

# Summary

---

The use of surface electromyography (sEMG) as a measure of human physiology and behaviour has grown recently, supported by developments in deep learning and wearable computing. Here, we present *EMGFlow*, an open-source Python package for preprocessing and extracting features from sEMG signals. *EMGFlow* has been designed to facilitate the analysis of large datasets through batch processing of signal files, a common requirement in machine learning. The package extracts an extensive set of features from both time and frequency domains. Regular expression matching provides additional flexibility in mapping files for selective preprocessing and extraction. The use of Pandas DataFrame throughout allows users to mix and match elements of the processing pipeline, supporting interoperability with other packages. An interactive dashboard supports human decision processes through a visual comparison of signals at each stage of preprocessing. *EMGFlow* is released under the GNU General Public License (v3.0) and can be installed from PyPI. Source code, documentation, and examples are accessible on GitHub (<https://github.com/Willson/EMGFlow-Python-Package>).

## Statement of Need

---

Although several packages exist for processing physiological and neurological signals, support for sEMG has remained limited. Many packages lack a comprehensive set of features that can be extracted from sEMG data, leaving researchers to use a patchwork of tools. Other packages are orientated around event detection in individual recordings and use a GUI-based workflow that requires more manual intervention. While this design works well for processing unedited continuous recordings of a single participant, it complicates the extraction of features from large datasets common to machine learning [abadi\_decaf\_2015; chen\_emotion\_2022; koelstra\_deap\_2012; schmidt\_introducing\_2018; sharma\_dataset\_2019; zhang\_biovid\_2016]. *EMGFlow*, a portmanteau of EMG and Workflow, fills this gap by providing a flexible pipeline for extracting a wide range of sEMG features, with a scalable design suited for large datasets.

## Comparison to Other Packages

---

Compared to other toolkits, *EMGFlow* extracts a comprehensive set of 32 statistical features from sEMG signals (Bota et al, 2024; Makowski et al, 2021; Sjak-Shie et al, 2022; Soleymani et al, 2017). An interactive dashboard visualizes batch processed files rather than individual recordings, allowing the operator to efficiently view the effects of preprocessing stages across all files. Adjustable filter settings and smoothing functions support cleaning of data collected in North America or internationally (50 vs 60 HZ mains AC), a subtle difference overlooked in some packages.

## Features

---

### Processing Pipeline

Extracting features from large datasets is a common task in machine learning and quantitative domains. *EMGFlow* supports this need through batch-processing, allowing users to either semi- or fully automate the treatment of sEMG recordings. To demonstrate, we use data from PeakAffectDS (Greene, Livingstone, &

Szymanski, 2022), a collection of physiological signals that includes two channels of facial sEMG, labelled Zyg and Cor, capturing Zygomaticus major and Corrugator supercilii muscle activity respectively. We begin by defining the path to the directory containing our raw, uncleaned files stored in plaintext (.csv) format. We then apply a notch filter to remove the AC mains noise introduced by the recording system's power source, a common initial step in preprocessing raw sEMG signals.

```
import EMGFlow

# Paths for sEMG files
raw_path = 'Data/01_Raw'
notch_path = 'Data/02_Notch'

# Sampling rate
sampling_rate = 2000

# Columns containing sEMG
cols = ['EMG_zyg', 'EMG_cor']

# Notch filter parameters
notch_vals = [(50,5)]

# Apply notch filter to raw sEMG files
EMGFlow.NotchFilterSignals(raw_path, notch_path, sampling_rate, notch_vals, cols)
```

Additional arguments allow users to customize which files are selected and how they are processed. Filtering functions accept an optional regex argument, allowing users to apply filters to specific files. Most functions use common sense defaults, which can be modified task-wide or for select cases. For example, in North America, mains electricity is nominally supplied at 120 VAC 60 Hz, while other countries may supply power at 200-240 VAC 50Hz. This variation in frequency requires different notch filter settings depending on where the data were recorded. *EMGFlow* accommodates this need by allowing the user to specify the frequency and quality factor of the applied filter. Extending our first example, we now apply an additional notch filter to a subset of files exhibiting noise at 150 Hz, the 3rd harmonic of the mains source.

```
# Paths for sEMG files
notch_s_path = 'Data/02_Notch_Special'

# Filter parameters for files that start with "08" or "11"
notch_sc = [(150,25)]
reg = '^(08|11)'

# Column to apply to
cols = ['EMG_zyg', 'EMG_cor']

# Apply notch filter to file subset
EMGFlow.NotchFilterSignals(notch_path, notch_s_path, sampling_rate, notch_vals_s,
cols, expression=reg_str, exp_copy=True)
```

