpy.array is provided

# OMsignal - Block 1 goals

- Build a simple supervised baseline, using only the limited amount of labeled data

# Block 1 - Supervised learning

## Instructions / Expected timeline

	2019/01/14 week	2019/01/21 week	2019/01/28 week	2019/02/04 week
Tasks / Homework	<ul style="list-style-type: none"><li>• Data visualization</li><li>• Data augmentation</li></ul>	<ul style="list-style-type: none"><li>• Code the data loader for the provided dataset</li><li>• (optional) Implement a supervised single-task model for the identification task</li></ul>	<ul style="list-style-type: none"><li>• Implement a supervised multi-task model</li></ul>	<ul style="list-style-type: none"><li>• Write a short report summarizing the work, and results</li><li>• (Peer-) Review of other teams' code</li></ul>
Objectives / Deliverables	<ul style="list-style-type: none"><li>• Have a clear understanding of the data</li></ul>	<ul style="list-style-type: none"><li>• Data loader</li><li>• (optional) Single task model</li></ul>	<ul style="list-style-type: none"><li>• Multi-task model</li></ul>	<ul style="list-style-type: none"><li>• Produce documented code and report summarizing the experimental work</li><li>• Provide model for blind test set evaluation</li><li>• Complete the peer code review</li></ul>

# Deadlines

- Each team needs to provide the deliverable (report + code + best model) corresponding to a block at the latest on Friday noon (12:00pm) of the last week of the block.
- Any block deliverable that is provided past Friday noon (12:00pm) of the last week of a block will automatically get 0% for the peer evaluation.

# Deadlines

- Any block deliverable that is provided past Tuesday 11:59pm following the last week of a block will automatically get 0% for the UdeM evaluation.
- Peer evaluation must be completed by Monday 11:59pm following the last week of a block.

# Evaluation for Block 1

## ***25% of the final score***

- 10% Code review [5% of averaged peer evaluation + 5% UdeM]
- 12% Report evaluation [UdeM]
- 3% Model performance evaluation on blind test set [UdeM]



# Code review - Peer evaluation

- **10% of the final score** [5% of average peer evaluation + 5% UdeM]
- Random assignation of code reviews
- The code provided by a team will be evaluated by at **least 2 other teams**

Code quality (peer evaluation + UdeM evaluation)	8
Coherent and modular code/file organization (e.g. data processing, model definition, model training, model inference are in different files/modules; no code duplication)	
Code respects the <a href="#">PEP8 standard</a>	
Comments are relevant (see <a href="#">article</a> )	
Proper management of input arguments in the training script (see argparse, python fire, configparser)	
Proper utilization of GitHub (e.g. branching, relevant commits and messages, usage of pull request)	
Meaningful variable and function names	
Executable scripts with a “main” function (see <a href="#">article</a> )	
Reproducible experiments (e.g. seed)	

Introduction	2
Introduction to the project	
Brief introduction to the methods that will be used in the report	
Methodology	6
Description of the algorithms and the experiments (including hyperparameter fine tuning (if appropriate), etc.)	
Data description and data selection (train/valid/test, number of samples, shape/structure of data points)	
Results and discussion	6
Presentation of results (tables, figures, etc.)	
Discussion of results	
Conclusion	2
Recommendation for next steps	
Summary of project state (what was done, what needs to be done)	
Quality of the report	2
Report format (title with team member names, clear sections, flow between sections, figures and tables titled, axes titled, etc.)	
Report is short and to the point (5-7 pages including references, font size 11)	

# Report Evaluation

- **12% of the final score**
- 5-7 pages
  - Including figures, tables, and references
- Single column
- Font size 11
- Use the [NeurIPS](#) LaTeX format

# Blind Test Set Evaluation

- **3% of the final score.**
- If the best model provided by a team crashes or provides results that are statistically worse than those of the baseline model provided by the TAs, the team gets 0%.
- Otherwise, if the best model provided by a team is statistically equivalent to the baseline model, the team gets 1%.
- Otherwise, if the best model provided by a team is statistically better than the baseline model:
  - The team gets 3% if the model is the best performing one or is statistically equivalent to the best performing model provided by another team.
  - Otherwise, the team gets 2%.

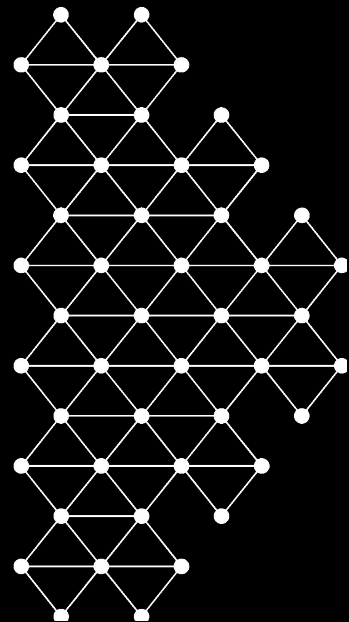
# Code Execution - Blind Test Set Evaluation

- The **test dataset** will be structured as follows:
  - A **numpy.memmap** file containing unlabeled samples of shape **160X3750**
- You will not have access to the **test set** and we will be executing your code on the **test set** ourselves.
- We will provide explicit instructions and examples for you to enhance an evaluation skeleton script that will be provided to you. You will need to complete this script. We reserve the right to give 0 if we cannot execute your code.
- Tips: **do not shuffle the test set.** Predictions should be made sequentially.

