

# HDA COURSE PROJECTS

Michele Rossi

[rossi@dei.unipd.it](mailto:rossi@dei.unipd.it)

Francesca Meneghelli

[meneghelli@dei.unipd.it](mailto:meneghelli@dei.unipd.it)



DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE

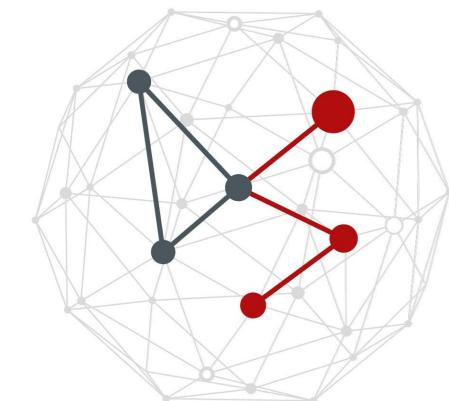


DIPARTIMENTO  
**MATEMATICA**

1222-2022  
**800** ANNI



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA



# Outline



- Exam
  - Project based - working in groups (2 people group max.)
    - 1) Delivery of a written report
    - 2) 20 minutes (max) talk, possibly showing running code
  - Proposed projects
    - 1) P1: Speech command recognition (keyword spotting)
    - 2) P2: Activity recognition 1: human motion
    - 3) P3: Activity recognition 2: human motion & interaction with env.
    - 4) P4: Sleep posture monitoring
    - 5) P5: Environmental sound classification
    - 6) P6: Lymphoma subtype classification
    - 7) P7: ECG beat classification
  - Exam dates for 2020

# Groups

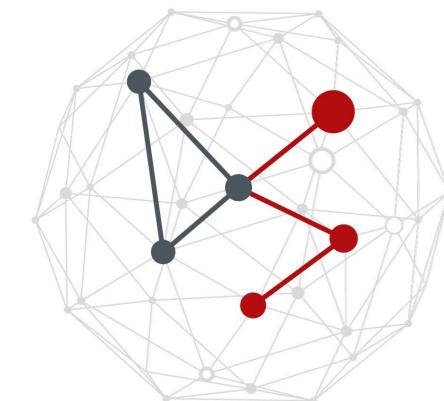
- Max. two people per group
  - You choose your work mate
  - Contribution of each group member
    - Should be stated in project report
    - Should emerge from the final talk

# Delivery of technical report

- Ideally
  - Should happen before the talk
- Final grade
  - Assigned only upon evaluating:
    - 1) written report,
    - 2) presentation

# PROJECT 1

---



# Project no. 1 “speech recognition”

- Reference paper

[Sainath15] Tara N. Sainath, Carolina Parada, Convolutional Neural Networks for Small-footprint Keyword Spotting, INTERSPEECH, Dresden, Germany, September 2015.

[Warden18] Pete Warden, Speech Commands: A Dataset for Limited-Vocabulary Speech Recognition, arXiv:1804.03209, April 2018.

<https://arxiv.org/abs/1804.03209>

- The authors are from Google Inc.
- Reference dataset recently released by Google [Warden18]

# Project no. 1 “speech recognition”

- Reference dataset for small-footprint keyword spotting (KWS)
  - Released on [August 2017](#)
  - **65,000** one second long utterances of **30 words**
  - by thousands of different people
  - released under creative commons 4.0 license
  - collected by AYI (you can also contribute to make it grow):  
[https://aiyprojects.withgoogle.com/open\\_speech\\_recording](https://aiyprojects.withgoogle.com/open_speech_recording)

## Google blog

<https://research.googleblog.com/2017/08/launching-speech-commands-dataset.html>

## Speech dataset (2.11 GB uncompressed)

[http://download.tensorflow.org/data/speech\\_commands\\_v0.02.tar.gz](http://download.tensorflow.org/data/speech_commands_v0.02.tar.gz)

# Project no. 1 “speech recognition”

- Approaches for implementing a **KWS** engine
  - **LVCSR based KWS** - This approach uses a two-stage process. In the first stage, the transcription of the speech into words is done using a **Large Vocabulary Continuous Speech Recognition (LVCSR)** engine, outputting formatted text. In the second stage, a textual search for the key-words within the text is performed. Using this approach, results from LVCSR and the text search are combined to spot the key-words
  - **Phoneme Recognition based KWS** - This approach also uses a two-stage process. In the first stage, the speech is transformed to a sequence of phonemes. In the second stage, the application searches for phonetically transcribed key-words in the phoneme sequence obtained from the first stage
  - **Word Recognition based KWS [Sainath15]** - This approach searches for the key-words in a **one stage operation**. The recognition is phoneme-based and the KWS engine looks for the keyword in the speech stream based on a target sequence of phonemes representing the key-word

# Project no. 1 “speech recognition”

- CNN model from [Sainath15]
  - Features are obtained from raw audio data
  - 40 dimensional log Mel filterbanks coefficients
    - audio frame length 25 ms
    - with a 10 ms time shift
  - At every new audio frame
    - Feature vector is obtained
    - And stacked with 23 frames to the left and 8 to the right (32 frames total)
    - This returns 32 frames at a time, spanning over  $31 \times 10 \text{ ms} + 25 \text{ ms} = 0.335 \text{ s}$
  - A Convolutional Neural Network (CNN) is used to detect words
  - Input to the CNN is a matrix of size  $t \times n = 32 \times 40 = 1,280$  elements
    - t represents the number of elements in time (number of audio frames)
    - n represents the number of elements in the frequency domain (Mel features)

# Project no. 1 “speech recognition”

- CNN model from [Sainath15]
  - 27-44% improvement for KWS with respect to traditional neural networks
- Paper focus is on
  - Devise CNN architectures with small memory footprint
  - Playing with CNN parameters (number of kernels, strides, pooling, etc.)

# Project no. 1 “speech recognition”

- Possible project developments
  - Experiment with different audio features (+)
    - Type of coefficients (e.g., discrete Wavelet transform)
    - Design of Mel filterbanks, e.g.,  
<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.59.1426&rep=rep1&type=pdf>
  - Play with a standard/deep CNN using (++)
    - dropout
    - regularization
    - different gradient descent techniques
  - Investigate recent/new ANN architectures (+++)
    - Autoencoder-based (CNN/RNN autoencoder + following SVM)
    - Attention mechanism and/or inception-based CNN networks
    - Comparison of different architectures: memory vs accuracy

# Project no. 1 - useful resources

## Recent developments

[Chorowski15] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, Y. Bengio, [Attention-Based Models for Speech Recognition](#), Conference on Neural Information and Processing Systems (NIPS), Montréal, Canada, 2015.

[Tang18] R. Tang and J. Lin, [Deep residual learning for small-footprint keyword spotting](#), in IEEE ICASSP, Calgary, Alberta, Canada, 2018.

[Andrade18] D. C. de Andrade, S. Leo, M. L. D. S. Viana, and C. Bernkopf, [A neural attention model for speech command recognition](#), arXiv:1808.08929, 2018. <https://arxiv.org/pdf/1808.08929.pdf>

White Paper: “Key-Word Spotting - The Base Technology for Speech Analytics”

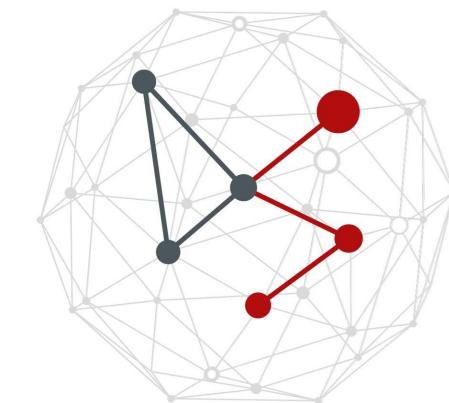
<https://pdfs.semanticscholar.org/e736/bc0a0cf1f2d867283343faf63211aef8a10c.pdf>

Example code:

[https://github.com/tensorflow/tensorflow/tree/master/tensorflow/examples/speech\\_commands/](https://github.com/tensorflow/tensorflow/tree/master/tensorflow/examples/speech_commands/)

# PROJECT 2

---



# Project no. 2 “activity recognition”

- **Reference paper**

[Frank10] K. Frank, M.J.V. Nadales, P. Robertson, M. Angermann,  
**Reliable Real-Time Recognition of Motion Related Human Activities  
Using MEMS Inertial Sensors**, Proceedings of the International  
Conference ION GNSS, Portland, Oregon, USA, 2010.

- **Basic sensor:** Inertial Measurement Unit (IMU)

- **Reference dataset (303 MB uncompressed)**

- German Aerospace Center (DLR, Deutsches Zentrum für Luft und Raumfahrt)  
[http://www.dlr.de/kn/desktopdefault.aspx/tabcid-8500/14564\\_read-36508/](http://www.dlr.de/kn/desktopdefault.aspx/tabcid-8500/14564_read-36508/)

# Why is activity recognition important

- **Navigation systems**
  - adapt to user movement
  - e.g., predict direction and only use that portion of the map(s)
  - put the system into power saving mode when there is no mobility
- **First responders**
  - security personnel, firefighters
  - e.g., who has to be assisted first
- **Assisted living**
  - react to reduced activity levels
  - unusual mobility patterns
  - user motion-aware services and/or environments
- **Rehabilitation**
  - measure recovery of motor functions
  - measure effectiveness of rehabilitation

# Activity recognition

- Application domains
  - Navigation, first responders, assisted living, rehabilitation
  - In most of these cases the use of cameras is not possible
  - The system has to be
    - unobtrusive
    - lightweight
    - portable
- [Frank10] explores using wearable IMU
  - To detect basic activities, such as, e.g., sitting, standing, walking, running, Jumping, falling, lying

# Project no. 2 “activity recognition”

- **IMU sensor:** xsens MTx IMU
  - <https://www.xsens.com/products/mtw-awinda/>
- **Dataset:** collected from 16 males & females subjects aged between 23 and 50
- Collected from a single IMU sensor on the belt (front) of the user
- Contains 4.5 hours of labeled activities
  - **9 axes** = 3 accelerometer, 3 gyroscope, 3 magnetometer

## Tracked activities

'RUNNING' = "running"

'WALKING' = "walking"

'JUMPING' = "jumping"

'STNDING' = "standing"

'SITTING' = "sitting"

'XLYINGX' = "lying"

'FALLING' = "falling"

'TRANSUP' = "getting up" i.e.: from lying (or sitting) to standing

'TRANSDW' = "going down" i.e.: from standing to sitting

'TRNSACC' = "accelerating"

'TRNSDCC' = "decelerating"

'TRANSIT' = "other transition or irrelevant information"

Activity	Duration (minutes)
Standing	107
Sitting	55
Lying	25
Walking	70
Running	15
Jumping	7
Falling	2