## Visualization of Preprocessing Stages

The application of a bandpass filter is often the second stage in preprocessing sEMG signals, as it isolates the frequency spectrum of human muscle activity. Signals are commonly filtered to the 10-500 Hz range (Livingstone et al., 2016; McManus, De Vito, & Lowery, 2020; Sato et al., 2021; Tamietto et al., 2009), though precise filter corner frequencies vary by research domain and approach (Abadi et al, 2015). After filtering, data can be further smoothed to remove high-frequency noise and outliers in preparation for the extraction of temporal features. The default smoother is RMS, equal to the square root of the total power in the sEMG signal and commonly used to estimate signal amplitude (McManus, De Vito, & Lowery, 2020). Additional filter options are provided, including boxcar, Gaussian, and LOESS. *EMGFlow* provides an interactive Shiny dashboard to visualize the effects of preprocessing on sEMG signals. Preprocessing stages can be displayed simultaneously or shown individually with options for Notch, Bandpass, and Smoothing steps. Users can select the file for visualization using the Files dropdown box. The dashboard is generated from a list of file paths containing files at different stages of preprocessing. Here, our example shows how signals are further bandpass filtered and smoothed, with results visualized using the dashboard.

```
# Paths for sEMG files
band_path = 'Data/03_Bandpass'
smooth_path = 'Data/04_Smoothed'

# Filter and smoothing parameters
band_low = 20
band_high = 450
smooth_window = 50

# Apply bandpass smoothing filters
EMGFlow.BandpassFilterSignals(notch_s_path, band_path, sampling_rate, band_low,
band_high, cols)
EMGFlow.SmoothFilterSignals(band_path, smooth_path, sampling_rate, smooth_window,
cols)

# Paths for dashboard generation
in_paths = [smooth_path, band_path, notch_path]
labels = ['Smooth', 'Bandpass', 'Notch']

# Column to visualize, and units of measurement
col = 'EMG_zyg'
units = 'mV'

# Create plot
EMGFlow.GenPlotDash(in_paths, sampling_rate, col, units, labels)
```

![[figure1.png]] *EMGFlow's interactive dashboard visualizing effects of different preprocessing stages on batch processed files.*

## The nature of electromyographic recordings

To better understand the range of features extracted by *EMGFlow*, we begin with a review of surface electromyography as a recording instrument. Nearly all body movement occurs by muscle contraction. During

contraction, nerve impulses sent from motoneurons cause muscle fibers innervated by the axon to discharge, creating a motor unit action potential (McManus, De Vito, & Lowery, 2020; Stepp, 2012). The speed at which action potentials propagate down the fibre is called muscle fiber conduction velocity. Each motor unit firing results in a force twitch. The superposition of these twiches over time produces a sustained force that enables functional muscle activity, such as lifting or smiling (De Luca, 2008). Surface electromyography measures voltage difference across muscle fibers generated by action potentials, producing a voltage timeseries that quantifies muscle activity (Fridlund & Cacioppo, 1986). It is from this voltage timeseries that statistical features are extracted.

## Feature Extraction Routines

Following data preprocessing, the signal files are ready for feature extraction. *EMGFlow* extracts 32 features that capture information in both time and frequency domains. The set of 18 time-domain features capture standard statistical moments, including mean, variance, skew, and kurtosis, along with sEMG-specific measures. These include features such as Willison amplitude, an indicator of motor unit firing calculated as the number of times the sEMG amplitude exceeds a threshold, and log-detector, an estimate of the exerted muscle force (Tkach, Huang, & Kuiken, 2010). A set of 12 frequency-domain features are also extracted, providing information on the shape and distribution of the signal’s power spectrum. Measures such as median frequency (Phinyomark, Limsakul, & Phukpattaranont, 2009) provide insight into changes in muscle fibre conduction velocity and are used in the assessment of muscle fatigue (van Boxtel, Goudswaard, van der Molen & van den Bosch, 1983; Lindstrom, Kadefors, & Petersen, 1977; McManus, De Vito, & Lowery, 2020). Standard frequency measures include spectral centroid, flatness, entropy, and roll-off. One novel sEMG feature introduced here is Twitch Ratio, an adaptation of Alpha Ratio from speech signal analysis (Eyben et al., 2016). Twitch Ratio is defined as the ratio of energy contained in the upper versus lower power spectrum, with a threshold of 60 Hz to delineate slow- and fast-twitch muscles fibres (Hegedus et al., 2020). Here, we demonstrate feature extraction in *EMGFlow*. After specifying locations of processed files, features are summarized into a single CSV file, containing rows for each file analyzed, as shown below.

```
# Data paths
feature_path = 'Data/05_Feature'

# Extracts features
EMGFlow.AnalyzeSignals(band_path, smooth_path, feature_path, sampling_rate, cols)

# Load feature file
df = read_csv('Data/05_Feature/Features.csv')

# Print first few rows
df.head()
"""
      File_ID  EMG_zyg_Min  ...  EMG_cor_Spec_Rolloff  EMG_cor_Spec_Bandwidth
0  01-01-01.csv    0.000826  ...                0.040222             1424.933862
1  01-01-02.csv    0.000740  ...                0.019559             2651.987804
2  01-01-03.csv    0.000780  ...                0.065183             2021.345274
3  01-01-04.csv    0.000660  ...                0.087384             1755.834836
4  01-01-05.csv    0.000697  ...                0.057368             1174.562467

[5 rows x 61 columns]
"""
```

--

*An example of part of a feature file. The "File\_ID" column contains the names of the files extracted, and the additional columns take the format "[Column name]\_[Feature name]".*

# Community Guidelines

---

We welcome contributions to the project. These can be initiated through the project's issue tracker or via a pull request. Suggestions for feature enhancements, tips, as well as general questions and concerns, can also be expressed through direct interaction with contributors and developers.

# Declaration of Generative AI and AI-Assisted Technologies in the Writing Process

---

During the preparation of this work, the authors used GPT-4o to edit a final draft of the manuscript for flow, tone, and grammatical correctness. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

# Acknowledgements

---

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC), (#2023-03786), and from the Faculty of Science, Ontario Tech University.

# References

---