# Enabling naturalistic, long-duration and continual animal experimentation

March 18, 2025

# Contents

# 1 Vision

For over four years, at the Sainsbury Wellcome Centre and Gatsby Computational Neuroscience Unit, we have been developing the AEON platform, a set of hardware and software tools that support a new type of experimentation, where animals are allowed to express ethologically-relevant behaviors, in naturalistic environments, and in long-duration experiments, while their behavior and neural activity is monitored continuously for weeks to months. We have used this platform to characterize foraging behavior in both solitary and groups of mice (**?**) (Figure **??**).

Our US partner, the Allen Institute for Neural Dynamics, is using the AEON platform in continuous learning experiments, where mice freely explore odors continuously for days to weeks (**?**).

This is an unprecedented type of experimentation that . . .

Several groups around the world are performing this new type of experimentation .

We have built the AEON platform, and have used it to collect weeks-to months-long NaLoDuCo experimental data. We next propose to develop advanced machine learning methods and intelligent visualisations to extract meaning from this data (Aim 1).

A central aim of both the SWC/GCNU and AIND is to contribute to open science. We thus propose to create software infrastructure to openly disseminate NaLoDuCo recordings, visualisation and data analysis methods (Aim 2).

Over more than four years we have developed the AEON platform following high software engineering practices. It is an open source platform that anybody can use and modify (**?**). We want it to become the standard platform for the collection of NaLoDuCo experimental data. We are currently using AEON on two new NaLoDuCo experiments: (1) odor learning experiments, lasting for days to weeks, lead by Dr. Carl Schoonover at the AIND, and (2) foraging experiments in very large arenas (eight meters in diameter), lead by Prof. Tiago Branco at the SWC. We will extend and validate the functionality of the AEON platform by applying it to these and new NaLoDuCo experiments. A key functionality that we propose to add as part of this project is real-time machine learning, to allow to control AEON experiments with live inferences (Aim 3).
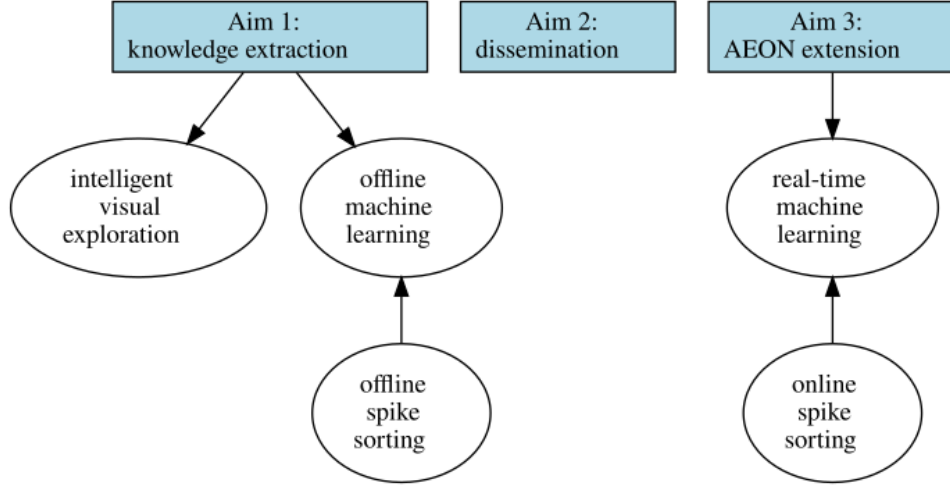
Figure 1: Proposal aims

## Aim 1: create infrastructure for open dissemination of NaLoDuCo experimental recordings

The dissemination of NaLoDuCo recordings is not trivial, as datasets generated by this new type of experimentation are enormous. For instance, the size of a dataset generated from a one week recording of behavioral and neural activity from a foraging mouse in SWC experiments exceeds 200 terabytes. It will take users several days to download these datasets over standard Internet connections.

Instead of bringing data to users, we will bring users to data, by storing datasets in the cloud (or in institutional clusters), and providing **cloud software to allow users to visually explore and statistically analyse behavioural and neural NaLoDuCo datasets where they live** (Figure 1, left box).

Our statistical analysis of neural time series will require knowledge of the spiking activity of single units; i.e., spike sorting. In long-duration experiments with freely moving animals spike sorting is a challenging problem, because movements of recording probes change the shape of spike waveforms over time and complicate the assignment of spikes to units based on their waveforms. We will address this problem by developing **spike sorting methods for long-duration and continual, long-duration and high-**

3

**channel-count recordings** (Figure 1, left box).

## Aim 2: create real-time machine learning methods for intelligent experimentation

In small-animal Neuroscience, most often statistical processing of neural time series is performed offline; i.e., experimental data is collected, saved to files, which are later statistically processed, with no runtime constraints. Most often all experimental data is processes at the same time; i.e., batch processing.

