

1	Contents	
2	1 Summary	2
3	2 Core team	3
4	3 Application questions	3
5	3.1 Research theme	3
6	3.2 Vision	4
7	3.2.1 Context	4
8	3.2.2 Focus areas	4
9	3.3 Approach	5
10	3.4 US applicants	7
11	3.5 Resources	7
12	A More details about the vision	8
13	A.1 Context	8
14	A.2 Focus areas and their challenges	8
15	A.2.1 Data acquisition and management	9
16	A.2.2 Data dissemination	10
17	A.2.3 Data visualisation	10
18	A.2.4 Spike sorting	10
19	A.2.5 Data analysis	11
20	A.2.6 Experiments driven by real-time machine learning inference	12

1 Intention to submit document for the Work with
2 US researchers BBSRC-NSF/BIO lead agency
3 2024 funding opportunity

4 Enabling Naturalistic, Long-Duration and
5 Continual Neuroscience Experimentation with
6 Advanced Machine Learning

7
8 October 19, 2024

9 **1 Summary**

10 Word limit: 2 A summary is not required for this section, please write 'N/A' in
11 the textbox. Please still include a title for your project.
12 N/A

1 2 Core team

2 List the key members of your team and assign them roles from the following:

- 3 • project lead (PL)
- 4 • project co-lead (UK) (PcL)
- 5 • specialist
- 6 • professional enabling staff
- 7 • research and innovation associate
- 8 • technician
- 9 • researcher co-lead (RcL)

10 Only list one individual as project lead.

11 The core team section must only contain details of the UK applicants. The
12 US applicant information should be listed in the ‘US applicants’ section.

13 Find out more about [UKRI's core team roles in funding applications](#).

14 **project lead (PL)** Prof. Maneesh Sahani

15 **project co-lead (UK) (PcL)** Prof. Tiago Branco, Prof. Thomas Mrsic-Flogel

16 **researcher co-lead (UK) (RcL)** Dr. Joaquin Rapela, Dr. Dario Campagner

17 3 Application questions

18 3.1 Research theme

19 Word limit: 5 Please state the research theme you are applying under. Choose
20 one of the following research themes:

- 21 1. biological informatics
- 22 2. understanding host-microbe interactions
- 23 3. synthetic cells and cellular systems
- 24 4. synthetic microbial communities

25 biological informatics

1 3.2 Vision

2 Word limit: 500

3 What are you hoping to achieve with your proposed work?

4 What the assessors are looking for in your response

5 Your vision should clearly address:

- 6 • one of the opportunity research themes (biological informatics, under-
7 standing host-microbe interactions, synthetic cells and cellular systems or
8 synthetic microbial communities)
- 9 • the remit of the BBSRC and the NSF/BIO division associated with your
10 chosen research theme

11 References may be included within this section, but this will count towards
12 your word count.

13 Images are not required for this section.

14 3.2.1 Context

15 Conventional systems neuroscience experiments are typically short in duration
16 and often place significant constraints on subject behavior to simplify data anal-
17 ysis. However, these restrictions may limit our ability to observe critical aspects
18 of brain function and behavior that only manifest in more naturalistic and ex-
19 tended conditions.

20 At the Sainsbury Wellcome Centre (SWC) for Neural Circuits and Be-
21 haviour, we are pioneering Naturalistic, Long-Duration, and Continual (NaLo-
22 DuCo) foraging experiments in mice that span weeks to months. During these
23 extended experiments, we collect high-resolution recordings of both behavioral
24 and neural activity in naturalistic settings.

25 This novel experimental approach will enable researchers to explore neu-
26 ral mechanisms underlying naturalistic behavior over extended periods for the
27 first time, offering the possibility of uncovering insights across a wide range of
28 phenomena, including long-term behavioral adaptation, neural plasticity, and
29 learning. The data generated from NaLoDuCo experiments represent an en-
30 tirely new resource in neuroscience, with the potential to drive breakthroughs
31 and discoveries that are beyond the reach of traditional experiments.