# Project no. 2 “activity recognition”

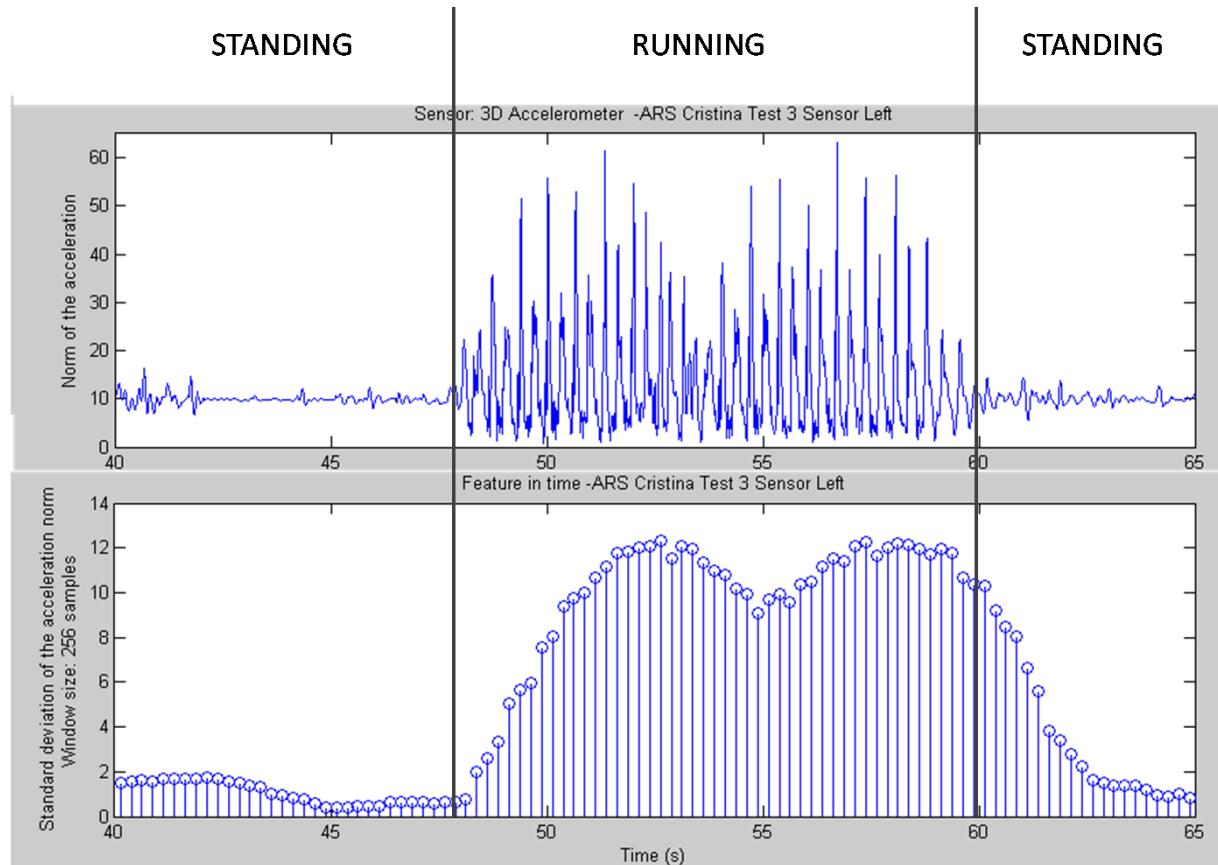
- **Dataset format**
  - Text file: 10 columns
  - 1st column: **TIME** extracted from the sensor in seconds
  - 2nd column (**accX**) is the acceleration in the X axis measured by the sensor
  - 3rd column (**accY**) is the acceleration in the Y axis measured by the sensor
  - 4th column (**accZ**) is the acceleration in the Z axis measured by the sensor
  - 5th column (**gyroX**) is the angular velocity in the X axis
  - 6th column (**gyroY**) is the angular velocity in the Y axis
  - 7th column (**gyroZ**) is the angular velocity in the Z axis
  - 8th column (**magX**) is the magnetic field in the X axis
  - 9th column (**magY**) is the magnetic field in the Y axis
  - 10th column (**magZ**) is the magnetic field in the Z axis

# Signal features

Example from [Frank10]

$$\mathbf{a}(t) = [a_x(t) \ a_y(t) \ a_z(t)]^T$$

$$|a(t)| = \sqrt{a_x(t)^2 + a_y(t)^2 + a_z(t)^2}$$



modulus of accelerometer  
signal,  $|a|$

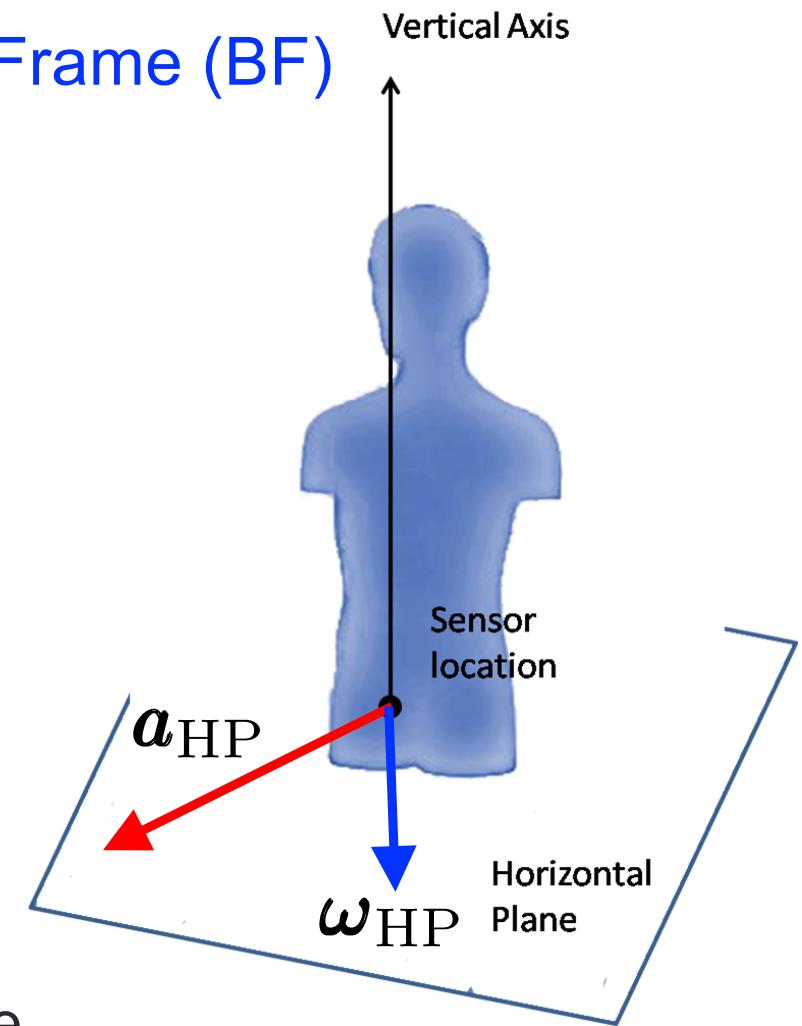
standard deviation of  $|a|$

# Reference coordinate system

- Referred in [Frank10] to as Body Frame (BF)
- Vertical axis
  - Aligned with body
- Horizontal Plane (HP)
  - Orthogonal to vertical axis

Acceleration and angular velocity

- Projected onto horizontal plane



# Reference coordinate system

- Sensor Frame (SF)
  - Reference system defined by IMU orientation on the body
- Body Frame (BF)
  - As in previous slide, aligned with body from feet to head
  - In [Frank10] SF data is moved onto BF
- Global frame
  - Vertical axis aligned with (constant) gravity vector
  - Vertical axis pointing towards earth's center
  - Reference system remains unchanged
    - even when user changes their position in space

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed
  - “vert”: vertical gravity direction
  - “horiz”: horizontal plane
  - LPF: low pass filtered
  - BPF: band pass filtered
  - RMS: root mean square
  - IQR: interquartile range:
    - difference between 25<sup>th</sup> and 75<sup>th</sup> percentile
  - MFC: main frequency component
    - computed using FFT

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{\frac{a}{a_{vert}_{GF}}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6–4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6–4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Max Frequency Component
  - Computed through FFT
  - To identify walking and running
  - Helps discriminate between
    - *falling / jumping vs running*

Feature No.	Definition	Window size
1	$\text{MAX}_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$\text{MAX}_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$\text{RMS}_{a_{vert}_{BF}}$	128
8	$\text{IQR}_{ \omega_{horiz}_{BF} }$	128
9	$\overline{\text{a}_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$\text{IQR}_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 \text{ Hz}})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 \text{ Hz}})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 \text{ Hz}})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Inter Quartile Range

- Difference between 25<sup>th</sup> and 75<sup>th</sup> percentile
- Help distinguish between *jumping & falling*

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Features 15-17

- Represent the norm of the acceleration
  - Within some pre-define frequency bands

- Feature 15

- Can distinguish between
    - *walking / jumping vs running*

- Features 16 & 17

- Can distinguish between
    - *running and jumping*

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Horizontal acceleration of body frame

- Feature 1, helps distinguish between:
  - *static* and *dynamic* activities
  - it is particularly high for falling

- Mean value distinguishes reliably

- static activities

- Standard deviation

- helps distinguish *jumping*, *falling*, *running*

- Feature 8 (IQR)

- High values only reached for *falling*

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Vertical acceleration

- to distinguish between
- MAX:
  - *jumping, falling and walking*
- MEAN:
  - *standing, sitting and lying*
- STDEV:
  - all *dynamic* activities
- RMS:
  - all *static* activities

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

- In [Frank10], 19 predefined features are computed

- Correlation coefficient, feature 18

- *walking, running and jumping*: high value
- other activities: no consistent patterns

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Signal features used in [Frank10]

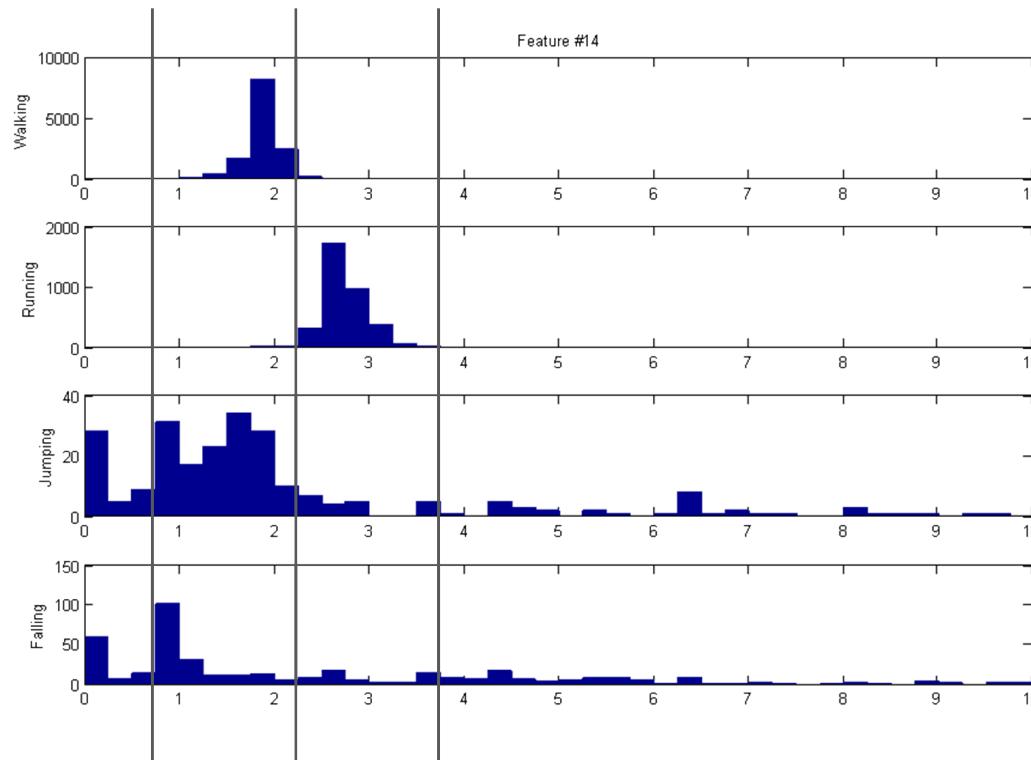
- In [Frank10], 19 predefined features are computed