# Official Evaluation Metrics

- Classification task
  - Macro Average Recall Score (`sklearn.metrics.recall_score`)
- Regression tasks
  - Kendall Correlation Score for each task (`scipy.stats.kendalltau`)
- Overall Score:
  - All individual scores are clipped at zero
  - Geometric mean of the scores of the 4 tasks
- The code of the scoring function will be **provided**

# OMsignal Project Data Visualization



# Visualization in ML

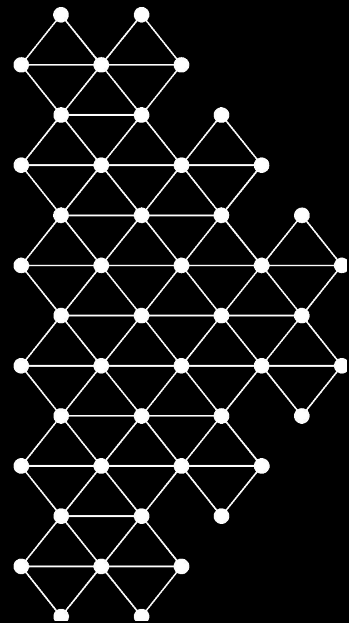
- It is always **good practice** to visualize data used in a ML (Machine Learning) project
- How to visualize high dimensional data (e.g., 3750) ?
  - Plotting all the samples as a line graph (**ok** in our case)
  - Techniques for dimensionality reduction
    - PCA,
    - t-SNE,
    - Auto-encoder,
    - ....

# Homework

- Choose randomly, from the **original Train dataset**, 1 window for **3** different users and plot each of them as a separate line graph figure.
- Use 2 dimensionality reduction techniques to project all the windows from the **original Train dataset** into a 2D space.
  - Make a scatter plot of the results. Use different colors per user.
  - Analyse the results.

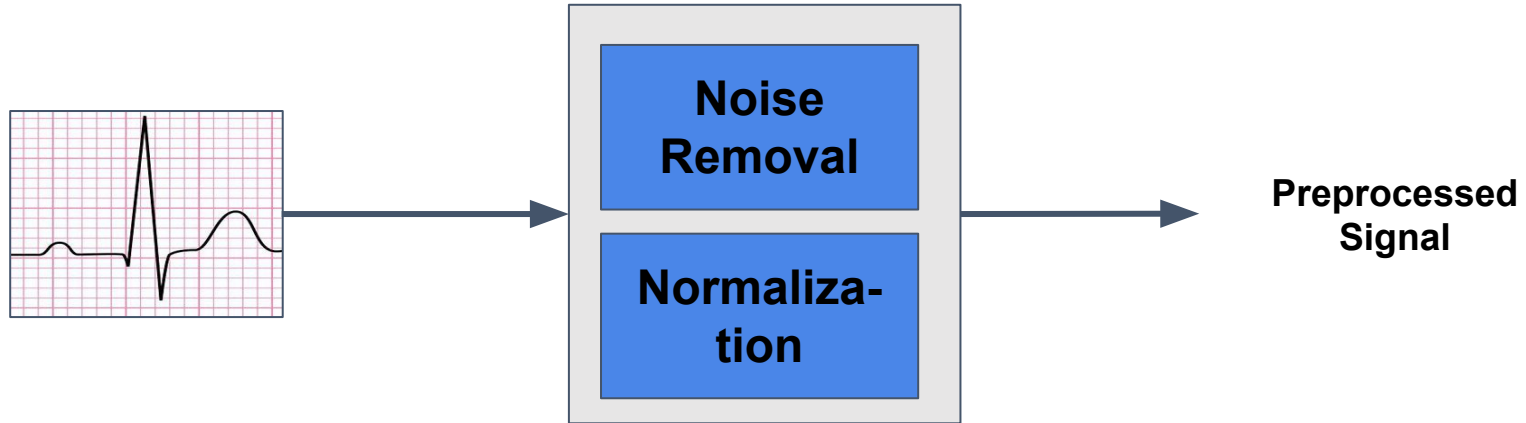


# OMsignal Project Data Preprocessing



# Data Preprocessing

- The collected ECG signals come with some noise and different levels of amplitude (running vs walking)
- **Needs preprocessing**



# Data Preprocessing

- **Noise Removal**

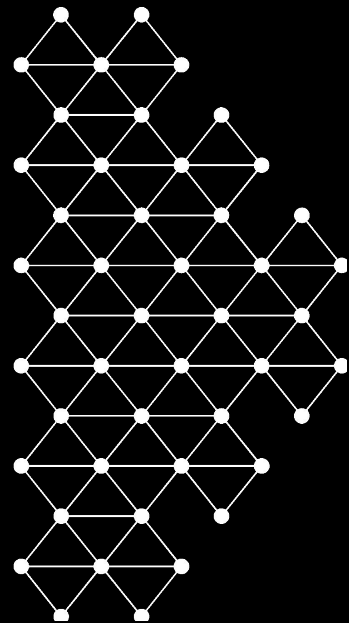
- Remove baseline wander (low frequency noise) with a moving average