32 Our vision is to empower research centers worldwide to adopt this ground-
33 breaking approach. However, the scale and complexity of the data generated
34 pose significant challenges in data acquisition, visualisation, and analysis. In
35 this proposal, we will address these challenges, developing and sharing openly
36 the necessary expertise, hardware, and software to enable this transformative
37 type of experimentation on a global scale.

38 3.2.2 Focus areas

39 Below, we outline the key focus areas we aim to address (Figure 3), along
40 with their associated challenges. These challenges primarily revolve around the

1 collection and analysis of continuously recorded, extremely large datasets—on
2 the order of hundreds of terabytes—gathered from experiments spanning weeks
3 to months.

4 While experiments in neuroscience that are naturalistic, long-duration, or
5 continuous have been conducted in the past (e.g., [Jhuang et al., 2010](#); [Mao](#)
6 [et al., 2021](#); [Voloh et al., 2023](#)), to the best of our knowledge, we are the first
7 to integrate all three of these features in a single experimental paradigm. This
8 combination introduces unprecedented complexities in data processing, as we
9 aim to capture behavior and brain activity in their most ecologically valid,
10 extended, and uninterrupted forms.

11 The focus areas of the proposed project are (Figure 3):

12 **Data Collection & Management** Efficiently gathering and organizing mas-
13 sive datasets over extended periods.

14 **Data Sharing** Providing easy access to large-scale datasets to researchers around
15 the globe using cloud-based technologies.

16 **Data Visualization** Developing efficient web-based tools to visualize very large
17 behavioral and neural datasets.

18 **Spike Sorting** Assigning spikes to neurons reliably, and tracking individual
19 neurons across long-periods of time in real time.

20 **Data Analysis** Evaluating existing methods, and developing new ones, when
21 necessary, to address key problems in behavioral and neural data analysis
22 (Figure 2).

23 **Inference-Driven Experimentation** Creating a new type of experimenta-
24 tion driven by real-time behavioral and neural inferences.

25 We have assembled a unique team to implement this project. The SWC is
26 a world leader in experimental neuroscience, working closely with the GCNU, a
27 renowned authority in computational neuroscience and machine learning. Both
28 institutions share the same building and collaborate extensively. NeuroGEARS
29 Ltd. leads the implementation of the NaLoDuCo experimental framework, while
30 Catalyst Neuro has played a pivotal role in developing and operating the DANDI
31 archive, in collaboration with Dr. Jeremy Magland, an expert in spike sorting,
32 data visualization, and cloud computing.

33 3.3 Approach

34 Word limit: 500

35 How are you going to deliver your proposed work?

36 What the assessors are looking for in your response

37 Your approach should give an overview highlighting:

- 38 • a clear description of the objectives and methodology for the proposed
39 work, including the contributions of the UK and US teams

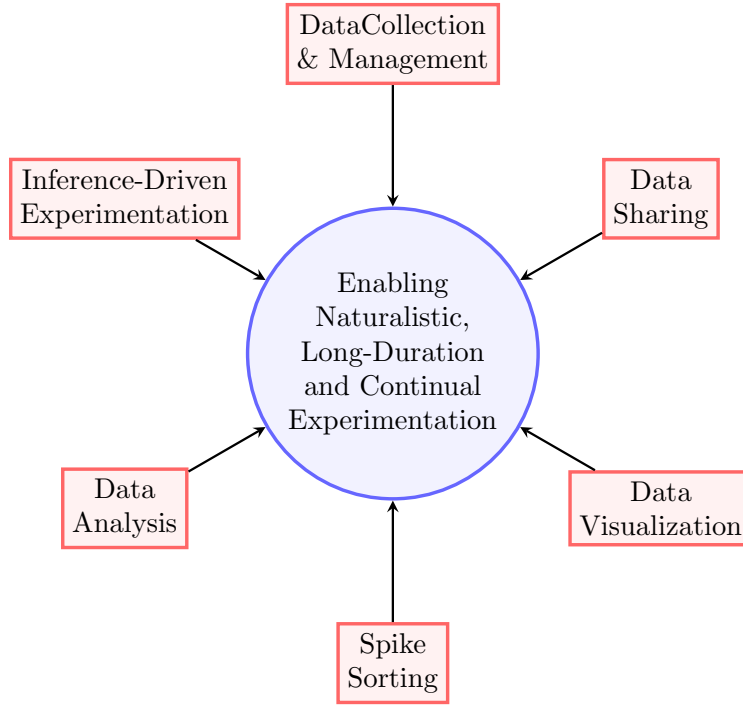


Figure 1: Project theme (blue) and focus areas (red).

