

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	5
9	3.3 Approach	5
10	3.3.1 Data collection & management	7
11	3.3.2 Data sharing	7
12	3.3.3 Data visualisation	7
13	3.3.4 Spike sorting	7
14	3.3.5 Data analysis	9
15	3.3.6 Inference-driven experimentation	9
16	3.4 US applicants	10
17	3.5 Resources	10
18	References	11
19	A More details about the vision	14
20	A.1 Context	14
21	A.2 Focus areas and their challenges	15
22	A.2.1 Data acquisition and management	15
23	A.2.2 Data dissemination	17
24	A.2.3 Data visualisation	17
25	A.2.4 Spike sorting	17
26	A.2.5 Data analysis	18
27	A.2.6 Experiments driven by real-time machine learning inference	18

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 23, 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

2 Core team

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

- project lead (PL)
- project co-lead (UK) (PcL)
- specialist
- professional enabling staff
- research and innovation associate
- technician
- researcher co-lead (RcL)

Only list one individual as project lead.

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

Find out more about UKRI's core team roles in funding applications.

project lead (PL) Prof. Maneesh Sahani

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

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

3 Application questions

3.1 Research theme

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

1. biological informatics
 2. understanding host-microbe interactions
 3. synthetic cells and cellular systems
 4. synthetic microbial communities
- 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. In collaboration with the Gatsby
25 Computational Neuroscience Unit (GCNU), we are developing novel analytical
26 methods to interpret this new class of data.

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

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

1 3.2.2 Focus areas

2 Below, we outline the key focus areas we aim to address (Figure 4). Challenges
3 addressing these areas primarily revolve around the collection and analysis of
4 continuously recorded, extremely large datasets—on the order of hundreds of
5 terabytes—gathered from experiments spanning weeks to months.

6 While experiments in neuroscience that are naturalistic, long-duration, or
7 continuous have been conducted in the past [e.g., 14, 19, 32], to the best of our
8 knowledge, we are the first to integrate all three of these features in a single ex-
9 perimental paradigm. This combination introduces unprecedented complexities
10 in data processing, as we aim to capture behavior and brain activity in their
11 most ecologically valid, extended, and uninterrupted forms.

12 The focus areas of the proposed project are (Figure 4):

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

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

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

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

21 **Data Analysis** Evaluating existing methods, and developing new ones, when
22 necessary, to study key behavioral and neural-coding problems with NaLo-
23 DuCo experimental data (Figure 2).

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

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

34 3.3 Approach

35 Word limit: 500

36 How are you going to deliver your proposed work?

37 What the assessors are looking for in your response

38 Your approach should give an overview highlighting:

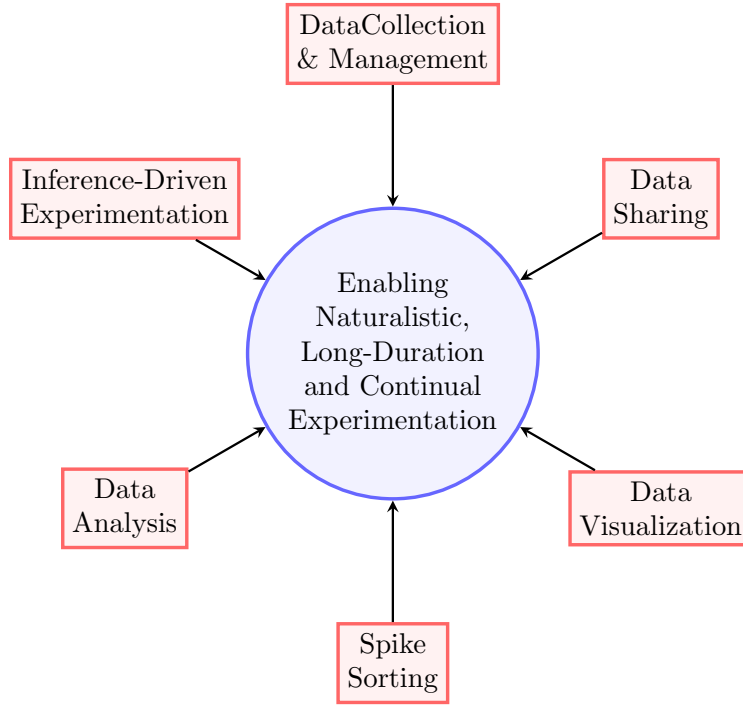


Figure 1: Project theme (blue) and focus areas (red).

