Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

# Study of a deep learning model for temporal sleep stage classification

## Bachelor Thesis Defense

Maëlys Solal[1], Dr Alexandre Gramfort[2] and Dr Olivier Pallanca[3]

[1]École polytechnique

[2]Inria Paris Saclay, Parietal team

[3]Laboratoire d'Informatique de l'École polytechnique (LIX), DaSciM team

January - March 2021

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

# Table of Contents

Sleep
Staging

Maëlys
Solal

Context
Sleep stages
Motivations
for sleep
stage
classification

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

# Table of Contents

Sleep
Staging

Maëlys
Solal

Context

Sleep stages
Motivations
for sleep
stage
classification

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

# Sleep stages

- Sleep cycle = mix between REM sleep and NREM sleep, lasts 90 to 120 minutes and occurs 4-6 times per night

- 5 main sleep stages: Wake (W), Rapid Eye Movement (REM) $\simeq$ paradoxical sleep, Non REM1 (N1) $\simeq$ light sleep. Non REM2 (N2) $\simeq$ deeper sleep, Non REM3/4 (N3/4) $\simeq$ deep sleep.

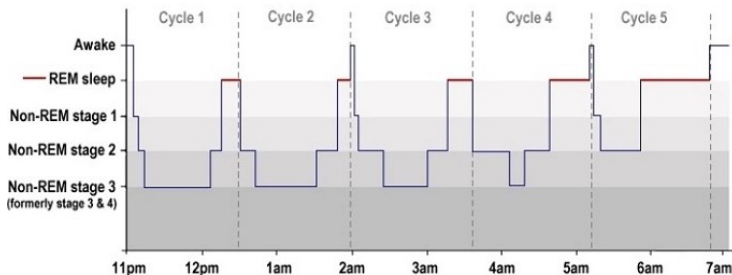- Have distinct characteristics as seen in the analysis of brain waves



Figure: A typical hypnogram showing sleep stages and cycles, by Luke Mastin
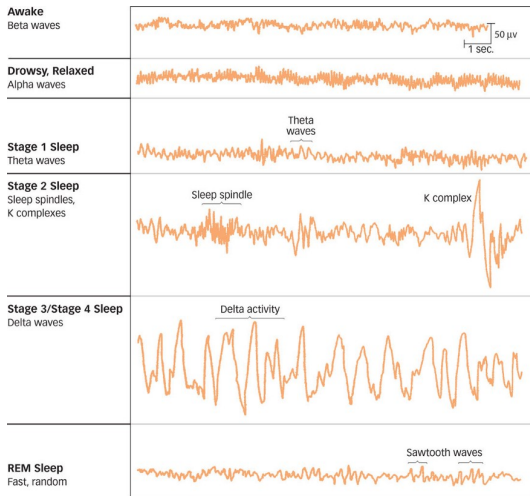
Sleep
Staging

Maëlys
Solal

Context
Sleep stages
Motivations
for sleep
stage
classification
Objectives
Methods
Results
and
Discussion
Conclusion
Annex

# Sleep stages

Polysomnography = sleep study

- Biophysical changes during sleep
- Includes EEG (brain's electrical activity), EOG (eyes), EMG (muscles), and ECG (heart)

Sleep stage classification = sleep scoring

- Visual investigation of the PSG, labelling 30s time segments with sleep stages
- According to a precise set of rules
- Usually done manually by sleep experts (scorers)



Figure: Brain waves during the different sleep stages, from MacMillan Learning

Sleep
Staging

Maëlys
Solal

Context

Sleep stages

Motivations
for sleep
stage
classification

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

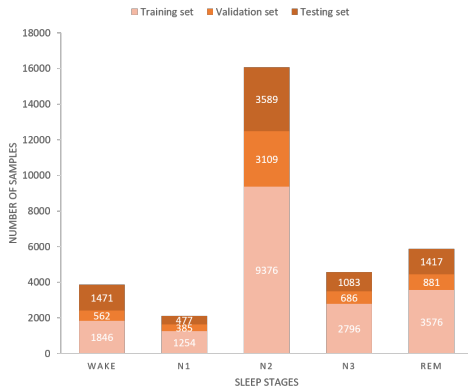# Motivations for sleep stage classification



Figure: Sleep stages imbalance in MASS dataset

From a clinical point of view:

- Used as a preliminary examination for clinical diagnosis of sleeping disorders
- Manual scoring is tedious

From a statistical learning point of view:

- Multiclass classification with imbalanced classes as shown in Figure 3
- Domain adaptation (i.e. differences of raw data between datasets)
- In general, quite noisy data (especially clinical data)