- Vertical acceleration

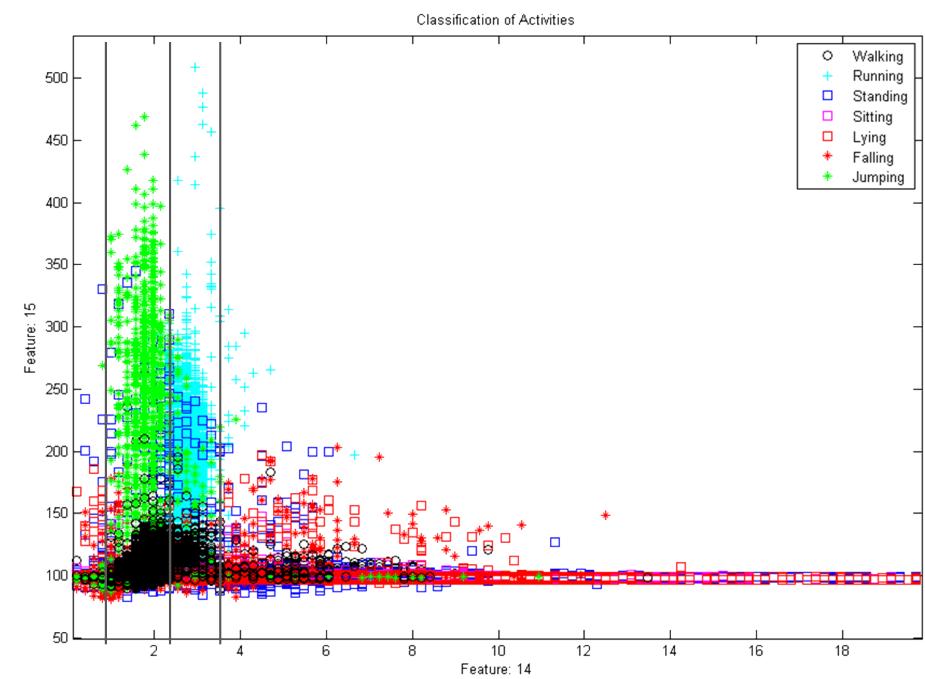
- To detect free fall phases during
  - *jumping and falling*

Feature No.	Definition	Window size
1	$MAX_{ a_{horiz}_{BF} }$	128
2	$\overline{ a_{horiz}_{BF} }$	128
3	$\sigma_{ a_{horiz}_{BF} }$	128
4	$MAX_{a_{vert}_{BF}}$	128
5	$\overline{a_{vert}_{BF}}$	128
6	$\sigma_{a_{vert}_{BF}}$	128
7	$RMS_{a_{vert}_{BF}}$	128
8	$IQR_{ \omega_{horiz}_{BF} }$	128
9	$\overline{a_{vert}_{GF}}$	32
10	$\overline{ a }$	32
11	$\overline{ a }$	512
12	$\sigma_{ a }$	256
13	$IQR_{ a }$	128
14	$MFC_{ a }$	128
15	$\hat{E}( a _{LPF < 2.85 Hz})$	128
16	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	64
17	$\hat{E}( a _{BPF 1.6-4.5 Hz})$	512
18	$\rho_{a_{vert}_{BF},  a }$	128
19	$att_{ a_{horiz}_{BF} ,  a_{vert}_{BF} }$	64

# Quantization of features



(a)



(b)

- **Histograms of feature values**
  - Example for Feature 14
    - visual inspection suggests splitting the feature range into 4 regions

# Approach in [Frank10]

- Uses manually crafted feature vectors
- Sampling (and label) frequency: 100 Hz
- Feature vector computation and classification: 4 Hz
- Uses *static* and *dynamic* Bayesian Networks
  - Dynamic BNs: Hidden Markov Models (HMMs)
  - Hidden states:
    - Corresponds to the monitored activities
  - Emission probabilities:
    - Prob. of generating a certain *feature vector* given activity
    - **Feature vector:** has 19 entries (defined at each sampling time)
  - Transition probabilities:
    - Prob. of *next activity* given *current one*

# Project proposal

- **Features (possibilities)**

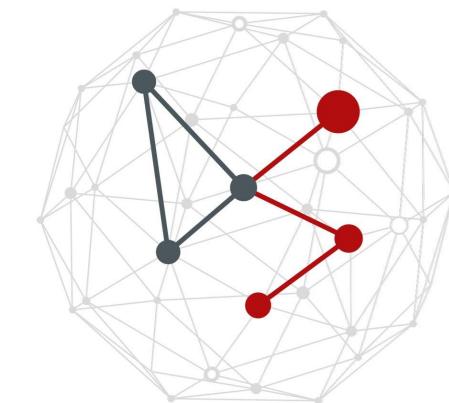
1. use features defined as in [Frank10] or,
2. compute features using FFT / DCT as we do for audio signals, or
3. use raw signals with automatic feature extraction

- **Classification architecture**

- Use a CNN architecture to classify the type of motion / activity type
  - plain Convolutional Neural Network (CNN)
  - CNN/RNN autoencoder + SVM
  - ...

# PROJECT 3

---



# Project no. 3 “advanced activity recognition”

## Reference paper

[Chavarriaga13] Chavarriaga, H. Sagha, A. Calatroni, S.T. Digumarti, G. Tröster, J.R. Millán, D. Roggen, [The Opportunity challenge: A benchmark database for on-body sensor-based activity recognition](#), Elsevier’s Pattern Recognition Letters, 2013.

## Large dataset (901 MB uncompressed)

<https://archive.ics.uci.edu/ml/datasets/opportunity+activity+recognition>

- Contains recordings of [on-body sensors](#)
- Subjects perform activities of [daily living](#)
  - Simple motion primitives and complex gestures

# High level description of the dataset

- More than **27,000 atomic motion activities collected**
  - In a sensor rich environment
- **Recordings involved**
  - **12 subjects,**
  - **15 networked sensor systems, with 72 sensors of 10 modalities**
  - Sensors integrated in the environment, in the objects and on the body
- **Aim**
  - Maximize the number of activities monitored
  - While keeping their execution naturalistic

## Monitored activities (1/2)

- Two types of activities: ADL and drill runs
- Activities of Daily Living (ADL)
  - For each subject, we recorded six different runs
  - Daily morning activities:
    - Activities start with the subject lying in the deckchair
    - Then she/he gets up, checking the object in the drawers and shelves
    - The subject leaves the room,
    - When she/he returns, the subject prepares coffee
    - Then prepares a sandwich with “salami” and cheese
    - Then she/he cleans the room,
    - Putting back objects in their original location or in the dishwasher
    - Then she/he goes back to the deckchair

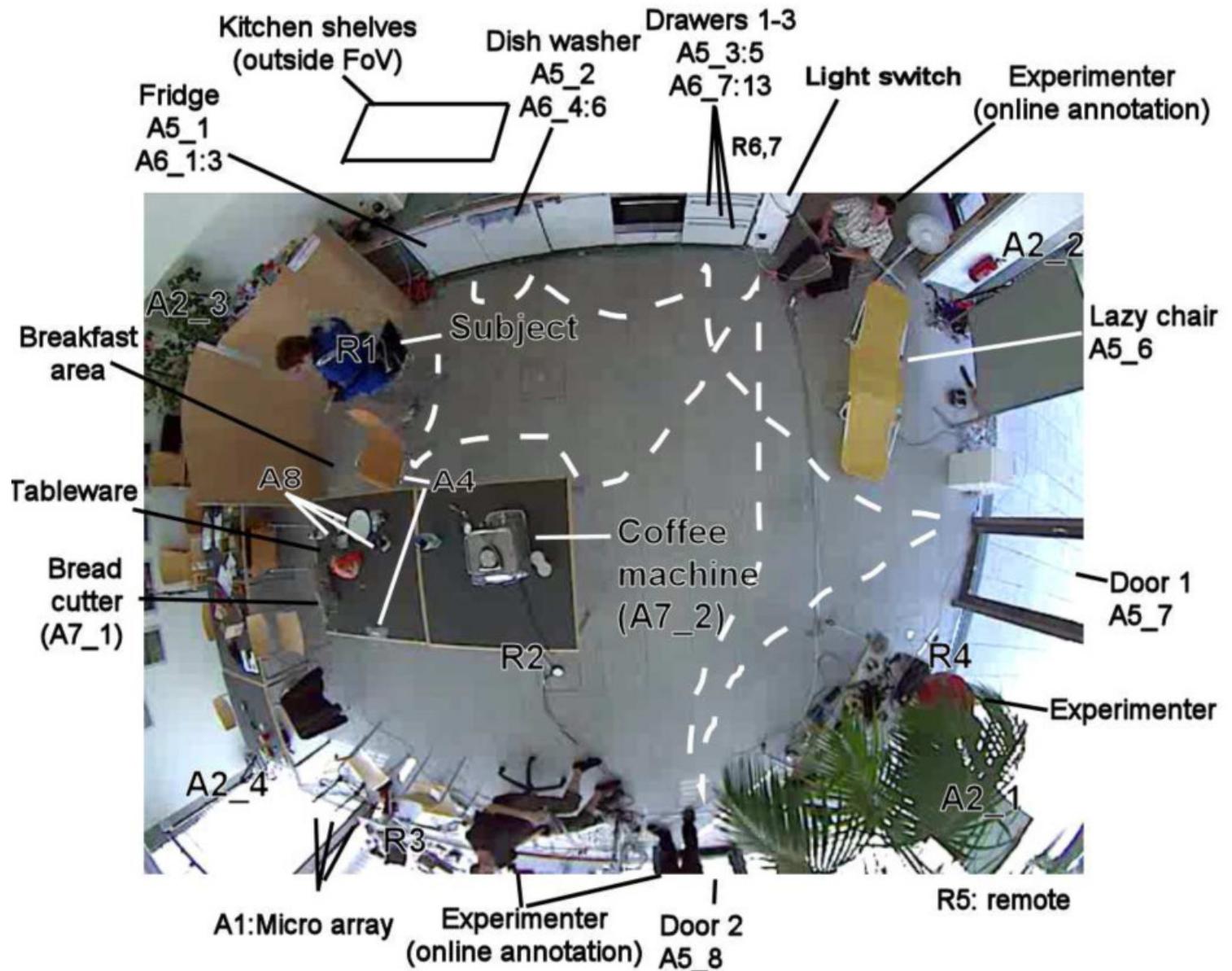
## Monitored activities (2/2)

