

## Section 1: Contact information for the Israeli PI

Prof. Elisha Moses, Weizmann Institute of Science, +917-8-934-3139, [elisha.moses@weizmann.ac.il](mailto:elisha.moses@weizmann.ac.il)

## Section 2: IOS Program Justification

BIO/IOS Neural Systems NSF 11-572. Our proposed work is most appropriate for the IOS neural systems section as it deals with fundamental questions relating to the computational capabilities of the central nervous systems of living organisms. The associated questions we intend to address are fundamental to an integrative understanding of structure-function physiology and co-evolution in all organisms with nervous systems and the evolutionary predecessors that gave rise to them. With respect to nervous system *organization*, we intend to deepen previous explorations into the interaction of developmental and environmental constraints on the architecture and complementary computational capacity of networks of living neurons. Focusing on flexible abstract models of computation, rather than those specific to digital electronics, will enable us to investigate potentially novel computational potential implemented via architectures that differ from digital electronics in their apparent ability to take advantage of synergy between stochastic and deterministic properties. We may therefore, in addition to contributing to fundamental biological understanding, identify computational principles unique to natural neural networks. It is essential in this project to link theoretical, computational and experimental methods of inquiry. This is necessary to address questions we ask regarding co-evolution and co-development between underlying biological organization and computational tasks such structures are capable of performing in a variety of environmental conditions. We have assembled a culturally diverse and interdisciplinary team of experienced and developing scientists to enable the success of the research program we propose.

## Section 3: Research description

### Motivating questions

One of the big mysteries regarding the brain is, why neurons that perform amazingly complex calculations become inefficient at computation once they are grown outside of it. In particular, once you take hippocampal neurons out of the rat brain and culture them on a dish, they lose their individuality. Living neuronal cultures are characterized by all-or-none network bursts, where practically all the neurons fire simultaneously, with varying degrees of synchrony. While in the brain they can compute the trajectory of the animal in a maze, but as an ensemble grown out of it they are unable to create new computations, and can only carry very few bits of information [1]. How is it then that individual neurons are organized in the brain as a splendid computer, but in the dish their computational repertoire seems limited? Obviously, their environment and growth conditions have changed, but can we isolate and pin down those elements that are necessary, and perhaps also sufficient, to support computation by a living neuronal network? In an effort to tackle such questions, we have embarked on an investigation of the computational abilities of living neuronal networks, and propose here to expand this investigation both theoretically and experimentally in several fundamental ways.

On the conceptual side, a number of deep questions arise. What are the structural properties at the level of networks of neurons, modules of networks of neurons, and perhaps higher order forms of organization necessary to support the capacity for abstraction, which is fundamental to computation [2] and may likewise be fundamental to *cognitive architecture*, development, and function [3]? Given that we have shown boolean logic devices capable of being implemented *in vitro* using neurons [4], is it possible to biologically engineer analogous devices capable of performing

computations that have natural embeddings within **first** or **higher-order logic**? Moreover, can we define environmental conditions that lead to the development of an extensionally equivalent machine, but whose architecture may differ from the human-engineered biological implementation? Would multiple developmental implementations conserve identifiable structural features? And, how do the structural properties of networks with the capacity for zeroth, first, or higher order computations compare?

We plan to experimentally design a number of novel logical devices, as well as use a number of new technologies that we have acquired for the construction of such computational constructs. Notable among these are the optogenetic toolbox that has recently become available. Our access to this toolbox includes light induced neuronal excitation using channelrhodopsins (courtesy of the Deisseroth lab in Stanford and the Yizhar lab at Weizmann) as well as genetically encoded fluorescent labels for optical imaging of neuronal excitation (courtesy of the Cohen lab at Harvard).

Finally, it seems clear that evolutionary forces have shaped the architecture, connectivity as well as the input that the neurons grow with inside the brain, all of which conspire to create its tremendous computational power. Our most novel proposal is that biological computation in general is set apart from that of electronic computers by the fact that the ‘hardware’ (e.g. the cell or the organism) co-evolved with the ‘software’ (DNA or the brain respectively). As a result, our ‘software’ naturally tends to the well being of the associated ‘hardware’, while in the computer world that is not necessarily, if ever, the case. This insight leads us to propose experiments that will connect the neuronal network with the ‘real’ world, and allow the structure and function to co-evolve in a range of different environmental conditions.

### ***Natural models of computation in living neural networks***

Since the development of the electronic digital computer, which makes use of the digital abstraction [5] from analog electrical circuits there has been a close heuristic association between boolean logic and computation. However, developments in formal logic over the past century have been largely motivated by its applications to computation and the theory of programming languages that sometimes build upon and sometimes go beyond boolean logic to support additional forms of abstraction. In fact, the capacity to support abstraction mirroring various systems of formal logic has become a means to order and thereby assign value to programming languages [2]. Electronic devices that are capable of supporting such forms of abstraction provide a concrete physical instantiation of the ideas inherent to the logical systems they are designed to faithfully implement.