Sleep
Staging

Maëlys
Solal

Context

Objectives

Study of a
deep learning
model for
sleep scoring

Data

Organising
the datasets:
BIDS
standard

Methods

Results
and
Discussion

Conclusion

Annex

# Table of Contents

Sleep
Staging

Maëlys
Solal

Context
Objectives
Study of a
deep learning
model for
sleep scoring
Data
Organising
the datasets:
BIDS
standard
Methods
Results
and
Discussion
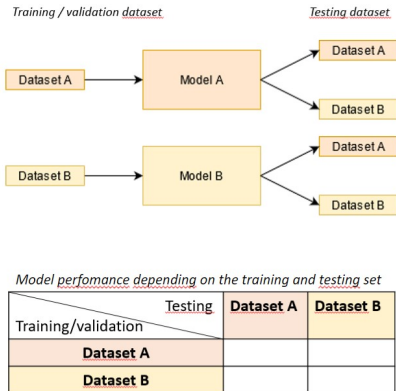Conclusion
Annex

# Study of a deep learning model for sleep scoring

- Our model of study: deep
  convolutional neural network,
  performs temporal sleep stage
  classification using multivariate and
  multimodal time series, by Chambon
  et al. in 2018[a]
- Study the transferability of the model
  i.e. performance depending on the
  training/validation sets and testing set

---

[a]Stanislas Chambon et al. "A deep learning
architecture for temporal sleep stage classification using
multivariate and multimodal time series". In: (2018).

*Training / validation dataset*          *Testing dataset*



*Model perfomance depending on the training and testing set*

| Testing Training/validation | Dataset A | Dataset B |
|---|---|---|
| **Dataset A** | | |
| **Dataset B** | | |

Figure: Schematic representation of our
experiment

Sleep
Staging

Maëlys
Solal

Context

Objectives

Study of a
deep learning
model for
sleep scoring

Data

Organising
the datasets:
BIDS
standard

Methods

Results
and
Discussion

Conclusion

Annex

# Data

- Datasets
  - Montreal Archive of Sleep Studies (MASS)[1]
  - SleepPhysionet[2]
  - Clinical dataset
- Varied number and types of EEG/EMG/EOG channels
  - Figure 5 shows the variety of EEG channels
  - For comparing the performance of the model across datasets: should have similar EEG/EMG/EOG channels
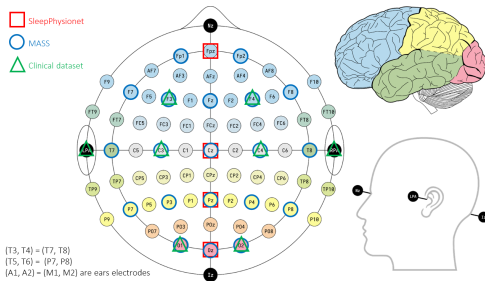


(T3, T4) = (T7, T8)
(T5, T6) = (P7, P8)
(A1, A2) = (M1, M2) are ears electrodes

SleepPhysionet
MASS
Clinical dataset

Figure: EEG electrodes positioning used in each dataset

---

[1] Christian O'Reilly et al. "Montreal Archive of Sleep Studies: an open-access resource for instrument benchmarking and exploratory research". In: (2014).

[2] B. Kemp et al. "Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the EEG". In: (2000), Ary L. Goldberger et al. "PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals". In: (2000).

Sleep
Staging

Maëlys
Solal

Context

Objectives

Study of a
deep learning
model for
sleep scoring

Data

Organising
the datasets:
BIDS
standard

Methods

Results
and
Discussion

Conclusion

Annex

# Organising the datasets: BIDS standard



Figure: BIDS standard

- Neuroimaging data is often complicated to arrange
  - Comes from various experiments / medical examinations
  - Outputs multiple files for a single patient
  - Little consensus about how to organise files
- BIDS standard proposes a way to organise this type of data
  - BIDS = Brain Imaging Data Structure[3]
  - subject > session > type / experiment
  - Compatible with existing software
  - Captures metadata

[3] Krzysztof J. Gorgolewski et al. "The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments". In: (June 2016). DOI: 10.1038/sdata.2016.44. URL: https://www.hal.inserm.fr/inserm-01345616, Stefan Appelhoff et al. "MNE-BIDS: Organizing electrophysiological data into the BIDS format and facilitating their analysis". In: (2019). URL: https://doi.org/10.21105/joss.01896, Cyril Pernet et al. "BIDS-EEG: an extension to the Brain Imaging Data Structure (BIDS) Specification for electroencephalography". In: (Jan. 2018).

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Main steps,
code
structure
Description
of the model
Pre-
processing,
training

Results
and
Discussion

Conclusion

Annex

# Table of Contents

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Main steps,
code
structure

Description
of the model

Pre-
processing,
training

Results
and
Discussion

Conclusion

Annex

# Main steps, code structure

**Main steps for running our experiment:**

1. Loading the datasets
2. Preprocessing the raw signals and extracting the 30s windows from events
3. Splitting the dataset into training, validations and testing sets
4. Creating / loading our model, training and testing
5. Visualising the results

**Code structure**, inspired by the `braindecode`[a] library

- `datasets` to load the datasets that were previously converted to BIDS
- `datautil` to take care of splitting the datasets
- `models` to load and save our models
- `visualisation` for visualising results and the 30s windows

---

[a]Robin Tibor Schirrmeister et al. "Deep learning with convolutional neural networks for EEG decoding and visualization". In: (). URL: http://dx.doi.org/10.1002/hbm.23730.