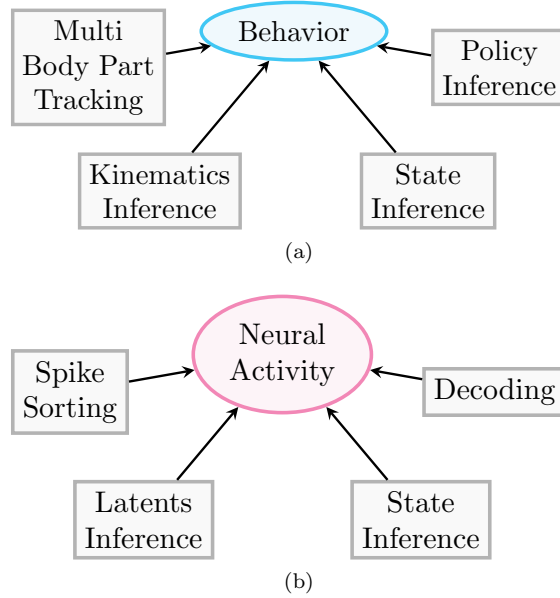


Figure 2: Behavioral (a) and neural (b) data analysis problems to address.

- 1 • a clear description of the objectives and methodology for the proposed
- 2 work, including the contributions of the UK and US teams
- 3 • the potential outputs and outcomes of the proposed work
- 4 References may be included within this section, but this will count towards
- 5 your word count.
- 6 Images are not required for this section.

7 **3.3.1 Data collection & management**

8 We have developed a new platform that allows housing of mice in large arenas
9 (>2m diameter), while manipulating and monitoring their behaviour at high
10 spatiotemporal resolution [Figure 3, 2]. We have openly shared software for
11 supporting data acquisition [11] and management [12] in this arena. Using this
12 platform we have collected several week long datasets both with single mouse
13 and multiple mice.

14 **3.3.2 Data sharing**

15 The large dataset sizes generated by NaLoDuCo experiments, on the order of
16 hundreds of terabytes, make it impractical to distribute data to users, and
17 require to bring users to data. Fortunately, cloud technologies are now mature
18 to allow this. We will store data in the Distributed Archives for Neuroscience
19 Data Integration (DANDI), which uses Amazon S3 buckets, and we will provide
20 software to visualize and analyze data in Amazon EC2 instances, to avoid costly
21 data transfers.

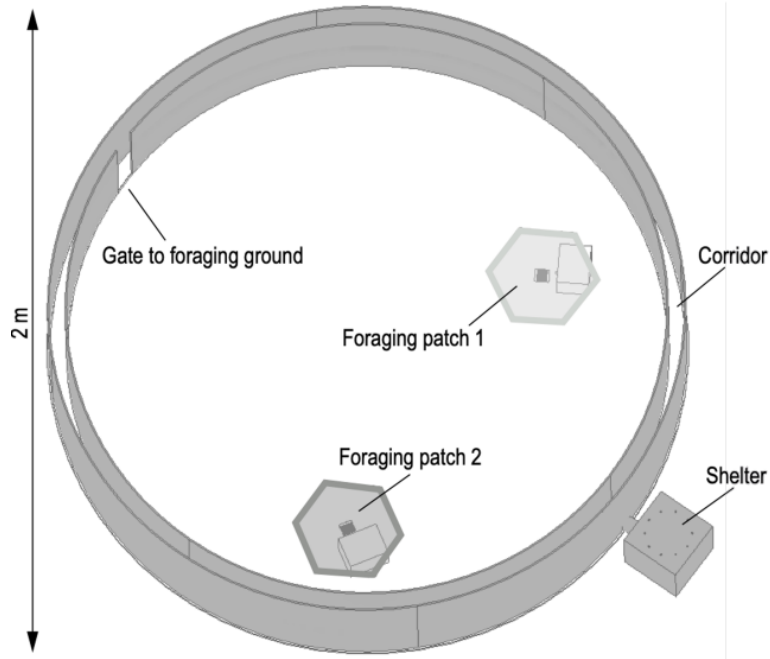
22 **3.3.3 Data visualisation**

23 Our visualisation tools need to display very large datasets at different temporal
24 scales, from milliseconds to weeks and months, and they need to be web based.
25 We will use multi-resolution visualization techniques, which store data at various
26 resolutions, and use the appropriate resolution for each zoom level. Web-based
27 visualisation will be optimized using web workers [8].