Neuronal logic devices (NLDs) represent an alternative physical modality to digital electronics for the purpose of performing computation. However, rather than attempting to directly parallel the history of the development of digital electronics via the digital abstraction, high-level descriptions of computation, such as the  $\lambda$ -calculus, serve as a specification of computation that is agnostic to the physical modality of implementation. If the specification that a neuronal computation device should be capable of implementing the  $\lambda$ -calculus, a natural first step toward this broad goal is to investigate simple neuronal systems that are capable of performing well-defined computations that require a capacity for first- or higher-order logic.

The first question to be addressed experimentally is whether a universal Turing machine can reliably be built out of central nervous system (CNS) neurons, which are essentially unreliable components. We have shown previously that complex devices such as a diode, an oscillator and an AND-gate can be engineered using CNS neurons grown on particular geometric configurations [4]. The production of a NOT-gate would complement this set and enable the construction of a universal Turing machine. The design we are proposing initially relies on the fact that individual NLDs can

be connected to form more complex structures. To form a NOT-gate, a delay line, oscillator AND gate and diode can be put together in a way that ensures that the output is continually excited (at the frequency of the oscillator) unless the Input is high. This relies on the fact that living neuronal networks have a latency period following their spike. The NOT gate can then be extended to build a NAND gate.

### A higher-order function to be implemented as a higher-order NLD

To design higher order computational devices, we will turn to the experimental techniques of microfluidics and of optogenetics [6, 7]. It is now possible to monitor with optical means the electrical activity of the network without any collateral damage caused by the fluorescent dye. This is done by incorporating genetically a fluorescent voltage-indicating protein into the neuron. It is furthermore possible to excite a region of the network optically by activating photosensitive channels that are genetically embedded as well. On top of this, the propagation velocity of a signal inside a one-dimensional neuronal network of the type we are using is constant, and can be reliably predicted. Thus it becomes possible to identify activity in one part of a device, and then excite another region co-incidentally with the arrival of the signal into that area. Different time delays, with the activation before, during or after the arrival of the synaptic input will allow the creation of several neuronal learning scenarios, and the comparison to learning in organisms with brains.

On the part of the theory, the  $\lambda$ -calculus notation is helpful in order to define any higher order function. Some of the notation may be implicitly familiar to users of imperative programming languages under the heading “anonymous functions”. We provide only an extremely informal set of examples necessary to explain our intended work; however, complete details can be found in [8]. We can define a standard binary boolean function like the “and” function as  $\lambda x.\lambda y.(\wedge x y)$ . Such a lambda expression can apparently be applied to any inputs; however the inputs to such a function are not necessarily restricted to booleans unless we infer that the standard logical operator “ $\wedge$ ” only accepts boolean arguments. Note that we have used the **prefix or Polish notation** for the  $\wedge$  operator, which in the perhaps more common infix notation is written with its arguments flanking the operator as  $x \wedge y$  to mean “ $x$  and  $y$ ”. We can provide explicit type annotations for the bound variables  $x$  and  $y$ , which indicate the types of the arguments

$$\lambda x : \text{bool}.\lambda y : \text{bool} . (\wedge x y) \quad (1)$$

We can now also provide a type annotation for this expression as a whole

$$[\lambda x : \text{bool}.\lambda y : \text{bool} . (\wedge x y)] : [\text{bool} \rightarrow \text{bool} \rightarrow \text{bool}] \quad (2)$$

Depending upon conventions with respect to **currying**, one could read the type annotation as “the function that takes two arguments each of type *bool* and returns a values of type *bool*” or “the function that takes an argument of type *bool* and returns a function that takes an argument of type *bool* and returns a value of type *bool*”. The first formulation may be easier to read, but the second is standard as a result of the way in which functional programming languages implement such functions.

Now we can imagine that if we wish to abstract from the particular binary boolean function implied by the  $\wedge$  operator, we need to introduce a functional variable, whose type will be explicitly denoted for concreteness despite the fact that it could be inferred, for which this operator can be substituted. Doing this results in the following second order function

$$[\lambda f : (\text{bool} \rightarrow \text{bool} \rightarrow \text{bool}).\lambda x : \text{bool}.\lambda y : \text{bool} . (f x y)] : [(\text{bool} \rightarrow \text{bool} \rightarrow \text{bool}) \rightarrow \text{bool} \rightarrow \text{bool} \rightarrow \text{bool}] \quad (3)$$

This function is intended to be read, in the less verbose uncurried form, as “the function that takes as its first argument (a function that takes two boolean values as arguments and returns a boolean value) and as its second and third arguments two boolean values and returns a boolean value”. It is perhaps remarkable that this simple abstraction is now capable of implementing any of the 16 possible boolean functions, provided that the proper binary boolean operator is submitted as the first argument to this function. Furthermore, the statement of this function in terms of a simply typed lambda calculus notation is agnostic to any physical implementation capable of realizing extensionally equivalent behavior.