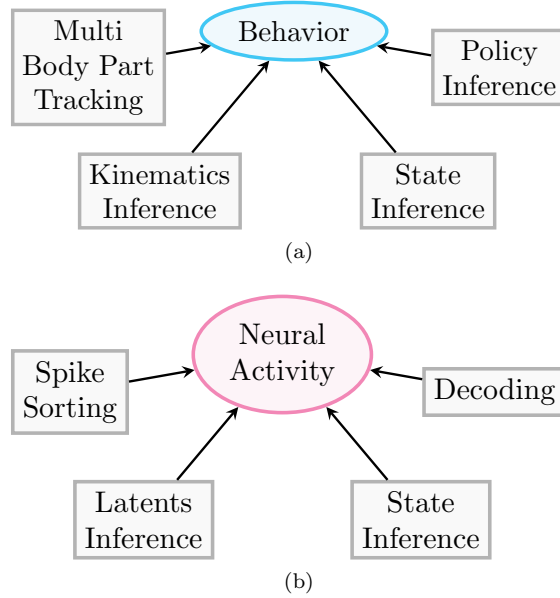


Figure 2: Behavioral (a) and neural (b) data analysis problems to address.

- 1 • the potential outputs and outcomes of the proposed work
- 2 References may be included within this section, but this will count towards
- 3 your word count.
- 4 Images are not required for this section.

5 **3.4 US applicants**

6 Word limit: 200

7 Please provide the following details of the US applicants on this application:

- 8 1. name
- 9 2. institute
- 10 3. job title
- 11 4. role in project (for example, project lead or project co-lead)
- 12 5. email address

13 Please also indicate who the lead US applicant will be.
14 NSF will use this information to confirm applicant eligibility.
15 Please do not include details of US applicants in the ‘Core team’ section.

16 **3.5 Resources**

17 Word limit: 200

18 Please provide the following:

- 19 • overall estimates for costings and staffing full time equivalent (FTE) for
- 20 both the UK and US components
- 21 • clear separation of UK and US costings, in pounds sterling and US dollars
- 22 (USD) respectively

23 The overall budget should be below the maximum £2 million combined fun-
24 der contribution

25 If there is more than one UK or US team associated with the application,
26 please combine their estimates together.

27 A detailed calculation and breakdown of resources is not required at this
28 stage, nor is a justification of costs.

29 The following is an example of how this might look.

30 UK Resources:

31 Total cost estimate: £600,000

32 Research council contribution: £480,000

33 0.2 FTE time, 1.0 FTE PDRA, 0.5 FTE technician

34 US Resources:

35 Total cost estimate: \$300,000

1 1.0 FTE PDRA or 1.0 FTE doctoral researcher
2 Total funder contribution estimate:
3 £716,475 (£480,000 + £236,475 (\$300,000 at exchange rate 0.79))

4 A More details about the vision

5 A.1 Context

6 Conventional systems neuroscience experiments are typically short in duration
7 and often place significant constraints on subject behavior to simplify data anal-
8 ysis. However, these restrictions may limit our ability to observe critical aspects
9 of brain function and behavior that only manifest in more naturalistic and ex-
10 tended conditions.

11 At the Sainsbury Wellcome Centre (SWC) for Neural Circuits and Be-
12 haviour, we are pioneering Naturalistic, Long-Duration, and Continual (NaLo-
13 DuCo) foraging experiments in mice that span weeks to months. During these
14 extended experiments, we collect high-resolution recordings of both behavioral
15 and neural activity in naturalistic settings. In collaboration with the Gatsby
16 Computational Neuroscience Unit (GCNU), we are developing novel analytical
17 methods to interpret this new class of data.