- Drill runs
  - Intended to record a large set of activities
  - The subject performs 20 repetitions of the sequence
    1. Open and close the fridge
    2. Open and close the dishwasher
    3. Open and close 3 drawers (at different heights)
    4. Open and close door 1
    5. Open and close door 2
    6. Turn on and off the lights
    7. Clean table
    8. Drink (standing)
    9. Drink (sitting)

# The monitored environment

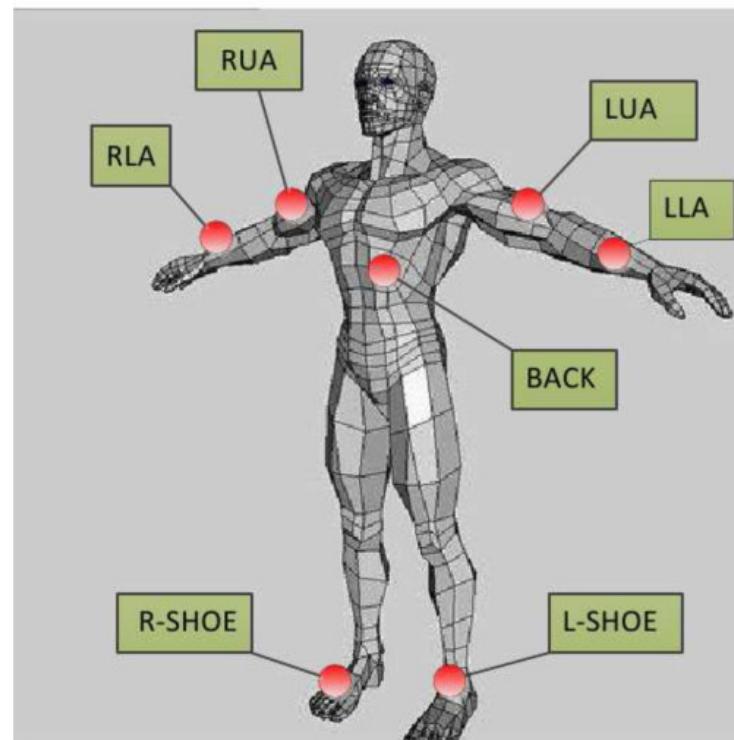
A room simulating  
a **studio flat** with:

A kitchen  
A deckchair  
A Table  
A chair  
A coffee machine  
Doors

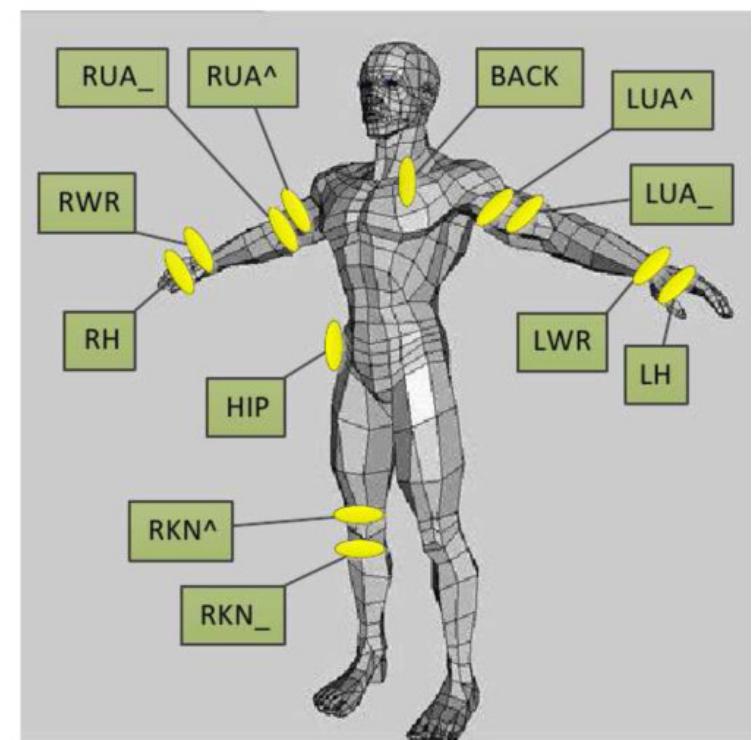


# On body sensors

- 5 Xsense IMUs (acc, gyro, mag) mounted on a custom-made motion jacket
- 2 InertialCube3 IMUs mounted on shoes
- 12 accelerometer (3 axes) sensors mounted on the limbs



● = Complete Inertial Measurement Unit



● = Triaxial Accelerometer

# Tasks

Four different tasks were defined, targeting

- Modes of locomotion (task A)
- Activity spotting (task B1)
- Arm gestures (task B2)
- Robustness to noise (task C)

Modes of locomotion	Gestures
Stand (1093)	Open Dishwasher (50)
Sit (1095)	Close Dishwasher (56)
Walk (90)	Open Fridge (129)
Lie (40)	Close Fridge (133)
<i>Null</i>	<i>Null</i>

More gestures		
Open Drawer1 (50)	Open Drawer2 (44)	Open Drawer3 (56)
Close Drawer1 (49)	Close Drawer2 (44)	Close Drawer3 (57)
Open Door1 (45)	Open Door2 (43)	Move Cup (184)
Close Door1 (39)	Close Door2 (41)	Clean Table (33)

# Task A

## Multimodal activity recognition: modes of locomotion

- **GOAL**
  - classify modes of locomotion from the full set of body-worn sensors
  - Possible modes of locomotion (classes) are:
    - standing
    - sitting
    - walking
    - lying
  - It is a **4-class** continuous classification problem

# Task B1

## Activity spotting & segmentation

- **GOAL**
  - The dataset is labeled, that is: activities are already segmented into classes
  - Realistic deployments must detect when no relevant action is performed
    - This is referred to as the NULL class
    - This is a 2-class segmentation problem (NULL vs any activity)
    - The full set of body-worn sensors is considered for this task

# Task B2

## Multimodal activity recognition: gestures

- **GOAL**
  - The task concerns the recognition of the 17 right arm gestures
  - This is a 17-classes segmentation and classification problem
  - The full set of sensors is used for this task
  - Labelled data is provided for the 17 gestures

# Task C

## Robustness to noise

- **GOAL**
  - Realistic applications are prone to noise
  - For this task *rotational noise* has been added to the testing dataset
  - The gestures to be recognized are the same as for task B2
  - Only the motion sensors from the jacket are considered

# Features and missing values (1/3)

- **Input space dimensionality**
  - Each sensor axis is treated separately
  - 12 accelerometers on the limbs:  $3 \times 12 = 36$  values
  - 5 IMUs on the jacket:  $9 \times 5 = 45$  values
  - 2 shoe sensors, 16 values each:  $16 \times 2 = 32$  values
  - **Total number of inputs is:  $36 + 45 + 32 = 113$  values**
- Missing data
  - Due to disconnections of wireless sensors
  - Previous value is used to replace missing data points
- Features
  - Data was recorded continuously and it is not segmented
  - Sliding windows of 500 ms, with time step of 250 ms are defined
  - **Feature:** average of an input within each time window

## Features and missing values (2/3)

- Input space dimensionality
  - Each sensor axis is treated separately
  - 12 accelerometers on the limbs:  $3 \times 12 = 36$  values
  - 5 IMUs on the jacket:  $9 \times 5 = 45$  values
  - 2 shoe sensors, 16 values each:  $16 \times 2 = 32$  values
  - Total number of inputs is:  $36 + 45 + 32 = 113$  values
- Missing data
  - Due to disconnections of wireless sensors
  - Previous value is used to replace missing data points
- Features
  - Data was recorded continuously and it is not segmented
  - Sliding windows of 500 ms, with time step of 250 ms are defined
  - **Feature:** average of an input within each time window

# Features and missing values (3/3)

- Input space dimensionality
  - Each sensor axis is treated separately
  - 12 accelerometers on the limbs:  $3 \times 12 = 36$  values
  - 5 IMUs on the jacket:  $9 \times 5 = 45$  values
  - 2 shoe sensors, 16 values each:  $16 \times 2 = 32$  values
  - Total number of inputs is:  $36 + 45 + 32 = 113$  values
- Missing data
  - Due to disconnections of wireless sensors
  - Previous value is used to replace missing data points
- Features
  - Data was recorded continuously, and it is not segmented
  - Sliding windows of 0.5 s, with time step of 0.25 s are defined
  - **Feature (def):** average of an input within each time window

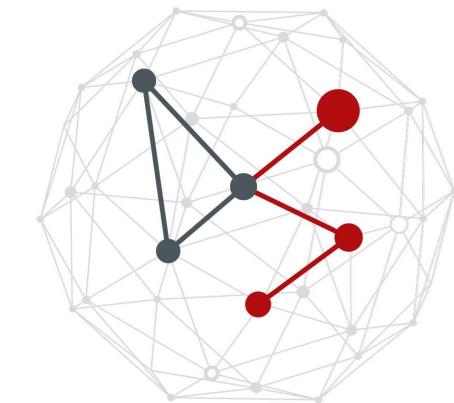
# Project no. 3 “advanced activity recognition”

- Project objective
  - Implement a method to **fill missing data**
  - Define more complex features
  - Use **CNN+RNN** for classification
  - Check classification performance
    - With all on-body sensors
    - Only using a subset of sensors

Ideally, with neural networks, we would like to get decent classification performance, i.e., at least comparable with that in **[Chavarriaga13]**, using as few IMUs as possible

# PROJECT 4

---



# Project no. 4 “sleep posture monitoring”

## Reference paper

[Pouyan17] M. B. Pouyan, J. Birjandtalab, M. Heydarzadeh, M. Nourani and S. Ostadabbas, [A pressure map dataset for posture and subject analytics](#), in Proceedings of the IEEE EMBS International Conference on Biomedical & Health Informatics (BHI), Orlando, FL, 2017.