- **Normalization**

- Normalize the amplitude of the signal within a moving window

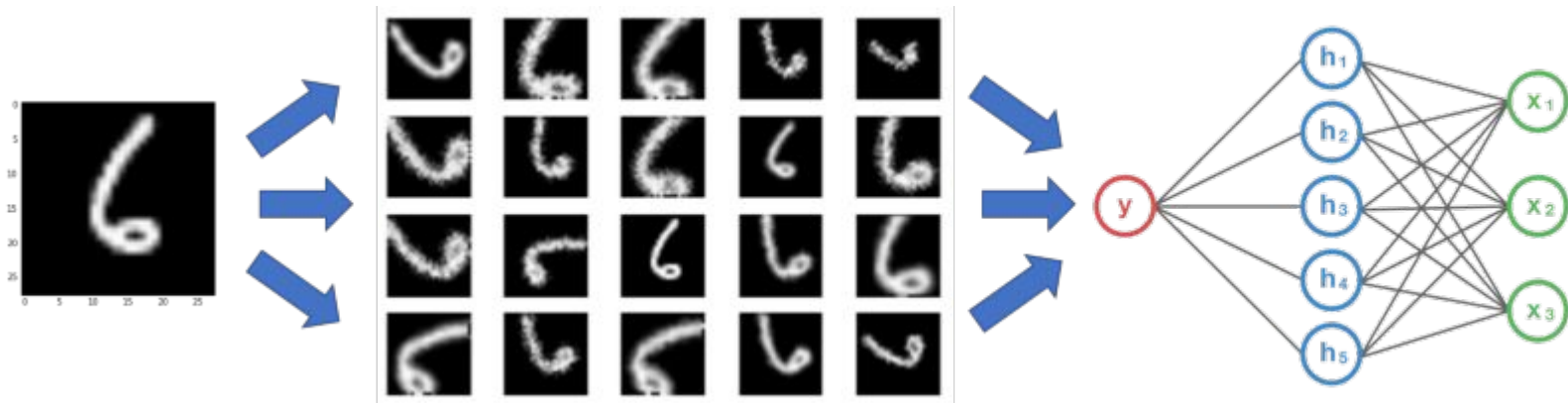
- The code for the preprocessing is **provided**.

# OMsignal Project Data Augmentation



# Data Augmentation

- DL (Deep Learning) in general requires large datasets for robustness
- How to get more data?



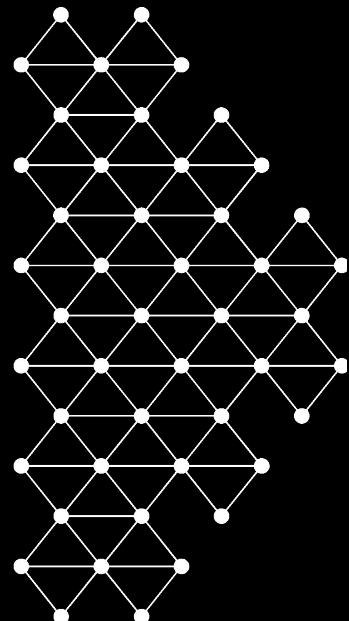
<https://medium.com/nanonets/how-to-use-deep-learning-when-you-have-limited-data-part-2-data-augmentation-c26971dc8ced>

- **offline** augmentation vs **online** augmentation

# (Possible) Data Augmentation for ECG

- **Signal shift**
  - Shift elements along the temporal axis (assuming periodic boundary conditions)
- **Partial noise addition**
  - Replace at most 2 second length sub-windows with controlled noise (matching the mean and the variance of the underlying samples)
- **Upside-Down inversion**
  - Negate the values of the signal
- Other ideas ...

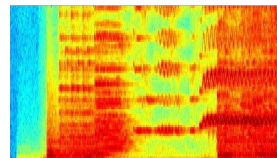
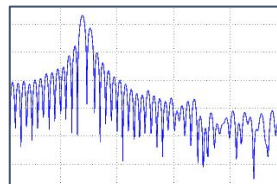
# OMsignal Project Data Transformation



# Data Transformations

It can be interesting to supply different representations of the data to a network to improve performance. For ECG signals, one can consider:

- **Fast Fourier Transform**
- **(log) Spectrogram**
- Other ideas ...
- Code for these transformations is **provided**.

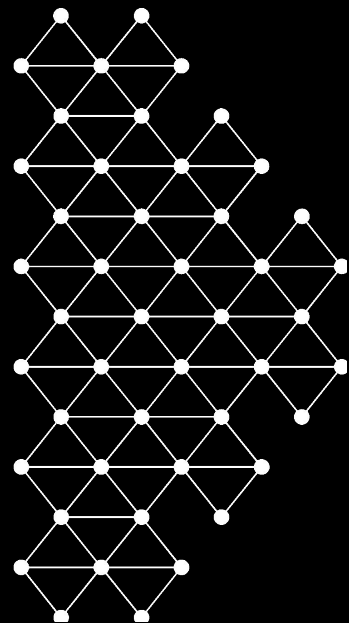




# Homework

- Implement your data loader for the labeled datasets
  - The code should provide the ability to perform **online** data augmentation if needed