18 This novel experimental approach will enable researchers to explore neural
19 mechanisms underlying behavior over extended periods for the first time, of-
20 fering the possibility of uncovering insights across a wide range of phenomena,
21 including long-term behavioral adaptation, neural plasticity, and learning. The
22 data generated from NaLoDuCo experiments represent an entirely new resource
23 in neuroscience, with the potential to drive breakthroughs and discoveries that
24 are beyond the reach of traditional experiments.

25 Our vision is to empower research centers worldwide to adopt this ground-
26 breaking approach. However, the scale and complexity of the data generated
27 pose significant challenges in data acquisition, visualisation, and analysis. In
28 this proposal, we will address these challenges, developing and sharing openly
29 the necessary expertise, hardware, and software to enable this transformative
30 type of experimentation on a global scale.

31 A.2 Focus areas and their challenges

32 Below, we outline the key focus areas we aim to address (Figure 3), along
33 with their associated challenges. These challenges primarily revolve around the
34 collection and analysis of continuously recorded, extremely large datasets—on
35 the order of hundreds of terabytes—gathered from experiments spanning weeks
36 to months.

37 While experiments in neuroscience that are naturalistic, long-duration, or
38 continuous have been conducted in the past (e.g., Jhuang et al., 2010; Mao
39 et al., 2021; Voloh et al., 2023), to the best of our knowledge, we are the first
40 to integrate all three of these features in a single experimental paradigm. This

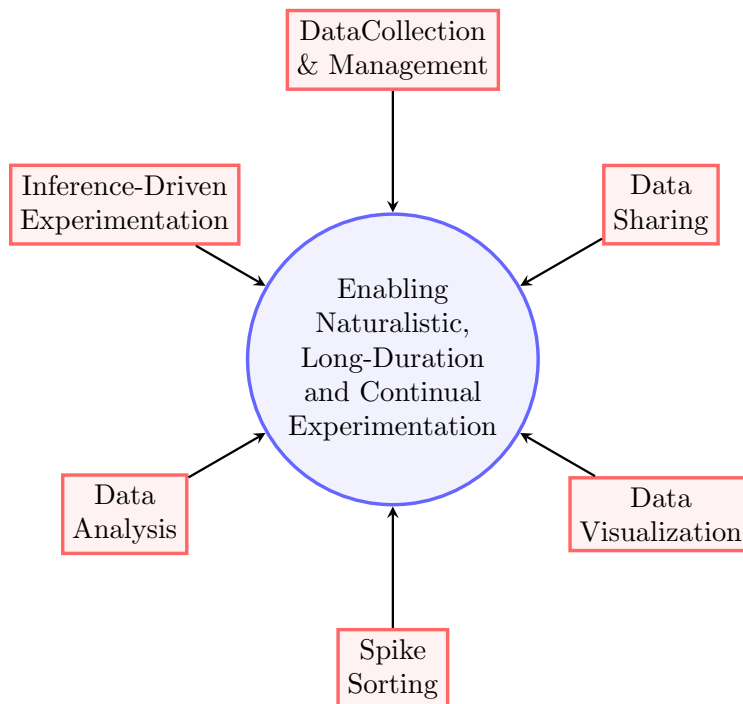


Figure 3: Project theme (blue) and focus areas (red).

1 combination introduces unprecedented complexities in data processing, as we
 2 aim to capture behavior and brain activity in their most ecologically valid,
 3 extended, and uninterrupted forms.

4 **A.2.1 Data acquisition and management**

5 At the SWC we have already performed foraging experiments in mice contin-
 6 uously collecting behavioral and experimental data 24 hours a day for seven
 7 days. We will share openly the specifications of the hardware used to build
 8 these experiments (e.g., instructions for building large foraging arenas, video
 9 cameras specifications, electrophysiological recording hardware), as well as the
 10 software we used for experimental control, data quality control, data access and
 11 management.

12 The data acquisition and management software used in our project is already
 13 publically available in GitHub¹. This software is already being used by scientists
 14 at the Allen Institue for Neural Dynamics and at Northwestern University. We
 15 will substantially improve its documentation to simplify its usage by external
 16 users.