## PmatData dataset (102.3 MB uncompressed)

<https://physionet.org/content/pmd/1.0.0/>

Contains in-bed posture pressure data

- multiple adult participants
- two different types of pressure sensing mats

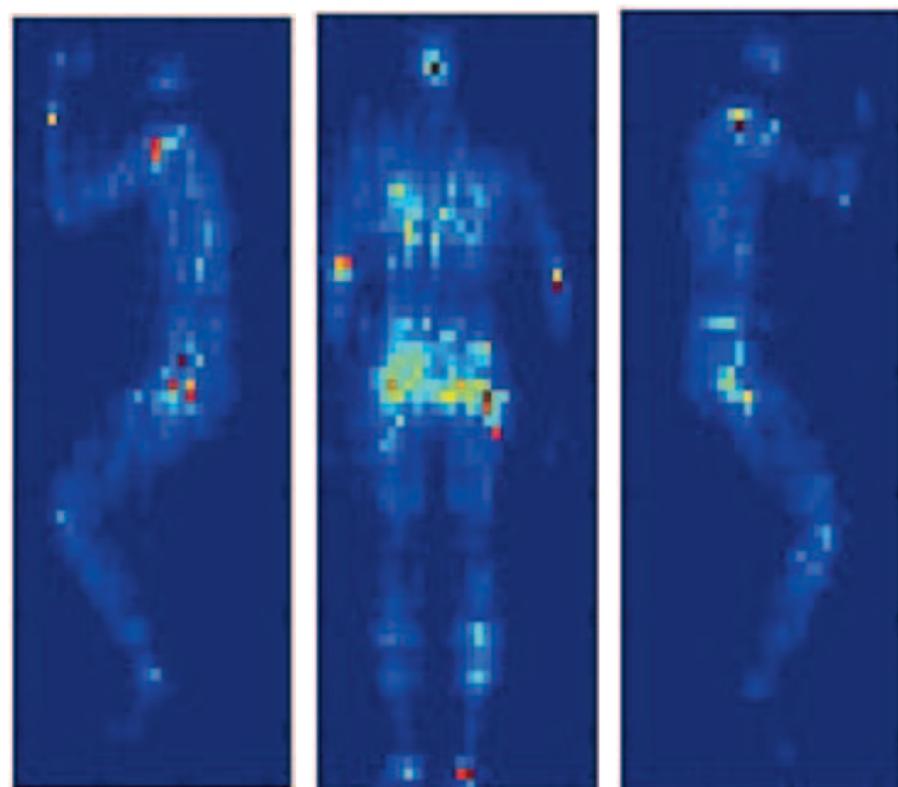
# High level description of the dataset

- Pressure data from two separate experiments
- **Experiment 1:** 13 participants
  - Size of pressure mat is 32\*64: 2048 values per acquisition
  - Sampling rate: 1Hz
  - Sample values range: [0-1000]
  - 17 files for each subject (8 standard postures and 9 additional states)
  - Each file includes around 2 minutes of acquisitions

Index	Posture	Bed Inclination (degree)	Body-roll (degree)	Symbol	Duration	Spec's of mat
1	Supine	0	0		2 mins	Vista
2	Right	0	0		2 mins	Vista
3	Left	0	0		2 mins	Vista
4	Right	0	30 (1 wedge)		2 mins	Vista
5	Right	0	60 (2 wedges)		2 mins	Vista
6	Left	0	30 (1 wedge)		2 mins	Vista
7	Left	0	60 (2 wedges)		2 mins	Vista
8	Supine	0	0		2 mins	Vista
9	Supine	0	0		2 mins	Vista
10	Supine	0	0		2 mins	Vista
11	Supine	0	0		2 mins	Vista
12	Supine	0	0		2 mins	Vista
13	Right Fetus	0	0		2 mins	Vista
14	Left Fetus	0	0		2 mins	Vista
15	Supine	30	0		2 mins	Vista
16	Supine	45	0		2 mins	Vista
17	Supine	60	0		2 mins	Vista

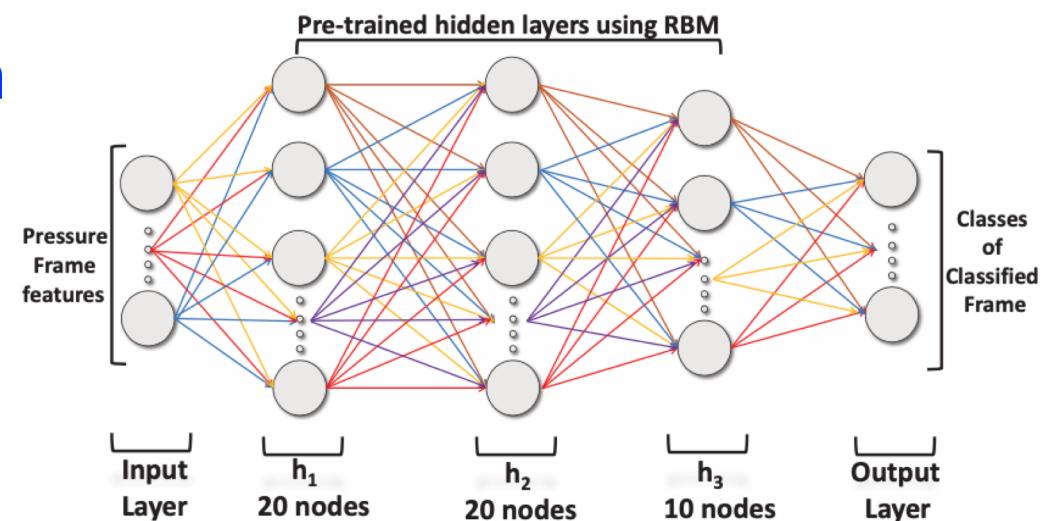
# High level description of the dataset

- **Experiment 2:** 8 participants
  - The data is collected for both sponge and air mattresses
  - Size of pressure mat is 27\*64: 1728 values per acquisition
  - Each file contains the average of around 20 acquisitions
  - 29 different states of 3 standard postures
  - Sample values range: [0-500]
  - Sampling rate: 1Hz



# In [Pouyan17]

- Subject identification in three standard postures:
  - right side
  - supine
  - left side
- Main idea: each subject has a personalized sleeping pattern in each posture
- Architecture: one FFNN for each posture
- Manual feature extraction
  - 18 statistical features



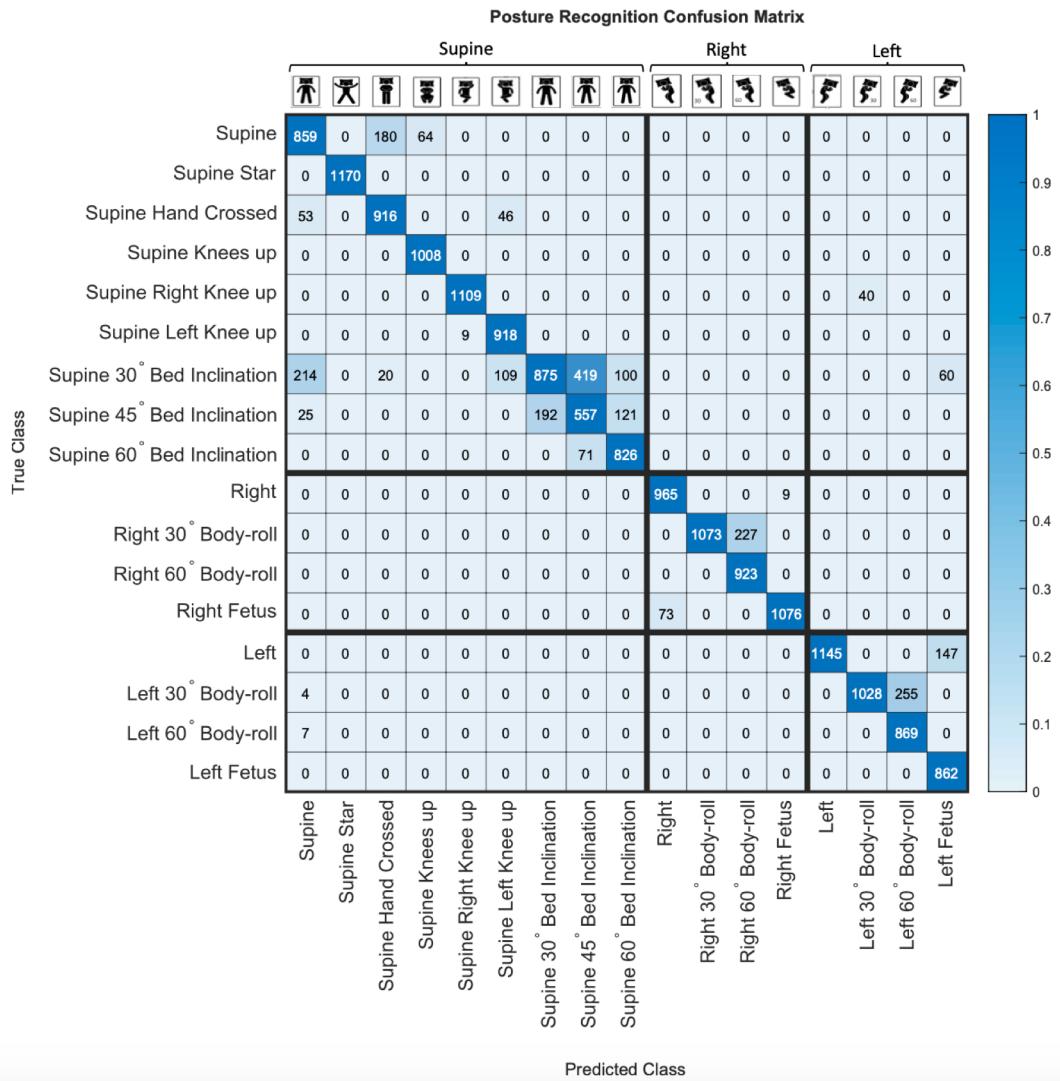
# In [Pouyan17]

- Results:

Posture	Predicted/Actual	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13
Supine	Recall	100	90	90	88.8	80	45.4	90	100	90	90	90	80	81.8
	Specificity	100	100	98.3	99.1	99.1	97.6	96.0	96.7	100	99.1	99.1	99.1	100
	Precision	100	100	81.8	88.8	88.8	62.5	64.2	71.4	100	90	80	88.8	100
	Accuracy									85.5				
Right Side	Recall	70	60	80	70	90	90	70	90	100	100	70	60	90
	Specificity	100	98.4	98.3	97.6	95.2	97.5	99.1	100	97.5	98.3	97.6	100	100
	Precision	100	75	72.7	70	60	75	87.5	100	76.9	8.3	70	100	100
	Accuracy									80.4				
Left Side	Recall	33.3	70	100	100	90	81.8	70	100	100	90	80	80	70
	Specificity	100	96.0	96.7	97.5	99.1	98.3	97.5	99.1	98.3	99.1	98.3	98.3	100
	Precision	100	58.3	71.4	76.9	90	75	87.5	90.9	83.3	90	80	80	100
	Accuracy									82.3				
Participants' Details	Age	19	23	23	24	24	26	27	27	30	30	30	33	34
	Height (cm)	175	183	183	177	172	169	179	186	174	174	176	170	174
	Weight (kg)	87	85	100	70	66	83	96	63	74	79	91	78	74

# Other reference [Davoodnia19]

- Subject identification and posture recognition using the same data of [Pouyan17]
- Architecture: CNN
- Automatic feature extraction



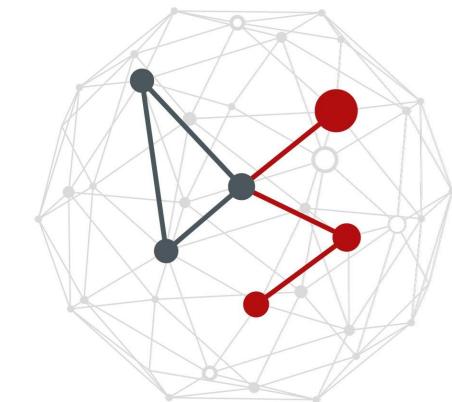
[Davoodnia19] V. Davoodnia and A. Etemad, Identity and Posture Recognition in Smart Beds with Deep Multitask Learning, in Proceedings of the IEEE International Conference on Systems, Man and Cybernetics (SMC), Bari, Italy, 2019.