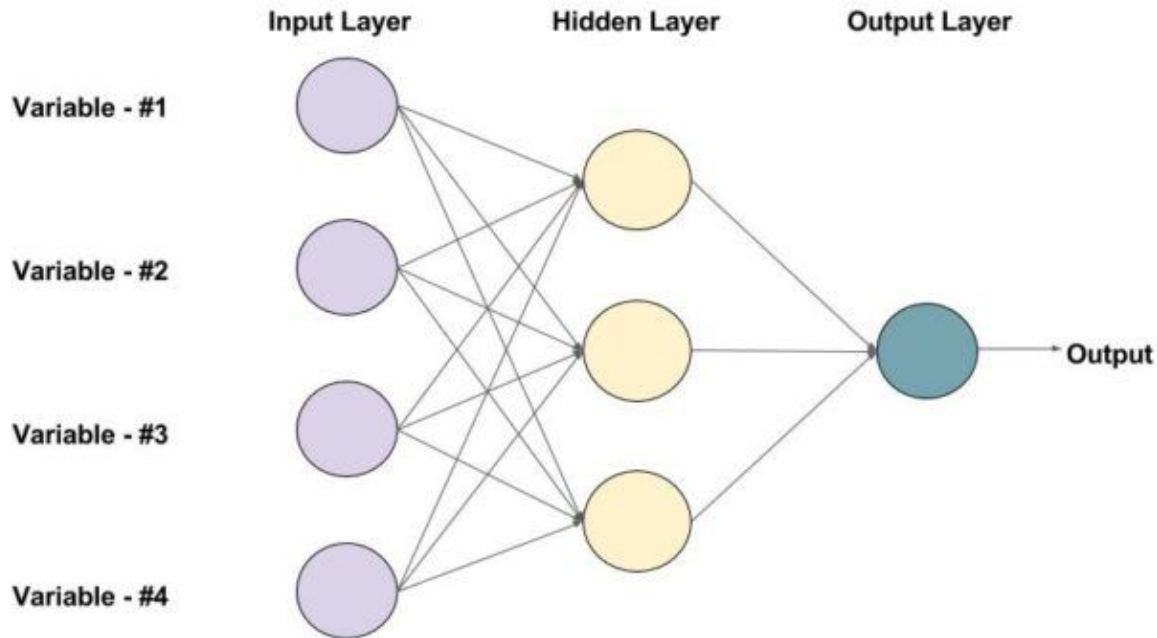
# OMsignal Project Models



# Basic Models for 1D Signals

- **MLP**

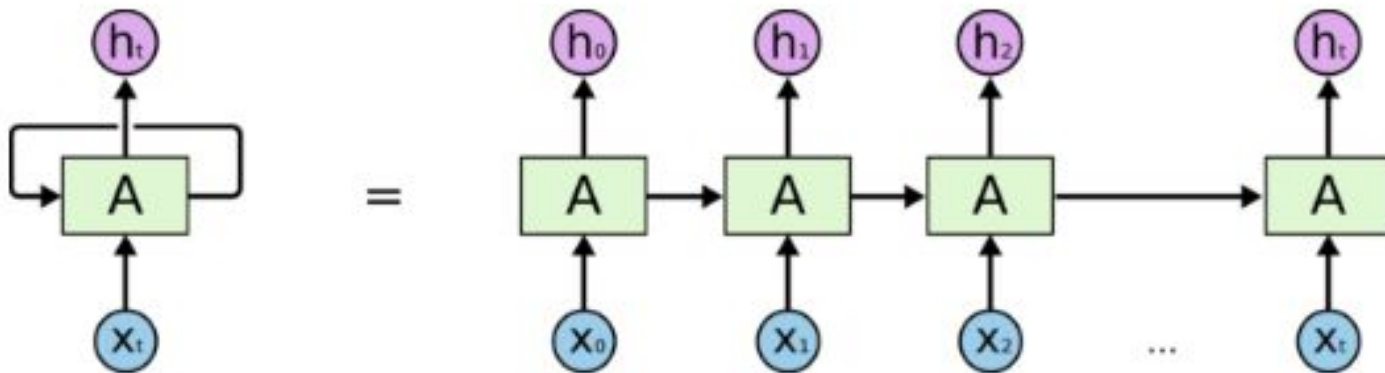
- The size of the input is fixed: 3750



An example of a Feed-forward Neural Network with one hidden layer ( with 3 neurons )

# Basic Models for 1D Signals

- **Recurrent Neural Networks**
  - Process the signal as a sequence

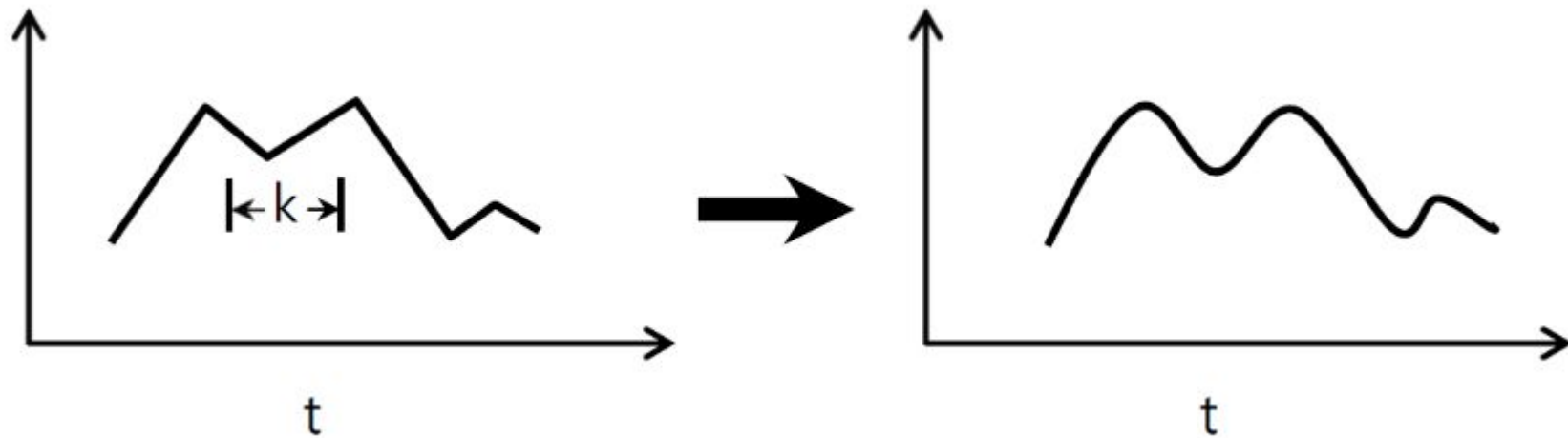


An unrolled recurrent neural network.

# Basic Models for 1D Signals

- **Convolution 1D**

- Process the signal on a windows basis. Apply the same filter on each window.