To make this more concrete, we can very simply implement the above function in any functional programming language, or any programming language supporting the passing of function handles to other functions. An implementation in the programming language OCaml appears as follows

```
1 let hobf = fun (bf : ('a 'a bool)) (i1 : bool) (i2 : bool)
    -> bf i1 i2
```

Listing 1: a higher order boolean function

What is required in order to evaluate this function are corresponding implementations of boolean functions to be substituted either for  $f$  in the lambda calculus notation or for  $bf$  in terms of the corresponding OCaml implementation. For example we can implement the XOR function using pattern matching to define a truth table.

```
1 let xor p1 p2 = match (p1, p2)
2 with (false, false) -> false
3      | (false, true) -> true
4      | (true, false) -> true
      | (true, true) -> false
```

Listing 2: implementation of an XOR boolean operator

Other binary boolean functions are implemented in an analogous fashion. In order to evaluate `hobf` we then simply provide the name of a binary boolean function and two boolean values as

```
1 >hobf xor 0 0 = 0
>hobf xor 0 1 = 1
3 >hobf xor 1 0 = 1
>hobf xor 1 1 = 0
```

Listing 3: Example output of `hobf`

## Assessment of abstraction potential in NLDs

An important consideration is to state precisely some criterion for determining that a particular NLD has achieved potential for an explicit form of abstraction such as is indicated in the relationship between expressions 2 and 3. In abstracting the binary boolean operator  $\wedge$  to the functional variable  $f$ , which is the fundamental transformation enabling the derivation of 3 from 2, we imply that any physical implementation must take at least three rather than two inputs and the first of these must specify a particular binary boolean operator to apply to the latter two boolean input values. This roughly means that any system that at least partially implements the lambda expression or function specified in equation 3 and listing 1 respectively, must be capable of interpreting the concept of “selection from a set” whose size is determined by the subset of the 16 binary boolean operators that is already implemented in a lower-level form. Indeed another representation of equation 3 as a partial or total set function could be written as  $hobf : Hex \times Bool \times Bool \rightarrow Bool$  where we interpret

*Bool* and *Hex* as two and sixteen element sets respectively. This point of view makes clear that the capacity for selection from two element sets is already apparent in the binary boolean operator written as a set function  $\wedge : Bool \rightarrow Bool$ . The difference between these is that the system must implement the typing constraints necessary to distinguish a set that takes on sixteen possible values from one that takes on two. For example, in the application of the function *hobf* the following evaluate as expected given the definition:  $hobf \wedge 1\ 0 = 0$  or  $hobf \wedge 1\ 1 = 1$ . However, what is to be expected given inputs such as:  $hobf\ 1 \wedge 0$  or  $hobf\ 1\ 1 \wedge$ ? In these cases an output type intuitively associated to *error* is required to indicate that the realization of a system implementing equation 3 has indeed correctly implemented the necessary typing constraints to claim that the function has been realized. In the case of neuronal logic devices, this alternative output should differ in a measurable way from those associated to 1 and 0.

### From NLDs to the CNS and principles of cognition

As noted above, we associate the failure of a cultured neuronal network to produce a valid computation with the fact that it is grown out of context, out of its natural surroundings. By using geometric constraints, we have shown that we can coax some of the networks connections to go in the direction we engineer, and a modicum of computation is restored [4]. However, the idea that the ‘software’ is growing without the presence of the ‘hardware’ it is accustomed to, leads us to the idea of co-culturing a different set of neurons along with the regular hippocampal ones. The obvious choice for us is that of sensory neurons, since those are the subset of neurons that communicate between the body of the organism and its brain. In this way we hope to recreate an information channel that exists in the developing brain, allowing neurons to attain functional connectivity according to the inputs it receives. We are thus communicating with a number of groups that study pain, with the intention of extracting sensory neurons that synapse onto CNS neurons to convey the signals associated with stress and pain. This is a natural initial choice, since the preception of pain signals is fundamental to survival in many cases and thus neurons may have evolved to be robustly sensitive to such signals.