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Main steps,
code
structure

Description
of the model

Pre-
processing,
training

Results
and
Discussion

Conclusion

Annex

# Description of the model

- General feature extractor, denoted by $Z : \mathbb{R}^{C \times T} \to \mathbb{R}^D$, where $C$ is the number of input channels, $T$ is the number of time steps and $D$ is the size of the estimated feature space
- **Linear spatial filtering**: to estimate virtual channels
- **Convolutional layers**: to capture spectral features
- **Separate pipelines**: to handle several modalities at the same time
- Performs **temporal** sleep staging, that is, takes into account the temporal context of the sample of interest to predict a label



Figure: Network general architecture

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Main steps,
code
structure

Description
of the model

Pre-
processing,
training

Results
and
Discussion

Conclusion

Annex

# Description of the model



Figure: Schematic representation of the sleep staging model's architecture



Figure: Schematic representation of the time distributed multivariate network

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods
Main steps,
code
structure
Description
of the model
Pre-
processing,
training

Results
and
Discussion

Conclusion

Annex

# Pre-processing, training

**Preprocessing steps**, using `mne-python`[a], same steps as Chambon et al.

- Low-pass filtered at 30Hz: to mitigate the impact of higher frequency noise
- Downsampled to 100Hz: SleepPhysionet was sampled at 100Hz
- Convert signals from V to $\mu$V: very small amplitude brain waves
- Cropped 30 minutes of wake events at the beginning and the end of the night
- Divided our signal in 30s samples (windows) corresponding to one specific sleep stage
- Standardised the windows (zero mean, unit variance): cope for the varying recording conditions

**Training specification**

- Implemented with PyTorch
- 60 subjects for each experiment, splitted using stratified cross-validation: roughly 60% of events in training set, 20% in validation set and 20% in testing set
- Weights initialised with a normal distribution ($\mu = 0$, $\sigma = 0.1$)
- Loss function: categorical cross entropy
- Optimizer: AdamOptimizer
- Minimisation: stochastic gradient descent, learning rate $lr = 5 \times 10^{-4}$, batch size 8, 10 training epochs

[a]Alexandre Gramfort et al. "MEG and EEG data analysis with MNE-Python". In: *Frontiers in Neuroscience* (2013). URL: https://www.frontiersin.org/article/10.3389/fnins.2013.00267.

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Typical
results for
our
experiment

Results of
our main
experiment
and
Discussion

Limitations
and Further
Works

Conclusion

Annex

# Table of Contents

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Typical
results for
our
experiment

Results of
our main
experiment
and
Discussion

Limitations
and Further
Works

Conclusion

Annex

# Typical results for our experiment



Figure: Confusion matrix for MASS dataset



Figure: Classification report for MASS dataset

# Results of our main experiment and Discussion

| Training set | Testing set MASS | SleepPhysionet | | Training set | Testing set MASS | Clinical |
|---|---|---|---|---|---|---|
| MASS | 0.802 | 0.487 | | MASS | 0.817 | 0.390 |
| SleepPhysionet | 0.560 | 0.634 | | Clinical | 0.530 | 0.635 |

Table: Balanced accuracy for main experiment

- Clinical-Clinical and SleepPhysionet-SleepPhysionet lower than expected
  - Pre-processing steps
  - Model initially benchmarked using MASS, model biased towards MASS?
  - Hyperparameters are not optimal
- Model transfers better from SleepPhysionet/Clinical to MASS than from MASS to SleepPhysionet/Clinical
  - Specificity in the MASS dataset?
- Domain adaptation remains one of the big challenges of sleep scoring algorithms
- Increasing number of EEG channels and adding EOG and EMG data greatly improves performance

# Limitations and Further Works

- SleepPhysionet scored according to the Rechtschaffen and Kales guidelines, MASS and Clinical scored according to the AASM guidelines
  - N4 sleep stage
  - Transition rules
- Tune hyperparameters and improve the model's architecture to get better performance on Clinical and SleepPhysionet
- Core differences in terms of data
  - Clinical datasets are in general more difficult to work with

Sleep
Staging

Maëlys
Solal

Context
Objectives
Methods
Results
and
Discussion
Conclusion
Perspectives
Acknowledge-
ments
Annex

# Table of Contents

# Perspectives

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Conclusion

Perspectives

Acknowledge-
ments

Annex

# Acknowledgements

Annex



Table: Table of confusion matrices, comparing the datasets MASS and SleepPhysionet

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

Table: Table of classification reports, comparing the datasets MASS and SleepPhysionet

Sleep
Staging

Maëlys
Solal

Context

Objectives

Methods

Results
and
Discussion

Conclusion

Annex

Table: Table of confusion matrices, comparing the datasets MASS and Clinical

Table: Table of classification reports, comparing the datasets MASS and Clinical