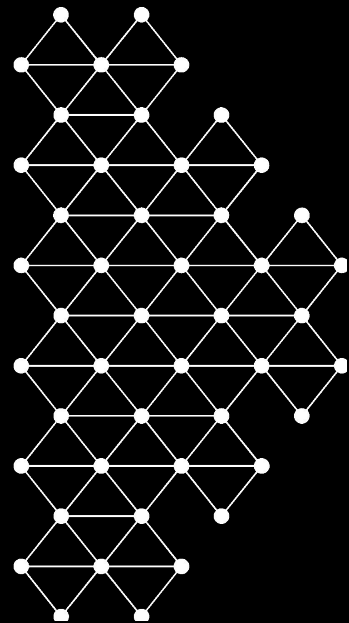


# Models

- **Other types (including combinations of the previously described models) can be considered**
- **Some transformed signals (e.g. spectrogram) might require other types of models**



# OMsignal Project Regularization

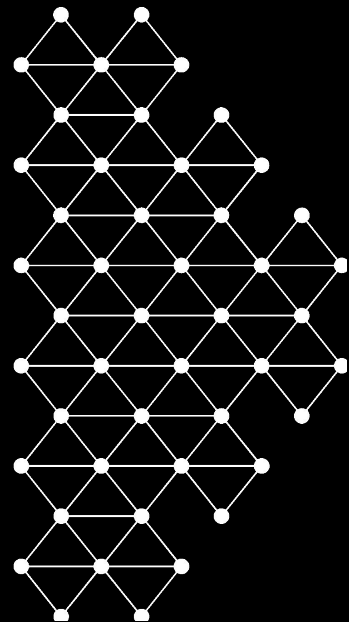


# Regularization

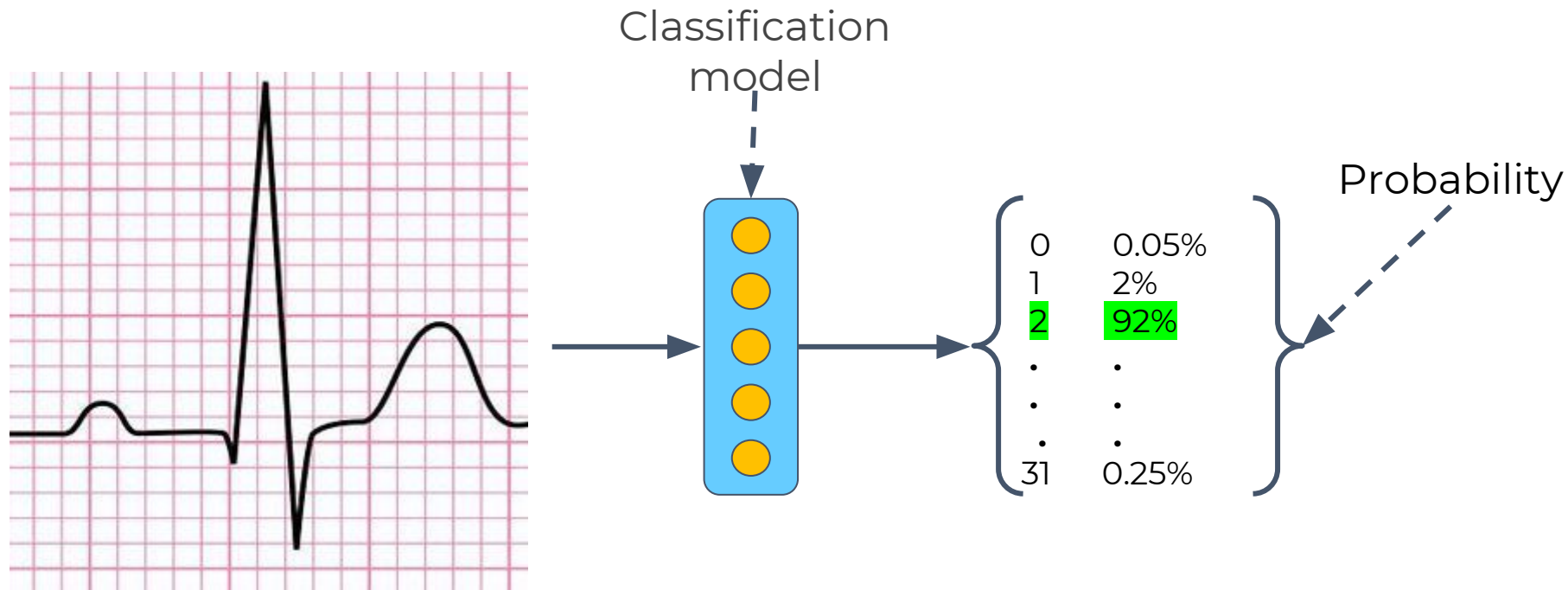
- **Useful** to avoid overfitting on the training data
- **Examples of techniques:**
  - Dropout
  - Batch Normalization
  - Instance Normalization
  - ...
- All are **readily available** in PyTorch.