A new online statistical processing approach is now emerging in small-animal Neuroscience, where data is processes while it is being collected, and at the speed of data generation (**?**).

Online methods are well suited for NaLoDuCo experimentation. In experiments extending for weeks to months animals learn and adapt, their motivation and fatigue may fluctuate, and experimental conditions (e.g., lighting) may change. Offline batch processing algorithms cannot model this type of changing data. They assume stationary data whose statistical properties do not change across time. Differently, most online processing algorithms are robust to these changes. Also, NaLoDuCo experimentation is well suited for online methods, as the long-duration of these experiments provide a large amount of data to accurately fit expressive online methods.

We will **optimize methods developed for Aim 1 so that they can operate in real time**, and focus on the following two applications of these online methods (Figure 1, right box).

### Intelligent neuromodulation

Brain activity can be modulated optically, chemically and electrically (). Most commonly these modulations are done at fixed experimental times, or based on simple behavioral or neural observations. We will guide optogenetic manipulations based on inferences from advanced machine learning methods.

For example, a scientists may hypothesize that a peak in a neural latent variable, inferred from a prefrontal cortex population, signals the moment when mice decide to begin a foraging bout. To test this, she estimates latent variables from prefrontal cortex activity online, and forecasts when this peak will occur. She then optogenetically inactivates the neural population at

4

the forecasted time. Because the inactivations prevented the mouse from initiating a foraging bout, her hypothesis is supported.

**Intelligent experimental data storage**

As the duration of NaLoDuCo experiments becomes longer, and the volume of the behavioral and neural recordings becomes larger, it will be unfeasible to store all raw data. We will be forced to intelligently decide, in real time, subsets of data to discard. Real-time machine learning methods can guide the decision of what subset of data to discard, as exemplified below.

If we are recording videos from a mouse foraging in a large arena with ten high-resolution cameras, it would save considerable storage if at any time we only save videos from cameras capturing the mouse at that time. This could be done by tracking the position of the mouse in real time with probabilistic machine learning methods. Then, when the confidence of the tracking is high, we would only save videos of cameras capturing the mouse at the tracked position, but when the confidence is low, we would save all videos.

# 2 Approach

## 2.1 Offline machine learning

Extracting meaning from long-duration and continual recordings opens a few challenges and opportunities that we will address and exploit in this project, as we describe below. Prior to this discussion we provide a list of core machine learning models that we plan to use and disseminate for the characterization of NaLoDuCo recordings.

### 2.1.1 Core Machine Learning Algorithms

### 2.1.2 Challenges

**Non-stationarities**

Conventional offline methods used to characterize neural time series assume that the statistical characteristics of the modeled data do not change with time (i.e., that the probability of the data is time invariant – stationarity). This assumption may be valid for shorter experiments. However, for

5

long-duration experiments, where animals learn and adapt, where their motivation fluctuates, and their activity is modulated by circadian, utradiem and peridiem rhythms, the stationarity assumption may not hold.

The field of adaptive signal processing (**?**) develops algorithms to characterize non-stationary systems. There is no unique solution to optimally process any non-stationary signals. Instead, adaptations to specific algorithms have been developed to improve their performance in non-stationary environments. For example, the recursive least-squares algorithm (**?**, Chapter 9) is an adaptation of the ordinary least square algorithm to perform linear regression with non-stationary data. In the field of artificial neural networks, a large number of strategies have been developed to address data non-stationarity

### Long processing times for very large datasets

Neural and behavioral data analysis is most effective when computations can be performed quickly, ideally in real time. Very slow computations discourage data analysis, and hurts scientific discovery. The large dataset sizes generated by NaLoDuCo experimentation are an important challenge for fast data analysis.

To overcome this limitation, we will leverage distributed computing, a paradigm in which tasks and data are divided across multiple computers. Instead of relying on a single powerful machine, distributed computing accelerates processing by executing multiple parts of a computation in parallel.

We will develop parallel implementations of the of core machine learning algorithms for behavioral and neural data analysis (Section 2.1.1). These implementations will use Apache Spark[1] to parallelise pre-processing and feature extraction, and Ray[2] to parallelise machine learning and deep learning functionality.

---

[1] https://spark.apache.org/
[2] https://docs.ray.io/

6

**2.1.3  Opportunities**

**2.2   Visual Exploration**

**2.3   Spike Sorting**

**2.4   Online Machine Learning**

**2.5   Software and Infrastructure**