From the point of view that identifies the capacity for abstraction with computation, the investigation of neuronal logic devices capable of such function provides a framework in which various hypotheses from cognitive science could begin to be evaluated at the level of well-defined neural circuits. By attempting to isolate minimal implementation criteria, this approach may serve to complement, enable simpler explanations of, or identify paradoxical results derived from studies that treat whole brains and their associated sensory apparatus as their object of study [9, 10].

## Section 4: Role and expertise of the PIs

Both Profs. Aviv Bergman and Elisha Moses have significant experience working in collaborative teams of experimentalists and theorists. The PIs intend this project to involve strong interaction between experiment and theory.

Prof. Bergman will lead the theoretical component of this project. Prof. Bergman’s expertise spans a number of theoretical areas including artificial neural networks, dynamical systems, mathematical evolutionary biology, and systems biology. The Bergman lab will develop the computational platform and analytical tools to predict neuronal architecture-function relationships, which in turn will help guide the construction of environmental conditions to guide the development of natural neural networks capable of performing fundamental computational tasks.

Prof. Moses will lead the experimental component of the project, heading a laboratory that focuses on the growth and measurement of neuronal activity in networks grown from CNS neurons from rat and mouse brain. The lab has pioneered the design of complex neuronal logical devices,

has generated several new experimental paradigms combining nonlinear dynamics and statistical physics with biological physics, and has participated in the formation of novel theoretical models for engineered neuronal networks.

## Section 5: Educational involvement

In the Bergman lab, three Ph.D. students (Mr. Cameron Smith, Mr. Daniel Biro and Mr. Ximo Pechaun) will be involved in developing the theoretical aspect of this project in close interaction with the PI. In the Moses lab two students (1 Ph.D., Ms. Shani Stern and 1 M.Sc., to be hired) and a postdoctoral fellow (Dr. Yaron Penn) will be involved. Students and postdoctoral fellows in the Bergman and Moses labs will interact on a weekly basis via video conference. More extensive interaction will be fostered by an exchange program that we plan to engage in at crucial theory-experiment integration stages throughout the project for the purpose of increasing the level of coherence between the theoretical and experimental aspects of this work. Virtual interaction among all involved, including the public, will be fostered by an open online Wiki that will be used to organize and collaborate on this project and, more generally, support the movement for [Open Notebook Science](#). All computer code will be made open source in accordance with the [MIT license](#) and the codebase history will be available to the public free and in real-time on [github](#). Despite the fact that both labs already combine theory and experiment, they do so in very different ways and exposure to each of these models for combining theory and experiment will be crucial for these students as they move on to make very important decisions in their careers with respect to postdoctoral training and ultimately building labs of their own.

## References Cited

- [1] Ofer Feinerman and Elisha Moses. Transport of information along unidimensional layered networks of dissociated hippocampal neurons and implications for rate coding. *The Journal of neuroscience*, 26(17):4526–34, April 2006.
- [2] Harold Abelson and Gerald Jay Sussman. *Structure and Interpretation of Computer Programs - 2nd Edition (MIT Electrical Engineering and Computer Science)*. The MIT Press, 1996.
- [3] J. B. Tenenbaum, C. Kemp, T. L. Griffiths, and N. D. Goodman. How to Grow a Mind: Statistics, Structure, and Abstraction. *Science*, 331(6022):1279–1285, March 2011.
- [4] Ofer Feinerman, Assaf Rotem, and Elisha Moses. Reliable neuronal logic devices from patterned hippocampal cultures. *Nature Physics*, 4(12):967–973, October 2008.
- [5] Stephen A. Ward and Robert H. Halstead. *Computation Structures*. The MIT Press, 1989.
- [6] Ofer Yizhar, Lief E Fenno, Thomas J Davidson, Murtaza Mogri, and Karl Deisseroth. Optogenetics in neural systems. *Neuron*, 71(1):9–34, July 2011.
- [7] Joel M Kralj, Adam D Douglass, Daniel R Hochbaum, Dougal Maclaurin, and Adam E Cohen. Optical recording of action potentials in mammalian neurons using a microbial rhodopsin. *Nature methods*, 9(1):90–5, January 2012.
- [8] H.P. Barendregt. *The Lambda Calculus, Its Syntax and Semantics*. North Holland, 1985.
- [9] James L McClelland, Matthew M Botvinick, David C Noelle, David C Plaut, Timothy T Rogers, Mark S Seidenberg, and Linda B Smith. Letting structure emerge: connectionist and dynamical systems approaches to cognition. *Trends in cognitive sciences*, 14(8):348–56, August 2010.
- [10] Thomas L Griffiths, Nick Chater, Charles Kemp, Amy Perfors, and Joshua B Tenenbaum. Probabilistic models of cognition: exploring representations and inductive biases. *Trends in cognitive sciences*, 14(8):357–64, August 2010.

## Biographical Sketch: Your Name