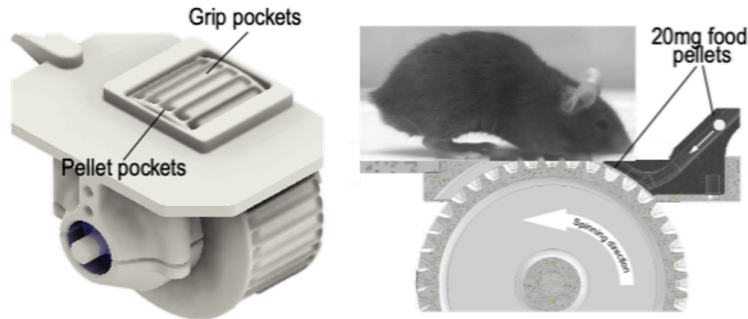
28 **3.3.4 Spike sorting**

29 Spike sorting is specially challenging in NaLoDuCo experimentation since we
30 want to track individual neurons of freely moving mice for weeks to months. In
31 addition, we need online spike sorting, to allow experiments driven by real-time
32 machine learning inference, as described below.

33 We will evaluate methods for tracking neurons over long periods of time [e.g.,
34 34, 31] and for online sorting [e.g., 26, 29].



(a)



(b)

Figure 3: Foraging arena (a) and feeder (b). The floors of the arenas are tessellated to form honeycombs of modular hexagonal tiles (a), each of which can be equipped with a newly designed underground feeder (b). Pellets are dispensed onto a foraging wheel once the mouse has spun it for a pre-defined programmable distance threshold using its forepaws (fictive digging). The arena contains up to six scale-equipped nesting modules that allows housing of mice in the arena and weight monitoring. Behavioural monitoring is achieved by an array of high-speed cameras (up to 15), by which mouse location, mouse identity and body parts can be track in real time. Long term monitoring of neural activity is performed using Neuropixels probes.

1 3.3.5 Data analysis

2 The very large size of NaLoDuCo experimental data, the fact that the statistics
3 of these data change across time, and the requirement for real-time and close-
4 loop inference create new challenges to conventional machine learning methods.
5 We will evaluate existing methods targeting the experimental problems in Fig-
6 ure 2 and, if necessary, modify them, or create new ones, to address the previous
7 challenges.

8 For behavioral data, we will evaluate methods to:

- 9 • track multiple body parts of animals [e.g., 20, 23, 1, and a switching-linear-
10 dyanamical method using RFIDs that we will develop],
- 11 • infer kinematics of foraging mice [e.g., 24, 3],
- 12 • segment behavior into discrete states [e.g., 33, 13, and a hierarchical HMM
13 that we will develop],
- 14 • infer the rules that govern mice behavior from behavioral observations
15 only (i.e., policy inference) [e.g., 36, 35].

16 For neural data, we will evaluate methods to:

- 17 • estimate low-dimensional continual representations of high-dimensional
18 spiking activity (i.e., latents inference) [e.g., 18, 9, 22, 28],
- 19 • segment neural activity into discrete states [e.g., 4, 10],
- 20 • decode environment variables from neural activity [e.g., 7, 15, 30].

21 3.3.6 Inference-driven experimentation

22 We call inference-driven experimentation to a type of experimentation driven
23 by machine learning inferences on neural or behavioral data, where the result
24 of these inferences can change the experiment in real time.

25 We will apply inference-driven experimentation to test if patterns of neural
26 activity are causally related to foraging behaviors. We would first check that
27 a pattern of neural activity always precedes a given foraging behavior. We
28 would then detect the occurrence of the pattern and in real time optogenetically
29 inactivate the neurons responsible for the pattern. If the behavior disappears
30 the causality argument would be supported.

31 For this we will use the Bonsai ecosystem for experimental control [17] and
32 online machine learning functionality that we are adding to Bonsai [25], funded
33 by a BBSRC award [21].