¹https://github.com/SainsburyWellcomeCentre/aeon_mecha

1 Challenges related to data acquisition and management include data index-
2 ing to allow fast access to very large amount of saved data, online quality control
3 and alert systems to guarantee that anomalies in data collection are detected
4 and corrected with minimal delay, and synchronization between multiple data
5 streams.

6 **A.2.2 Data dissemination**

7 Datasets of the scale of hundreds of terabytes cannot be practically down-
8 loaded from data repositories. This is specially true for contiguous experiments
9 where unique insights are extracted by characterizing full datasets, and not only
10 parts of them. Therefore, we will store data in DANDI, which uses Amazon S3
11 buckets, and provide software in Amazon EC2 instances to visualize and analyze
12 data on the cloud, avoiding costly data transfers. That is, the large dataset sizes
13 of NaLoDuCo experiments make it impractical to distribute data to users and
14 require to bring users to data. Fortunately, cloud technologies are now mature
15 to allows this.

16 Importantly, if we distributed these very large datasets to users, only those
17 in large research centers would have the computing power to process them. But,
18 by deploying data and computing in the cloud, any person with Internet access
19 around the world will be able to benefit from them. Storing large datasets in
20 DANDI is free.

21 Dr. Ben Ditcher, founder of CatalystNeuro, has played a pivotal role in
22 supporting the development and operations of the DANDI archive.

23 **A.2.3 Data visualisation**

24 Visualisations are essential for scientific discovery. For the proposed project
25 visualisation present two major challenges. First, they need to display very large
26 datasets at different temporal scales, from milliseconds to weeks and months.
27 Second, as data and software will be deployed in the cloud, visualisation need
28 to be web based. Standard visualization tools cannot display terabyte sized
29 datasets. We will build custom web-based visualization tools to do this.

30 We have substantial experience building web-based visualization tools for
31 neurophysiological data. Dr. Jeremy Magland is now developing Neurosift² a
32 web-based visualizer for DANDI datasets.

33 **A.2.4 Spike sorting**

34 When electrodes are placed in the brain, they typically record spikes from mul-
35 tiple nearby neurons. Spike sorting attributes spikes to individual neurons.

36 Spike sorting is specially challenging for NaLoDuCo experiments. First,
37 because these experiments require to track individual neurons of freely moving
38 mice for weeks to months. Second, because spike sorting needs to be done

²<https://github.com/flatironinstitute/neurosift>

1 online, to allow experiments driven by real-time machine learning inference, as
2 described below.

3 Prof. Sahani pioneered the use of Bayesian inference methods for spike sort-
4 ing (Sahani, 1999). Dr. Jeremy Magland has significantly advanced the field of
5 spike sorting, particularly through his development of MountainSort³ and his
6 contributions to SpikeInterface⁴.

7 A.2.5 Data analysis

8 Advanced data analysis methods are indispensable to extract meaning from
9 NaLoDuCo experimental data. However, analyzing this data is challenging for
10 at least three reasons. First, important insights will most probably come from
11 the characterization of complete datasets, and not from subsets extracted from
12 them. Conventional batch methods cannot be used with datasets of the size
13 produced by NaLoDuCo experiments. For instance, for learning, batch linear re-
14 gression cannot load into memory and invert a data matrix with high-resolution
15 observations from a one-month-long experiment. Thus, **online methods** that
16 can process infinite data streams become mandatory.

17 Second, a pervasive assumption in most ML algorithms is stationarity; i.e.,
18 the assumption that the statistics of data do not change over time. But in long-
19 duration and continuous experiments this assumption is most often violated
20 as, for example, the arousal of subjects changes. Hence, the analysis of data
21 generated by these experiments requires **adaptive methods**.

22 Third, statistical algorithms consist of two key stages: learning (or training)
23 and inference (or prediction). The learning stage identifies model parameters,
24 and the inference stage uses the learned model to make predictions, or infer
25 latent variables, from new unseen data. Frequently training is performed on a
26 small subset of a dataset, and inference is done on the remaining data. However,
27 since in long-duration and continual experiments behavior and neural activity
28 are generally not stationary, it is not optimal to train models on data subsets and
29 use them to make inferences on the remaining data, since the state of the animal
30 at training and inference times may be different. To overcome this difficulty we
31 will use **continual learning methods**.