# OMsignal Project Multi-Task Learning

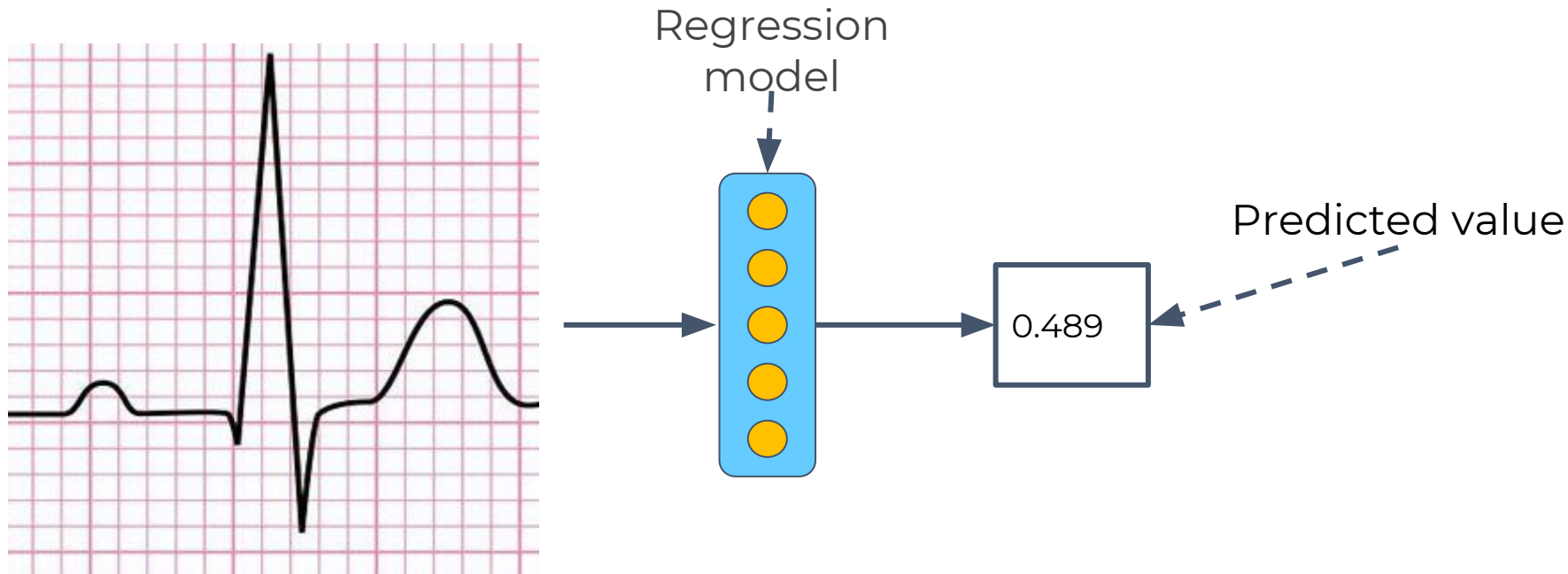


# User Identification Task

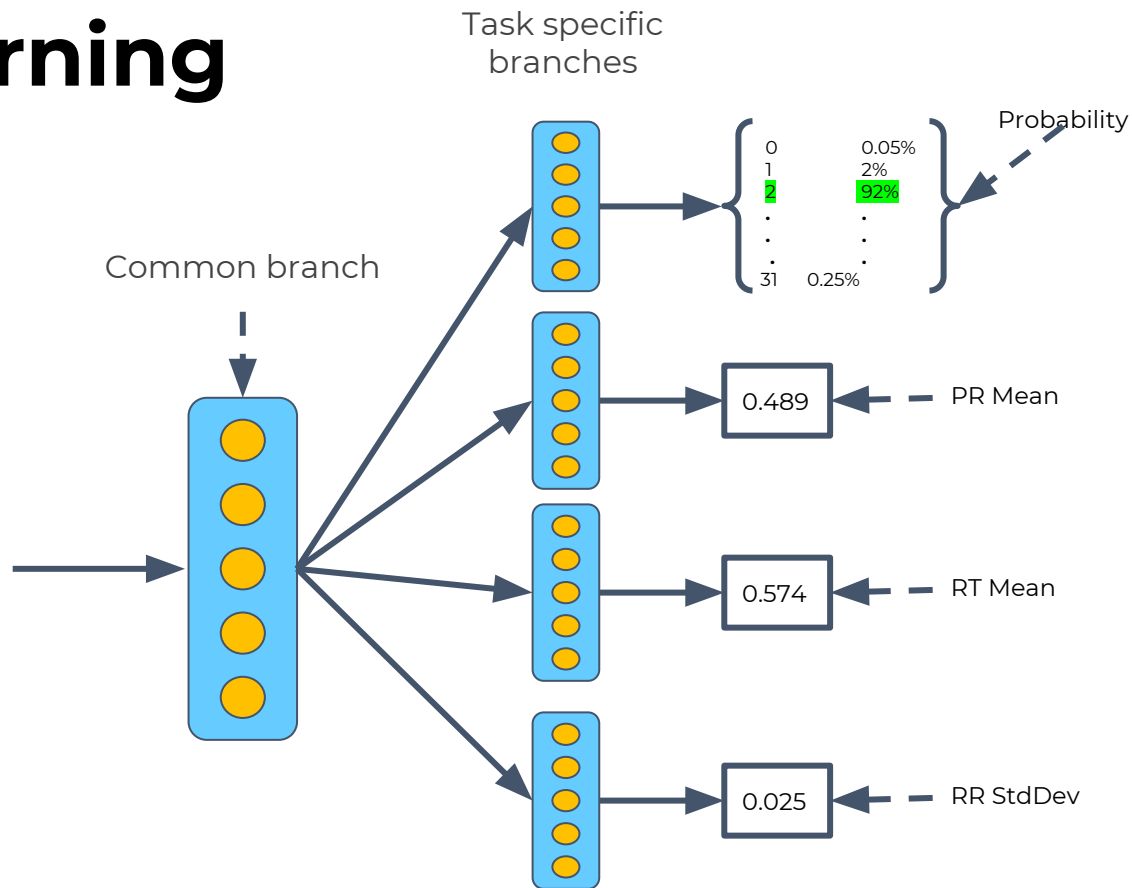
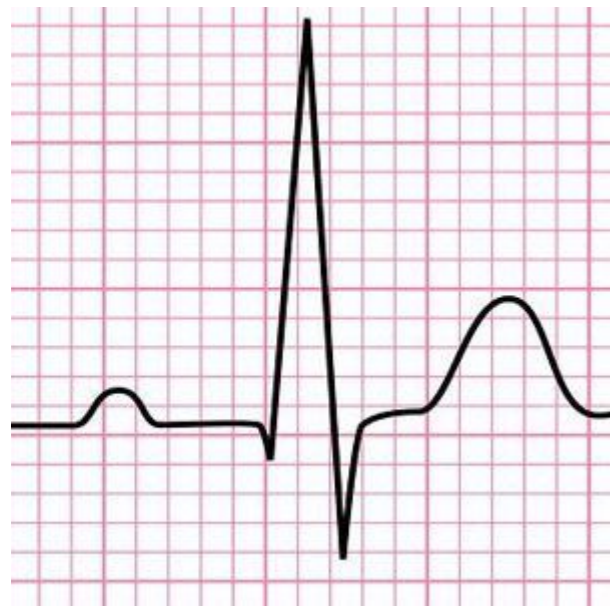