1 3.4 US applicants

2 Word limit: 200

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

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

9 Please also indicate who the lead US applicant will be.

10 NSF will use this information to confirm applicant eligibility.

11 Please do not include details of US applicants in the ‘Core team’ section.

12 3.5 Resources

13 Word limit: 200

14 Please provide the following:

- 15 • overall estimates for costings and staffing full time equivalent (FTE) for
16 both the UK and US components
- 17 • clear separation of UK and US costings, in pounds sterling and US dollars
18 (USD) respectively

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

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

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

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

26 UK Resources:

27 Total cost estimate: £600,000

28 Research council contribution: £480,000

29 0.2 FTE time, 1.0 FTE PDRA, 0.5 FTE technician

30 US Resources:

31 Total cost estimate: \$300,000

32 1.0 FTE PDRA or 1.0 FTE doctoral researcher

33 Total funder contribution estimate:

34 £716,475 (£480,000 + £236,475 (\$300,000 at exchange rate 0.79))

1 References

2 References

- 3 [1] Dan Biderman, Matthew R Whiteway, Cole Hurwitz, Nicholas Greenspan,
4 Robert S Lee, Ankit Vishnubhotla, Richard Warren, Federico Pedraja, Dil-
5 lon Noone, Michael M Schartner, et al. Lightning pose: improved ani-
6 mal pose estimation via semi-supervised learning, bayesian ensembling and
7 cloud-native open-source tools. *Nature Methods*, pages 1–13, 2024.
- 8 [2] D. Campagner, J. Bhagat, G. Lopes, L. Calcaterra, J. Ahn, A. Almeida,
9 F. J. Carvalho, B. Cruz, A. Erskine, C. Lo, T. T. Nguyen, A. Pouget,
10 J. Rapela, T. Ryan, J. Reggiani, and S. SWC Foraging Behaviour Work-
11 ing Group. Aeon: an open-source platform to study the neural basis of
12 ethological behaviours over naturalistic timescales. In *Society for Neuro-*
13 *science Abstracts*, page PST033.03 / I26, 2024. Presented at the *Society*
14 *for Neuroscience Annual Meeting*.
- 15 [3] Subhash Challa, Mark R. Morelande, Darko Mušicki, and Robin J. Evans.
16 *Fundamentals of Object Tracking*. Cambridge University Press, 2011.
- 17 [4] Zhe Chen, Sujith Vijayan, Riccardo Barbieri, Matthew A Wilson, and
18 Emery N Brown. Discrete-and continuous-time probabilistic models and
19 algorithms for inferring neuronal up and down states. *Neural computation*,
20 21(7):1797–1862, 2009.
- 21 [5] Kelly B Clancy and Thomas D Mrsic-Flogel. The sensory representation
22 of causally controlled objects. *Neuron*, 109(4):677–689, 2021.
- 23 [6] Karl Deisseroth. Optogenetics: 10 years of microbial opsins in neuroscience.
24 *Nature neuroscience*, 18(9):1213–1225, 2015.
- 25 [7] Xinyi Deng, Daniel F Liu, Kenneth Kay, Loren M Frank, and Uri T Eden.
26 Clusterless decoding of position from multiunit activity using a marked
27 point process filter. *Neural computation*, 27(7):1438–1460, 2015.
- 28 [8] MDN Web Docs. Web workers. [https://developer.mozilla.org/](https://developer.mozilla.org/en-US/docs/Web/API/Web_Workers_API)
29 [en-US/docs/Web/API/Web_Workers_API](https://developer.mozilla.org/en-US/docs/Web/API/Web_Workers_API), 2024. Web Workers makes it
30 possible to run a script operation in a background thread separate from
31 the main execution thread of a web application. The advantage of this is
32 that laborious processing can be performed in a separate thread, allow-
33 ing the main (usually the UI) thread to run without being blocked/slowed
34 down.
- 35 [9] Lea Duncker and Maneesh Sahani. Temporal alignment and latent gaussian
36 process factor inference in population spike trains. In *Advances in Neural*
37 *Information Processing Systems*, pages 10445–10455, 2018.