# Project proposal

- **Classification tasks**
  - Subject identification
  - Posture recognition
  - Joint subject identification and posture recognition
- **Datasets**
  - use one of both the available datasets
  - different mattresses
- **Features**
  - manual feature extraction or raw data
- **Architecture**
  - CNN, RNN, combinations...

# PROJECT 5

---



# Project no. 5

## “environmental sound classification”

### Reference papers

[Piczak15] K.J. Piczak, [ESC: Dataset for Environmental Sound Classification](#), in Proceedings of the 23rd ACM International Conference on Multimedia, Brisbane, Australia, 2015.

[Piczak15-1] K. J. Piczak, [Environmental sound classification with convolutional neural networks](#), in Proceedings of the IEEE 25th International Workshop on Machine Learning for Signal Processing (MLSP), Boston, MA, 2015.

### ECS-50 dataset (884 MB uncompressed)

<https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/YDEPUT>

- Annotated collection of 2000 short clips comprising 50 classes of various common sound events

# High level description of the dataset

- 5-second-long clips, 44.1 kHz, single channel
- Arranged into 5 uniformly sized cross-validation folds, ensuring that clips originating from the same initial source file are always contained in a single fold

dog - 5-231762-A-0.wav



# High level description of the dataset

- 50 classes in the dataset

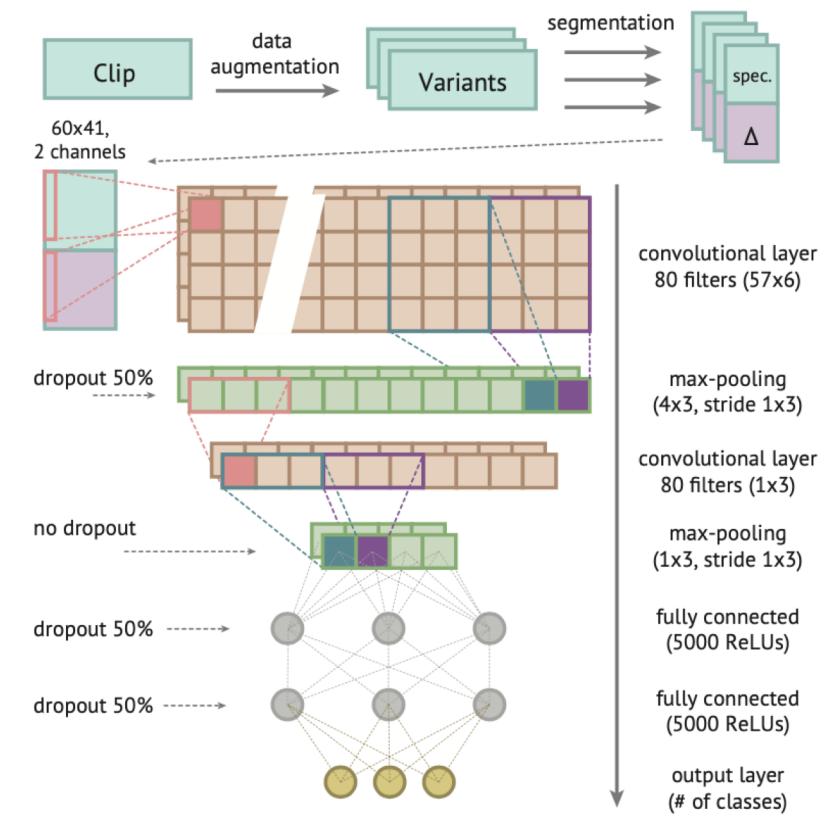
Animals	Natural soundscapes & water sounds	Human, non-speech sounds	Interior/domestic sounds	Exterior/urban noises
Dog	Rain	Crying baby	Door knock	Helicopter
Rooster	Sea waves	Sneezing	Mouse click	Chainsaw
Pig	Crackling fire	Clapping	Keyboard typing	Siren
Cow	Crickets	Breathing	Door, wood creaks	Car horn
Frog	Chirping birds	Coughing	Can opening	Engine
Cat	Water drops	Footsteps	Washing machine	Train
Hen	Wind	Laughing	Vacuum cleaner	Church bells
Insects (flying)	Pouring water	Brushing teeth	Clock alarm	Airplane
Sheep	Toilet flush	Snoring	Clock tick	Fireworks
Crow	Thunderstorm	Drinking, sipping	Glass breaking	Hand saw

# High level description of the dataset

- **ESC-10:** selection of 10 classes from the bigger dataset
  - The differences between classes are much more pronounced, with limited ambiguity
  - Classes: *sneezing, dog barking, clock ticking, crying baby, crowing rooster, rain, sea waves, fire crackling, helicopter, chainsaw*
- [meta/esc50.csv](#) data description, the “esc10” column indicates if a given file belongs to the *ESC-10* subset
- [meta/esc50-human.xlsx](#) contains the human classification accuracy

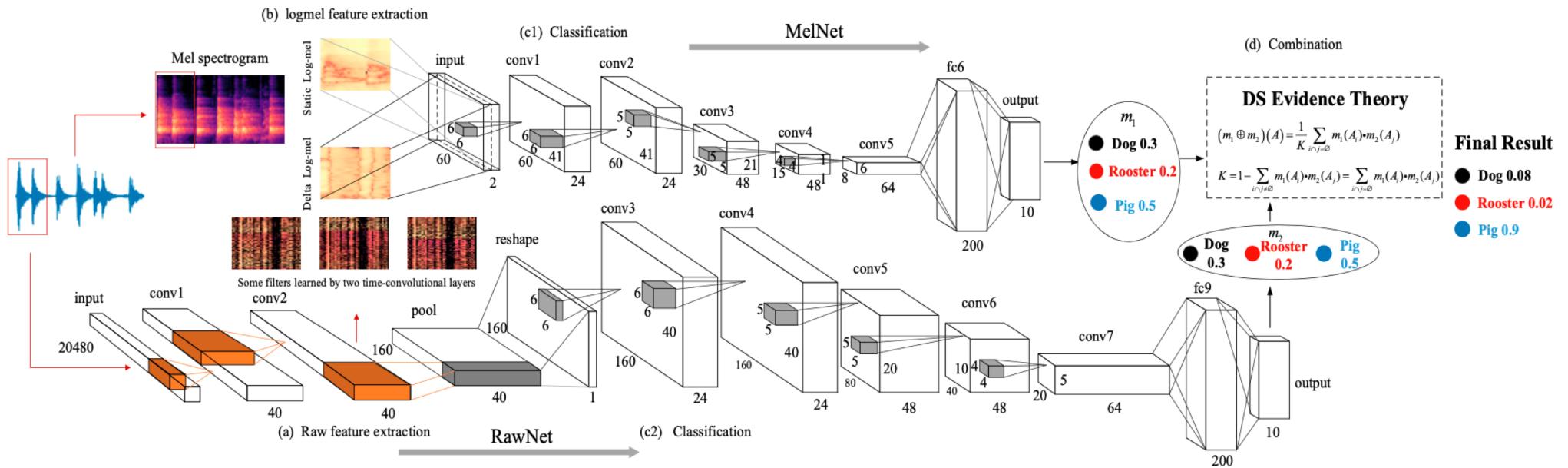
# Approach in [Piczak15-1]

- **Data augmentation:** apply random time delays to the original recordings
- **Feature extraction:** log-scaled mel-spectrograms with 60 mel-bands
  - resampled to 22,050 Hz
  - windows size 1024
  - hop length 512
- **Learning architecture:** CNN



# Other reference [Li18]

- Combines mel-spectrogram features and raw audio waveform



[Li18] S. Li, Y. Yao, J. Hu, G. Liu, X. Yao and J. Hu, *An Ensemble Stacked Convolutional Neural Network Model for Environmental Event Sound Recognition*, Applied Science, vol. 8, no. 1152, July 2018.

# Useful links

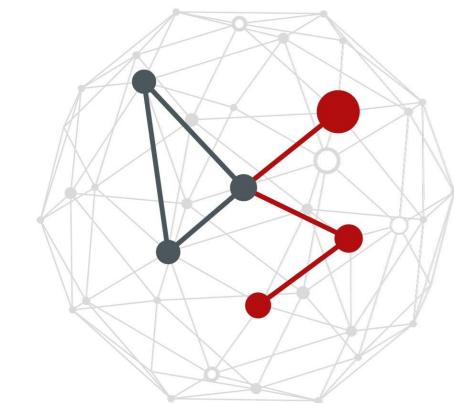
- Dataset GitHub repository  
<https://github.com/karolpiczak/ESC-50>
- Some useful functions  
<https://nbviewer.jupyter.org/github/karoldvl/paper-2015-esc-dataset/blob/master/Notebook/ESC-Dataset-for-Environmental-Sound-Classification.ipynb>

# Project proposal

- **Classification tasks**
  - on the entire ESC-50 dataset
  - on the restricted ESC-10 dataset
  - on each of the 5 groups of sounds:
    - animals
    - natural soundscapes & water sounds
    - human, non-speech sounds
    - interior/domestic sounds
    - exterior/urban noises
- **Features:** try with different approaches: mel-spectrogram, other manual-extracted features, raw data, combinations
- **Architecture**
  - different possibilities: CNN, RNN, ...