32 We will evaluate methods to analyze different aspects of behavior and neu-
33 ral activity (Figure ??). We will test how these methods process very large
34 datasets, how they handle non-stationary data, and how feasible is to retrain
35 them to adapt to changing conditions. We will adapt these methods so that they
36 better address these challenges and, when needed, develop new ones. We will
37 carefully report the outcomes of these evaluations so that researchers performing
38 NaLoDuCo experimentation can choose the best methods that suit their needs.

³<https://github.com/flatironinstitute/mountainsort5>

⁴<https://github.com/spikeinterface/spikeinterface>

1 A.2.6 Experiments driven by real-time machine learning inference

2 Small animal experiments are usually controlled by simple static rules or direct
3 behavioral observations. Funded by a BBSRC award⁵ we are developing soft-
4 ware to allow a new type of experimental control based on statistical inferences
5 made on behavioral and/or neural measurements.

6 For example, after inferring latent variables from neural activity and observ-
7 ing that one of these latents have crossed a threshold, we can deliver a reward (as
8 done in learning to control a BCI; Clancy and Masic-Flogel, 2021), or perform
9 an action (as done in motor imagery BCI; Lebedev and Nicolelis, 2006), or ma-
10 nipulate of neural activity (as done when studying the causal relation between a
11 pattern of brain activity and behavior; Deisseroth, 2015). We propose to further
12 develop the previous software and use it to test causal effects of neural activity
13 patterns on foraging decisions using our NaLoDuCo foraging experiments.

14 Building experiments driven by real-time machine learning inferences brings
15 at least two challenges. The first one is a machine learning problem, how to
16 build fast inferences that can operate in real time. The second one is a neuro-
17 science problem, how to identify neuroscience experiments suitable to real-time
18 control, and then perform the experiment with real-time control. Fortunately
19 at the Gatsby Unit we are experienced on building advanced machine learning
20 algorithms to address the first challenge. And at the SWC we perform many so-
21 phisticated animal experiments that could benefit from real-time experimental
22 control.

23 In summary, we are pioneering a new paradigm in neuroscience experimen-
24 tation, driven by advanced inferential methods applied to rich behavioral and
25 neural recordings. This innovative technology has the potential to transform
26 the field, enabling experiments that were previously unimaginable. By leverag-
27 ing these sophisticated inferences, we may unlock new dimensions of knowledge
28 that could not be achieved through simpler, conventional approaches. This
29 breakthrough could open doors to insights that redefine our understanding of
30 brain-behavior relationships.

31 References

- 32 Clancy, K. B. and Masic-Flogel, T. D. (2021). The sensory representation of
33 causally controlled objects. *Neuron*, 109(4):677–689.
- 34 Deisseroth, K. (2015). Optogenetics: 10 years of microbial opsins in neuro-
35 science. *Nature neuroscience*, 18(9):1213–1225.
- 36 Jhuang, H., Garrote, E., Yu, X., Khilnani, V., Poggio, T., Steele, A. D., and
37 Serre, T. (2010). Automated home-cage behavioural phenotyping of mice.
38 *Nature communications*, 1(1):68.

⁵<https://gow.bbsrc.ukri.org/grants/AwardDetails.aspx?FundingReference=BB%2FW019132%2F1>

- 1 Lebedev, M. A. and Nicolelis, M. A. (2006). Brain-machine interfaces: past,
2 present and future. *TRENDS in Neurosciences*, 29(9):536–546.
- 3 Mao, D., Avila, E., Caziot, B., Laurens, J., Dickman, J. D., and Angelaki,
4 D. E. (2021). Spatial modulation of hippocampal activity in freely moving
5 macaques. *Neuron*, 109(21):3521–3534.
- 6 Sahani, M. (1999). *Latent variable models for neural data analysis*. California
7 Institute of Technology.
- 8 Voloh, B., Maisson, D. J.-N., Cervera, R. L., Conover, I., Zambre, M., Hayden,
9 B., and Zimmermann, J. (2023). Hierarchical action encoding in prefrontal
10 cortex of freely moving macaques. *Cell reports*, 42(9).