- 1 [10] Sean Escola, Alfredo Fontanini, Don Katz, and Liam Paninski. Hidden
2 markov models for the stimulus-response relationships of multistate neural
3 systems. *Neural computation*, 23(5):1071–1132, 2011.
- 4 [11] SWC Foraging Behavior Working Group. aeon_acquisition repository.
5 https://github.com/SainsburyWellcomeCentre/aeon_acquisition,
6 2024. Task control and acquisition systems for Project Aeon.
- 7 [12] SWC Foraging Behavior Working Group. aeon_mecha repository. [https:](https://github.com/SainsburyWellcomeCentre/aeon_mecha)
8 [//github.com/SainsburyWellcomeCentre/aeon_mecha](https://github.com/SainsburyWellcomeCentre/aeon_mecha), 2024. Project
9 Aeon’s main library for interfacing with acquired data. Contains modules
10 for raw data file input/output, data querying, data processing, data quality
11 control, database ingestion, and building computational data pipelines.
- 12 [13] Alexander I Hsu and Eric A Yttri. B-soid, an open-source unsupervised
13 algorithm for identification and fast prediction of behaviors. *Nature com-*
14 *munications*, 12(1):5188, 2021.
- 15 [14] Hueihan Jhuang, Estibaliz Garrote, Xinlin Yu, Vinita Khilnani, Tomaso
16 Poggio, Andrew D Steele, and Thomas Serre. Automated home-cage be-
17 havioural phenotyping of mice. *Nature communications*, 1(1):68, 2010.
- 18 [15] Fabian Kloosterman, Stuart P Layton, Zhe Chen, and Matthew A Wilson.
19 Bayesian decoding using unsorted spikes in the rat hippocampus. *Journal*
20 *of neurophysiology*, 111(1):217–227, 2014.
- 21 [16] Mikhail A Lebedev and Miguel AL Nicolelis. Brain-machine interfaces:
22 past, present and future. *TRENDS in Neurosciences*, 29(9):536–546, 2006.
- 23 [17] NeuroGEARS Ltd. Bonsai. <https://bonsai-rx.org/>, 2024. A visual
24 language for reactive programming.
- 25 [18] Jakob H Macke, Lars Buesing, John P Cunningham, Byron M Yu, Kr-
26 ishna V Shenoy, and Maneesh Sahani. Empirical models of spiking in neu-
27 ral populations. *Advances in neural information processing systems*, 24,
28 2011.
- 29 [19] Dun Mao, Eric Avila, Baptiste Caziot, Jean Laurens, J David Dickman,
30 and Dora E Angelaki. Spatial modulation of hippocampal activity in freely
31 moving macaques. *Neuron*, 109(21):3521–3534, 2021.
- 32 [20] Alexander Mathis, Pranav Mamidanna, Kevin M Cury, Taiga Abe,
33 Venkatesh N Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge.
34 Deeplabcut: markerless pose estimation of user-defined body parts with
35 deep learning. *Nature neuroscience*, 21(9):1281–1289, 2018.
- 36 [21] Thomas Mrsic-Flogel. Machine intelligence for neuroscience experi-
37 mental control. [https://gow.bbsrc.ukri.org/grants/AwardDetails.](https://gow.bbsrc.ukri.org/grants/AwardDetails.aspx?FundingReference=BB%2FW019132%2F1)
38 [aspx?FundingReference=BB%2FW019132%2F1](https://gow.bbsrc.ukri.org/grants/AwardDetails.aspx?FundingReference=BB%2FW019132%2F1), 2023. BBSRC award
39 BB/W019132/1.