# PROJECT 6

---



# Project no. 6

## “lymphoma subtype classification”

### Reference paper

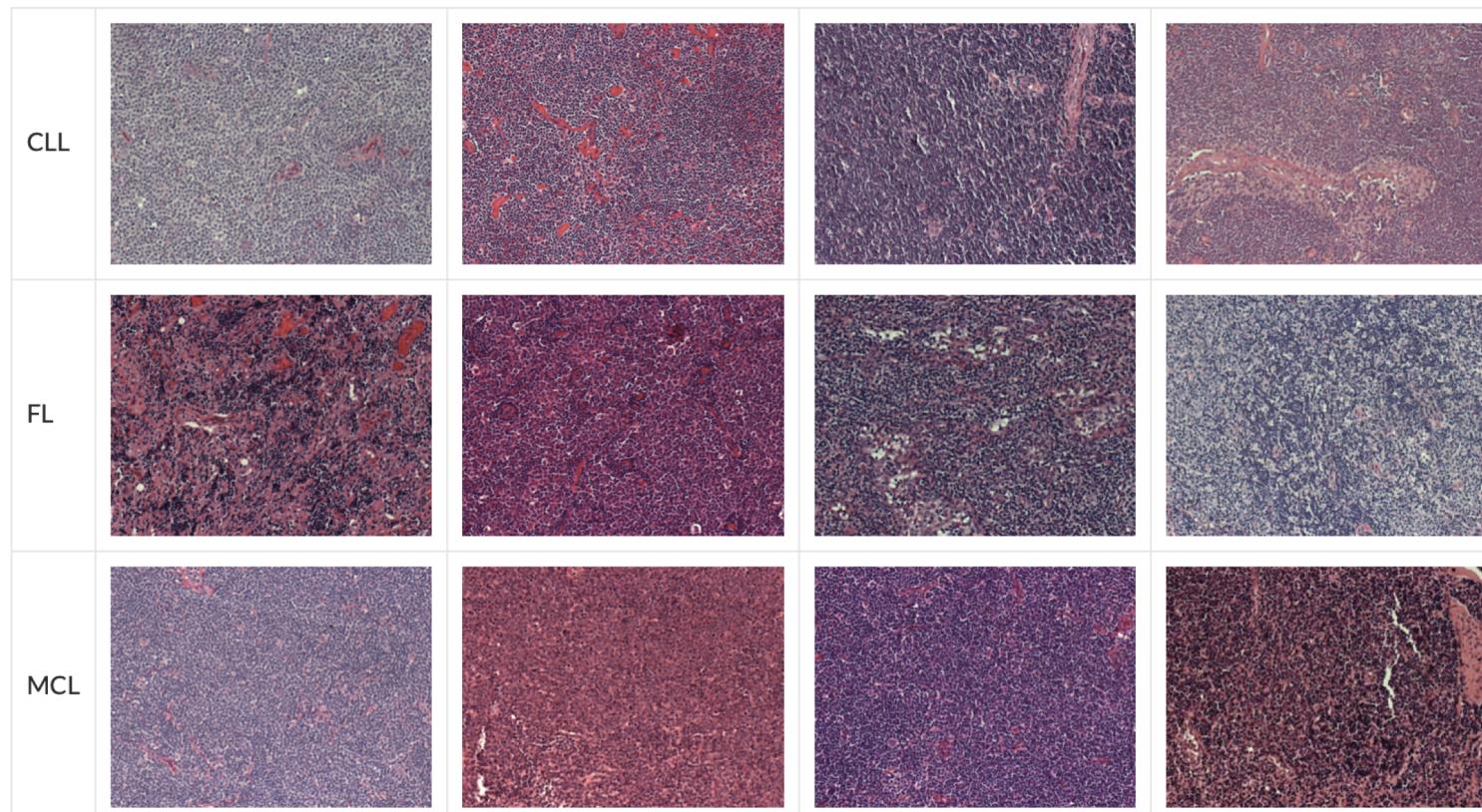
[Orlov10] N. V. Orlov, W. W. Chen, D. M. Eckley, T. J. Macura, L. Shamir, E. S. Jaffe, and I. G. Goldberg, [Automatic classification of lymphoma images with transform-based global features](#), IEEE transactions on information technology in biomedicine, vol. 14, no. 4, pp. 1003–1013, 2010.

### Dataset (1.62 GB uncompressed)

<https://ome.grc.nia.nih.gov/iicbu2008/lymphoma/index.html>

# High level description of the dataset

- 374 images of size 1388 x 1040:
  - 113 for the Chronic Lymphocytic Leukemia (CLL) class
  - 139 for the Follicular Lymphoma (FL) class
  - 122 for the Mantle Cell Lymphoma (MCL) class



# Approach in [Orlov10]

- Feature fusion: two-stage approach
  - compute spectral planes: simple (Fourier, Chebyshev, and wavelets) and compound transforms (Chebyshev of Fourier and wavelets of Fourier)
  - compute features for each pixel plane (raw data and spectral planes) as in [Orlov08]
- Classification
  - weighted neighbor distance (WND)
  - naïve Bayes network (BBN)
  - radial basis functions (RBF)
- Several color spaces were used: RGB, gray, Lab, H&E

[Orlov08] N. Orlov, L. Shamir, T. Macura, J. Johnston, D. M. Eckley, I. G. Goldberg, WND-CHARM: Multi-purpose image classification using compound image transforms, Pattern Recognition Letters, vol. 29, pp. 1684–1693, 2008.

# Other reference [Janowczyk16]

[Janowczyk16] A. Janowczyk and A. Madabhushi, Deep learning for digital pathology image analysis: A comprehensive tutorial with selected use cases, Journal of Pathology Informatics, vol. 7, no. 1, pp. 29, July 2016.

**'Lymphoma Subtype Classification Use Case' p.14**

- Author blog: <http://www.andrewjanowczyk.com/use-case-7-lymphoma-sub-type-classification/>

# Approach in [Janowczyk16]

- Approach
  - split the images in sub-patches
  - collect classification output for each patch
  - **winner-take-all**: the class with the highest number of votes became the designated class for the entire image
- Architecture
  - AlexNet ([\[Krizhevsky17\]](#))
- Accuracy: 96.58%

[\[Krizhevsky17\]](#) A. Krizhevsky, I. Sutskever, and G. E. Hinton, [ImageNet classification with deep convolutional neural networks](#), Communications of the ACM, vol. 60, no. 6, pp. 84–90, May 2017.

# Other reference [Tambe19]

- Inception V3 network ([Szegedy16])
  - several branches used to determine the appropriate type of convolution to be made at each layer

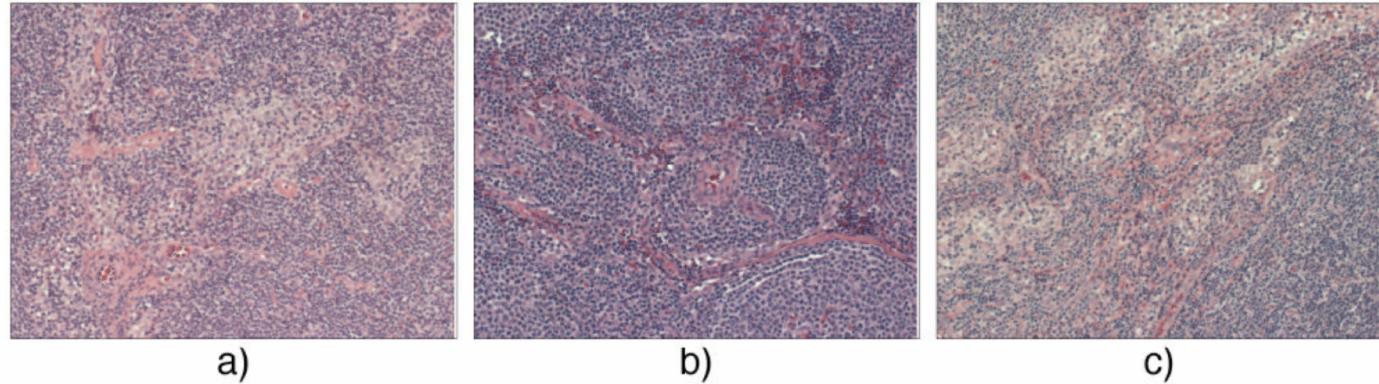
[Tambe19] R. Tambe, S. Mahajan, U. Shah, M. Agrawal, and B. Garware, Towards Designing an Automated Classification of Lymphoma subtypes using Deep Neural Networks, in Proceedings of the ACM India Joint International Conference on Data Science and Management of Data (CoDS-COMAD '19), Association for Computing Machinery, New York, NY, USA, 2019.

[Szegedy16] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens and Z. Wojna, Rethinking the Inception Architecture for Computer Vision, in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016.

# Project proposal

- Classification task

- CLL
- FL
- MCL



- Different color spaces can be used

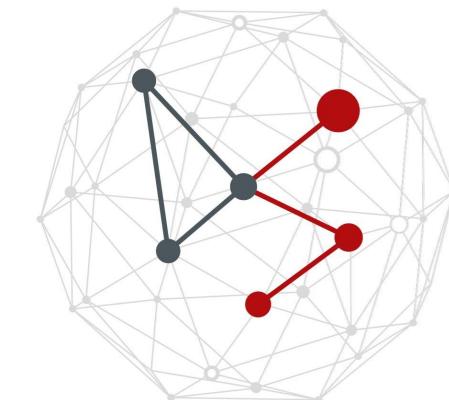
- RGB
- gray scale
- ...

- Architecture

- CNN, RNN, ...

# PROJECT 7

---



# Project no. 7

## “ECG beat classification”

### Reference papers

[Chen19] G. Chen, Z. Hong, Y. Guo, C. Pang, A cascaded classifier for multi-lead ECG based on feature fusion, Computer Methods and Programs in Biomedicine, vol. 178, pp. 135-143, September 2019.

### Dataset (794.5 MB uncompressed)

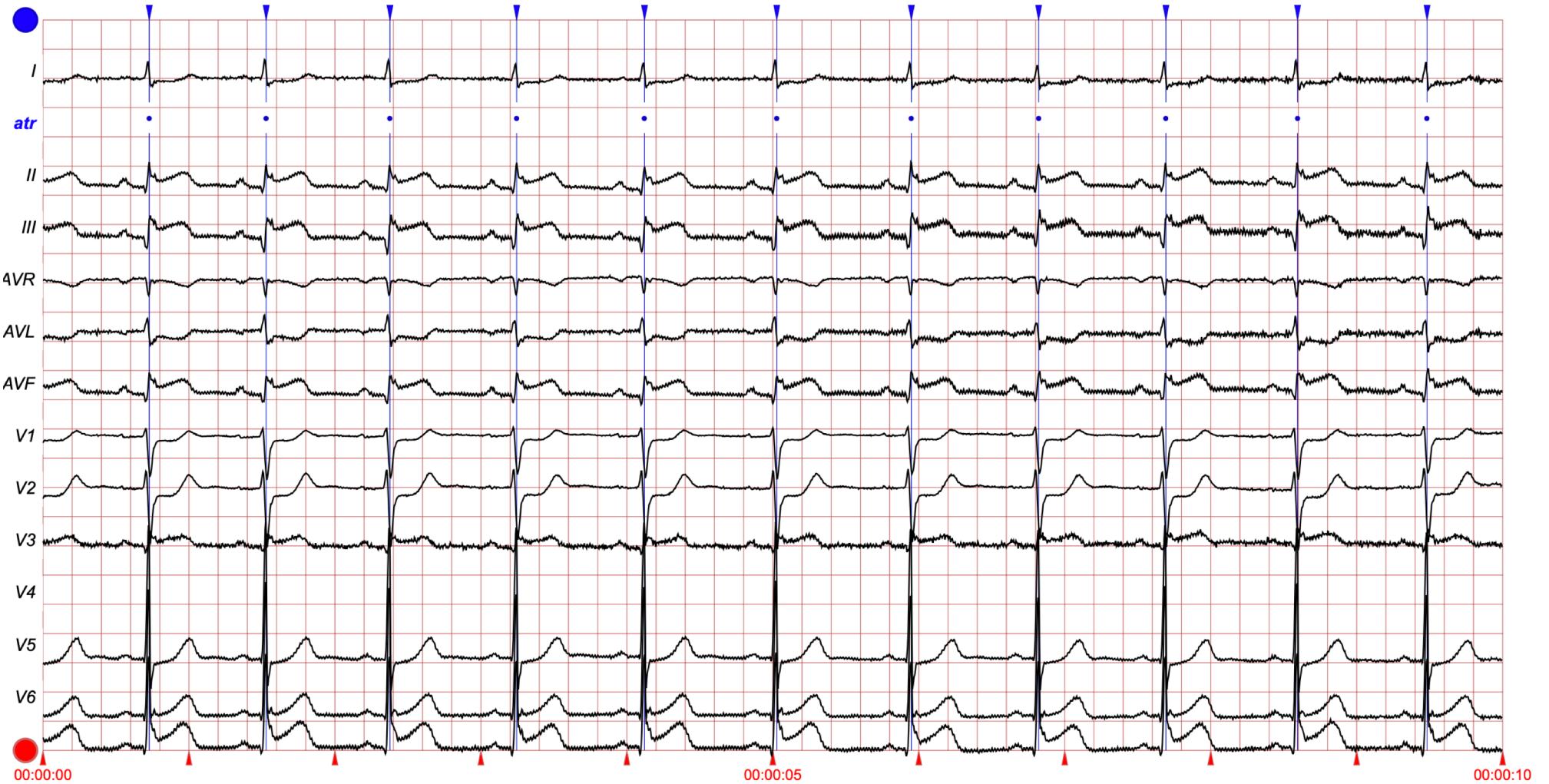
<https://physionet.org/content/incartdb/1.0.0/>