# Regression Tasks

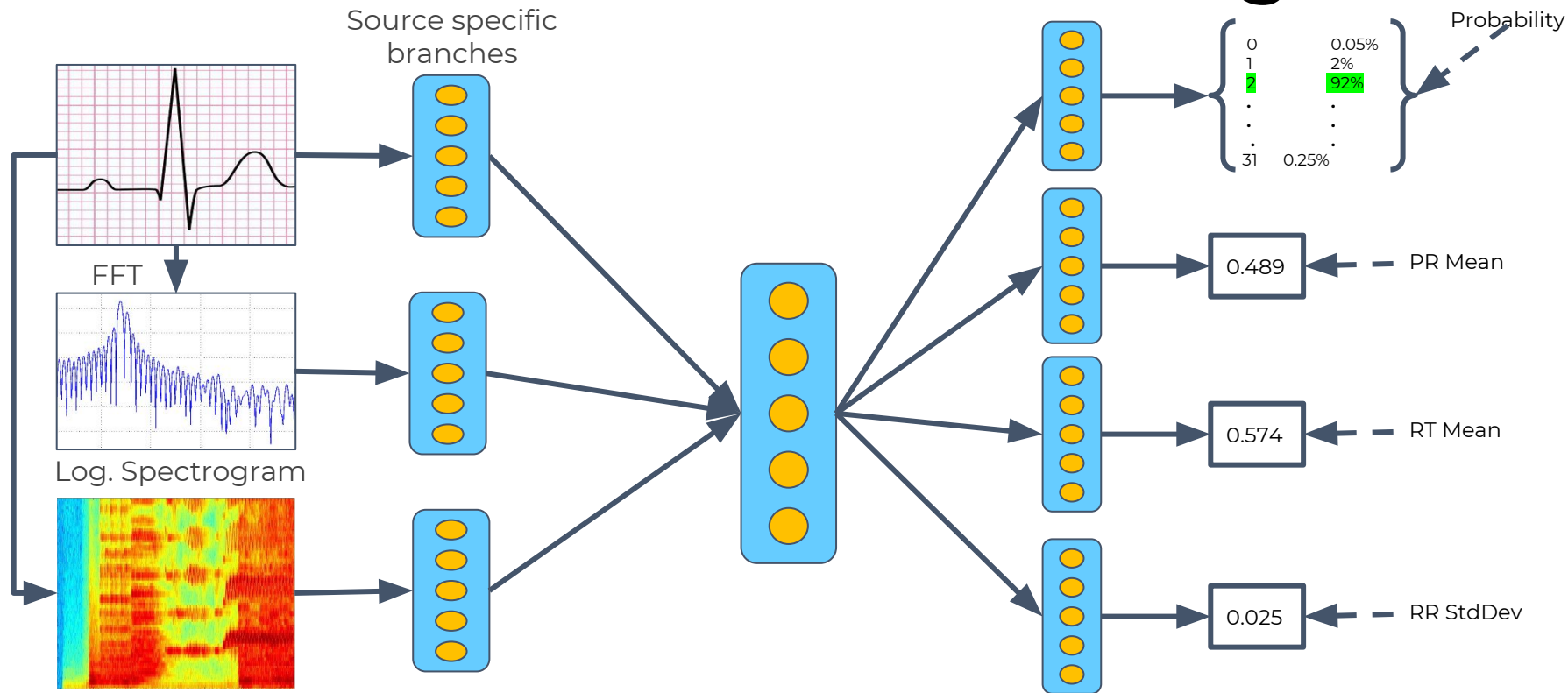
- Applicable for the prediction of the fiducial point statistics: **PR Mean, RT Mean, RR StdDev**