- 1 [22] Chethan Pandarinath, Daniel J O’Shea, Jasmine Collins, Rafal Jozefowicz, Sergey D Stavisky, Jonathan C Kao, Eric M Trautmann, Matthew T Kaufman, Stephen I Ryu, Leigh R Hochberg, et al. Inferring single-trial neural population dynamics using sequential auto-encoders. *Nature methods*, 15(10):805–815, 2018.
- 6 [23] Talmo D Pereira, Nathaniel Tabris, Arie Matsliah, David M Turner, Junyu Li, Shruthi Ravindranath, Eleni S Papadoyannis, Edna Normand, David S Deutsch, Z Yan Wang, et al. Slep: A deep learning system for multi-animal pose tracking. *Nature methods*, 19(4):486–495, 2022.
- 10 [24] Joaquin Rapela. Linear dynamical systems in python. https://github.com/joacorapela/lds_python, 2024. Python code to estimate Gaussian linear dynamical systems.
- 13 [25] Joaquin Rapela, Nicholas Guilbeault, and Goncalo Lopes. Bonsai machine learning. <https://bonsai-rx.org/machinelearning/>, 2024. Machine learning functionality for experimental control in Bonsai.
- 16 [26] Ueli Rutishauser, Erin M Schuman, and Adam N Mamelak. Online detection and sorting of extracellularly recorded action potentials in human medial temporal lobe recordings, in vivo. *Journal of neuroscience methods*, 154(1-2):204–224, 2006.
- 20 [27] Maneesh Sahani. *Latent variable models for neural data analysis*. Phd thesis, California Institute of Technology, Pasadena, CA, May 1999. Available at <https://www.proquest.com/docview/304498531>.
- 23 [28] Omid G Sani, Hamidreza Abbaspourazad, Yan T Wong, Bijan Pesaran, and Maryam M Shanechi. Modeling behaviorally relevant neural dynamics enabled by preferential subspace identification. *Nature Neuroscience*, 24(1):140–149, 2021.
- 27 [29] Gopal Santhanam, Maneesh Sahani, Stephen I Ryu, and Krishna V Shenoy. An extensible infrastructure for fully automated spike sorting during online experiments. In *The 26th Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, volume 2, pages 4380–4384. IEEE, 2004.
- 32 [30] Ardi Tampuu, Tambet Matiisen, H Freyja Ólafsdóttir, Caswell Barry, and Raul Vicente. Efficient neural decoding of self-location with a deep recurrent network. *PLoS computational biology*, 15(2):e1006822, 2019.
- 35 [31] Enny H van Beest, Célian Bimbard, Julie MJ Fabre, Sam W Dodgson, Flóra Takács, Philip Coen, Anna Lebedeva, Kenneth D Harris, and Matteo Carandini. Tracking neurons across days with high-density probes. *Nature Methods*, pages 1–10, 2024.

- 1 [32] Benjamin Voloh, David J-N Maisson, Roberto Lopez Cervera, Indirah
2 Conover, Mrunal Zambre, Benjamin Hayden, and Jan Zimmermann. Hi-
3 erarchical action encoding in prefrontal cortex of freely moving macaques.
4 *Cell reports*, 42(9), 2023.
- 5 [33] Alexander B Wiltchko, Matthew J Johnson, Giuliano Iurilli, Ralph E Pe-
6 terson, Jesse M Katon, Stan L Pashkovski, Victoria E Abaira, Ryan P
7 Adams, and Sandeep Robert Datta. Mapping sub-second structure in
8 mouse behavior. *Neuron*, 88(6):1121–1135, 2015.
- 9 [34] Augustine Xiaoran Yuan, Jennifer Colonell, Anna Lebedeva, Michael Okun,
10 Adam S Charles, and Timothy D Harris. Multi-day neuron tracking in high-
11 density electrophysiology recordings using earth mover’s distance. *Elife*,
12 12:RP92495, 2024.
- 13 [35] Hao Zhu, Brice De La Crompe, Gabriel Kalweit, Artur Schneider, Maria
14 Kalweit, Ilka Diester, and Joschka Boedecker. L(m)v-iql: Multiple intention
15 inverse reinforcement learning for animal behavior characterization. *arXiv*
16 *preprint arXiv:2311.13870*, 2023.
- 17 [36] Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al.
18 Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages
19 1433–1438. Chicago, IL, USA, 2008.

20 A More details about the vision

21 A.1 Context

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

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

34 This novel experimental approach will enable researchers to explore neural
35 mechanisms underlying behavior over extended periods for the first time, of-
36 fering the possibility of uncovering insights across a wide range of phenomena,
37 including long-term behavioral adaptation, neural plasticity, and learning. The
38 data generated from NaLoDuCo experiments represent an entirely new resource

1 in neuroscience, with the potential to drive breakthroughs and discoveries that
2 are beyond the reach of traditional experiments.

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

9 **A.2 Focus areas and their challenges**

10 Below, we outline the key focus areas we aim to address (Figure 4), along
11 with their associated challenges. These challenges primarily revolve around the
12 collection and analysis of continuously recorded, extremely large datasets—on
13 the order of hundreds of terabytes—gathered from experiments spanning weeks
14 to months.