# High level description of the dataset

- 75 annotated recordings extracted from 32 Holter records
- each record
  - is 30 minutes long
  - contains 12 standard leads, each sampled at 257 Hz
- Total: over 175,000 beat annotations

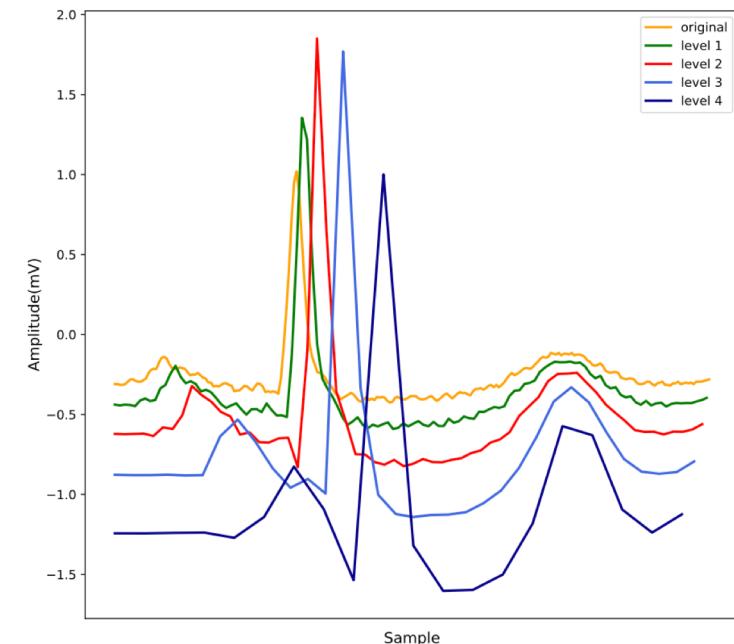
<b>Diagnosis</b>	<b>Patients</b>
Acute MI	2
Transient ischemic attack (angina pectoris)	5
Prior MI	4
Coronary artery disease with hypertension	7 (4 with ECGs consistent with left ventricular hypertrophy)
Sinus node dysfunction	1
Supraventricular ectopy	18
Atrial fibrillation or SVTA	3 (2 with paroxysmal AF)
WPW	2
AV block	1
Bundle branch block	3

# High level description of the dataset



# Approach in [Chen19]

- Classification task: 4 classes
  - normal beat (N)
  - atrial premature beat (A)
  - premature ventricular contraction (V)
  - right bundle branch block beat (R)
- Uses all the 12 leads
- 10 signal features (e.g., DWT)
- Cascade classification: random forest & MLP



Beat types	Precision	Recall	F1-score	Beat number
N	0.998	1.000	0.999	37,555
A	0.985	0.858	0.917	471
V	0.998	0.995	0.997	5033
R	0.999	0.997	0.998	772
<b>Avg</b>	<b>0.998</b>	<b>0.998</b>	<b>0.998</b>	<b>43,831</b>

# Other approaches

- See bibliography in [Chen19]

**Table 8**

Summary of recent ECG classification methods on INCART Database.

Literature	Method	Classes	Accuracy
M Llamedo et al. [31]	DWT,LDC	3	91
N Jannah et al. [42]	SIMCA,MSVM with PCA	4	76.83(SIMCA), 98.33( MSVM)
M Llamedo et al. [33]	PCA and WT	3	93
A Solosenko et al. [43]	DNN	3	98.6
<b>Proposed</b>	feature fusion, cascaded classifier	4	<b>99.8</b>

# Other reference [Yang19]

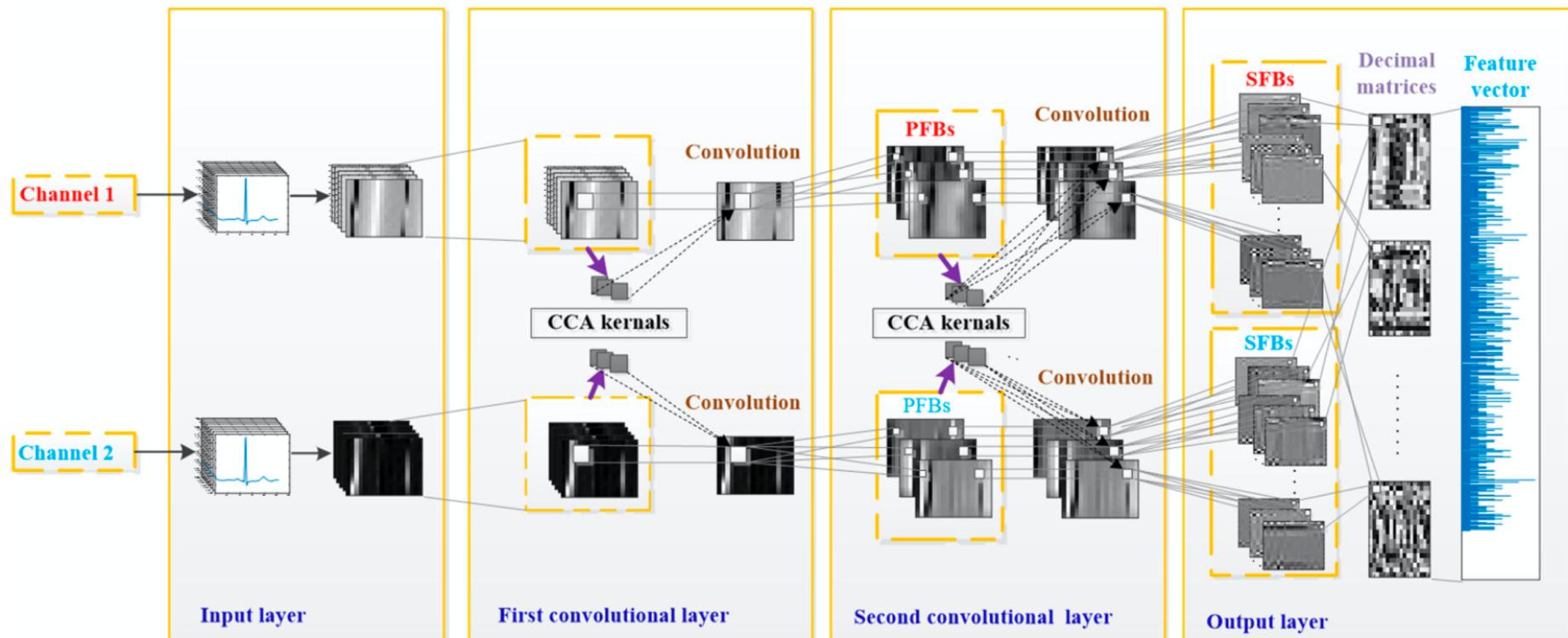
- Classification task: 7 classes
- Uses only three leads: II, V1 and V5

Types	Quantity	
N	Normal beat	500
V	Premature ventricular	500
A	Atrial premature beat	200
F	Fusion of ventricular contraction	200
n	Supraventricular espace beat	30
R	Right bundle branch block beat	200
j	Nodal (junction) escape beat	90
Total		1720

[Yang19] W. Yang, Y. Si, D. Wang and G. Zhang, [A Novel Approach for Multi-Lead ECG Classification Using DL-CCANet and TL-CCANet](#), Sensors, vol. 19, no. 14, pp. 3214, 2019.

# Approach in [Yang19]

- Feature extraction from dual-lead and three-lead ECG through DL-CCANet and TL-CCANet ([\[Yang17\]](#))
- Classification: SVM



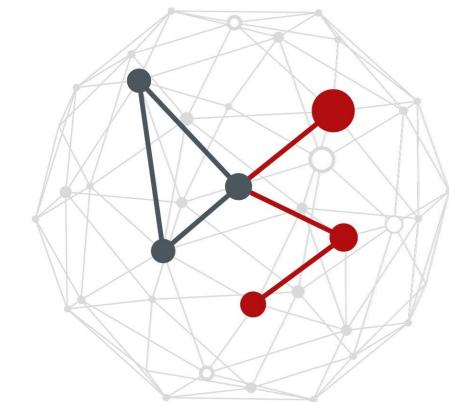
[\[Yang17\]](#) X. Yang, W. Liu, D. Tao and J. Cheng, [Canonical correlation analysis networks for two-view image recognition](#), *Information Sciences*, vol. 385–386, pp. 338-352, 2017.

# Project proposal

- Solve the classification task
  - with 3, 4 or more classes
- Number of leads
  - use all the leads available or apply a lead-selection strategy to retain only the most representative ones
- Features
  - raw data or manual feature extraction
- Architecture
  - different possibilities: CNN, RNN, ...

# EXAM DATES

---



# Exam dates - summer session 2020

Exams will be held online via Zoom Meeting

- Instructions and Zoom meeting link will be provided
- Subscription to the exam will be via **UNIWEB**
- Groups are required to register via Moodle questionnaire
  - in the HDA course Moodle site
  - to specify group composition (max. 2 people per group), selected topic for the final project, project title

**Exam dates** for the present academic year

- Monday **June 29, 2020** - Friday **July 3, 2020**
- Monday **July 20, 2020** - Friday **July 24, 2020**

Subscribe to the exams via **UNIWEB**

<https://uniweb.unipd.it/>

# HDA COURSE PROJECTS

Michele Rossi

[rossi@dei.unipd.it](mailto:rossi@dei.unipd.it)

Francesca Meneghello

[meneghello@dei.unipd.it](mailto:meneghello@dei.unipd.it)



DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE



DIPARTIMENTO  
**MATEMATICA**

1222-2022  
**800** ANNI



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