# Multi-task Learning

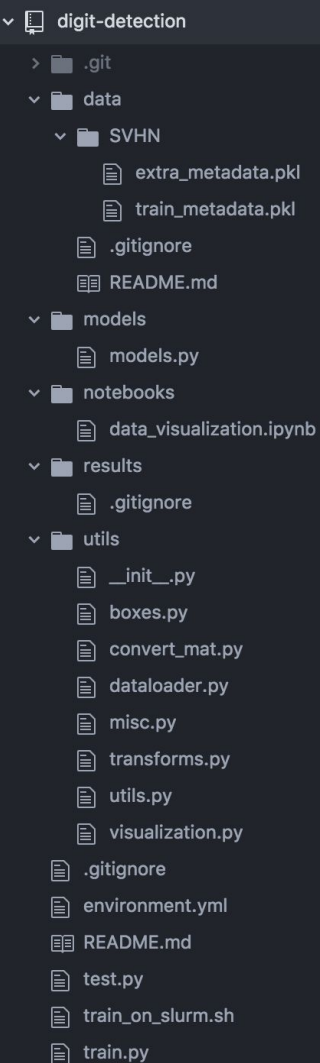


# Multi-source Multi-task Learning



# Homework

**Implement your model**

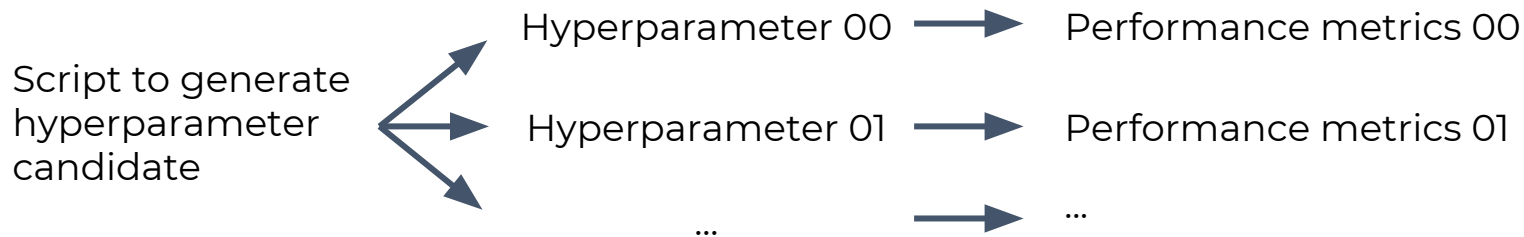


# How the code should be organized

- Have separate folders: config, data, models, notebooks, results, trainer, utils ...
- In each folder, separate the code into several files to simplify reading
- Use meaningful names for your directories and files
- Avoid code duplication => use object-oriented programming (e.g. parent and child classes)

# Running a Deep Learning experiment

- Exploration of hyperparameters to obtain the best model



Each experiment must produce enough logs to:

- diagnose errors and suspicious behaviour
- allow the selection of a set of hyperparameters that is considered the best to generalize to new data



# Implementation - To keep in mind (1)

## ***General advice***

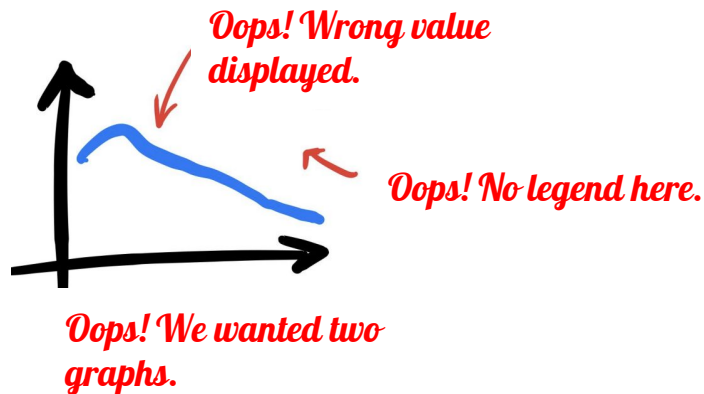
- See what other people are doing and use it as inspiration.
- Develop your own style from this.

## ***Implementation of a deep learning model***

- Reuse as much as possible models already available online. Search for "model zoo".

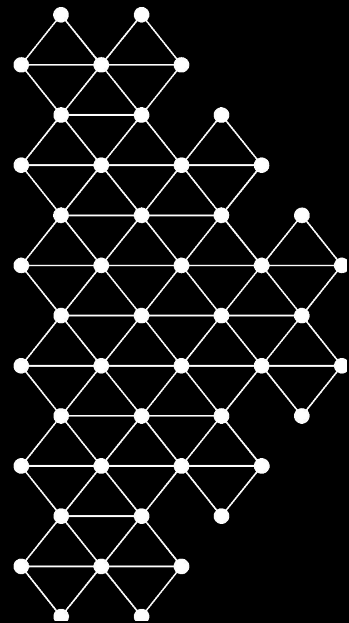
# Implementation - To keep in mind (2)

- When you run several large-scale experiments, you must store the data necessary to generate graphs / figures in case you need to modify / update them.



- An experiment that can be interrupted in the middle and then restarted without being affected by the interruption is much easier to manage.

# OMsignal Project To go further

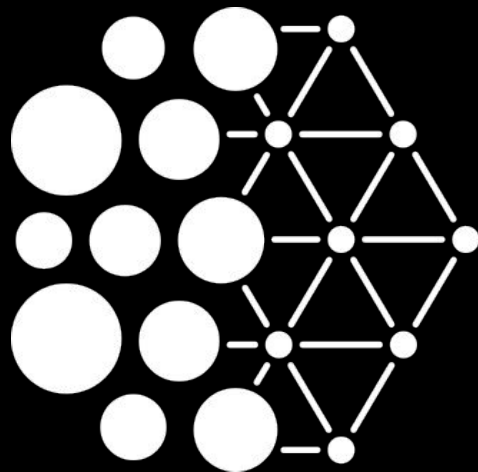


# Dealing with Unlabeled data

**Goal:** Efficient way to integrate knowledge from the unlabeled data

- **Unsupervised + Supervised Learning**
  - **Step 1:** Auto-Encoder like model training
  - **Step 2:** Supervised training based on features extracted from the trained encoder
- **Semi-supervised Learning**
  - One step process.
  - Possible approaches (combined with a supervised loss):
    - Reconstruction loss (unlabeled data) - auto encoder
    - Regularization loss (unlabeled data) based on some assumptions (e.g. invariance of the output to small amounts of noise added to the input signal)

Quebec  
Artificial  
Intelligence  
Institute



Mila

Université   
de Montréal