15 While experiments in neuroscience that are naturalistic, long-duration, or
16 continuous have been conducted in the past [e.g., 14, 19, 32], to the best of our
17 knowledge, we are the first to integrate all three of these features in a single ex-
18 perimental paradigm. This combination introduces unprecedented complexities
19 in data processing, as we aim to capture behavior and brain activity in their
20 most ecologically valid, extended, and uninterrupted forms.

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

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

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

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

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

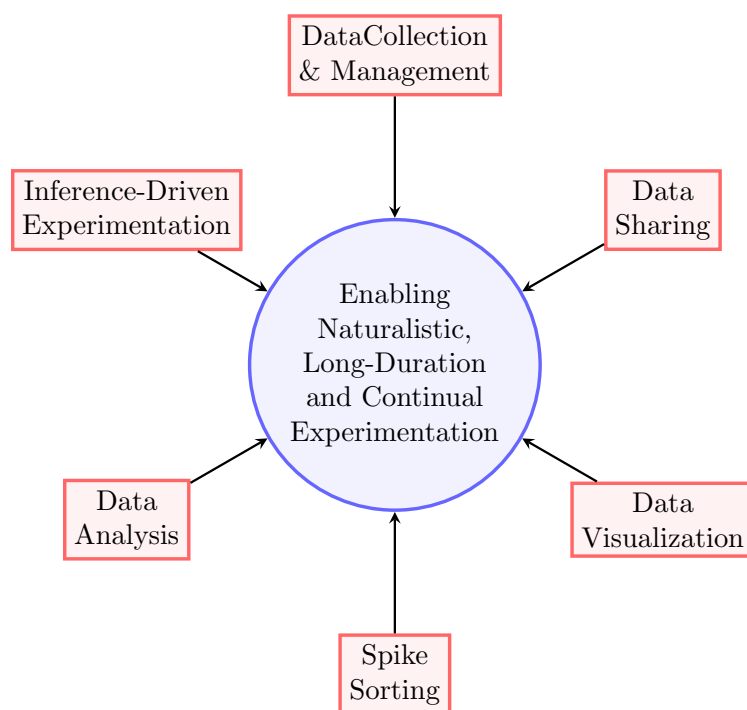


Figure 4: Project theme (blue) and focus areas (red).

1 **A.2.2 Data dissemination**

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

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

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

18 **A.2.3 Data visualisation**

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

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

28 **A.2.4 Spike sorting**

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

31 Spike sorting is specially challenging for NaLoDuCo experiments. First,
32 because these experiments require to track individual neurons of freely moving
33 mice for weeks to months. Second, because spike sorting needs to be done
34 online, to allow experiments driven by real-time machine learning inference, as
35 described below.

36 Prof. Sahani pioneered the use of Bayesian inference methods for spike sort-
37 ing [27]. Dr. Jeremy Magland has significantly advanced the field of spike sort-
38 ing, particularly through his development of MountainSort³ and his contribu-
39 tions to SpikeInterface⁴.

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

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

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

1 A.2.5 Data analysis

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

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

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

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

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

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

38 For example, after inferring latent variables from neural activity and observ-
39 ing that one of these latents have crossed a threshold, we can deliver a reward
40 [as done in learning to control a BCI; 5], or perform an action [as done in motor
41 imagery BCI; 16], or manipulate of neural activity [as done when studying the

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

1 causal relation between a pattern of brain activity and behavior; 6]. We pro-
2 pose to further develop the previous software and use it to test causal effects
3 of neural activity patterns on foraging decisions using our NaLoDuCo foraging
4 experiments.

5 Building experiments driven by real-time machine learning inferences brings
6 at least two challenges. The first one is a machine learning problem, how to
7 build fast inferences that can operate in real time. The second one is a neuro-
8 science problem, how to identify neuroscience experiments suitable to real-time
9 control, and then perform the experiment with real-time control. Fortunately
10 at the Gatsby Unit we are experienced on building advanced machine learning
11 algorithms to address the first challenge. And at the SWC we perform many so-
12 phisticated animal experiments that could benefit from real-time experimental
13 control.

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