
A Universal Abstraction for Hierarchical Hopfield Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Conceptualized as Associative Memory, Hopfield Networks (HNs) are powerful
2 models which describe neural network dynamics converging to a local minimum
3 of an energy function. HNs are conventionally described by a neural network with
4 two layers connected by a matrix of synaptic weights. However, it is not well
5 known that the Hopfield framework generalizes to systems in which many neuron
6 layers and synapses work together as a unified Hierarchical Associative Memory
7 (HAM) model: a single network described by memory retrieval dynamics (conver-
8 gence to a fixed point) and governed by a global energy function. In this work we
9 introduce a universal abstraction for HAMs using the building blocks of neuron
10 layers (nodes) and synapses (edges) connected within a hypergraph. We imple-
11 ment this abstraction as a software framework, written in JAX, whose autograd
12 feature removes the need to derive update rules for the complicated energy-based
13 dynamics. Our framework, called HAMUX (**HAM** User e**X**perience), enables any-
14 one to build and train hierarchical HNs using familiar operations like convolutions
15 and attention alongside activation functions like Softmaxes, ReLUs, and Layer-
16 Norms. HAMUX is a powerful tool to study HNs at scale, something that has
17 never been possible before. We believe that HAMUX lays the groundwork for a
18 new type of AI framework built around dynamical systems and energy-based as-
19 sociative memories.

20 1 Introduction

21 Non-linear ordinary differential equations (ODEs) are extensively used in modern AI architectures,
22 leading to impressive results. An important subset of general ODEs are systems with an under-
23 lying global Lyapunov function, often called an energy function, which decreases in time as the
24 non-linear dynamical system approaches the fixed point state. Paradigmatic examples of such sys-
25 tems are **Hopfield Networks** (HNs), introduced in 1982 by John Hopfield as models of associative
26 memory retrieval [1, 2]. The core idea is that the non-linear (discrete or continuous) dynamical
27 system, governed by the energy function, is designed to have multiple fixed points (memories) with
28 substantial basins of attraction around them. Given an initial prompt (a query), the network picks
29 one of the basins of attraction and follows the gradient of the energy function to “retrieve” the fixed
30 point state stored at the bottom of that basin. There has been a resurgence of interest in HNs in
31 the past few years thanks to novel results pertaining to their memory storage capacity, relationship
32 to transformer’s attention, and possible ways of integrating these ideas in a wide variety of deep
33 learning architectures [3, 4, 5, 6, 7, 8].

34 Hopfield Networks are conventionally seen as shallow, two-layer systems, e.g., [9, 6, 7, 10]. How-
35 ever, this restriction of shallowness is one imposed by historical usage and not one of the energy-
36 based equations themselves. It turns out that the same energy rules that govern the simple two-layer
37 system generalize to HNs composed of any number of layers and synapses, an architecture called

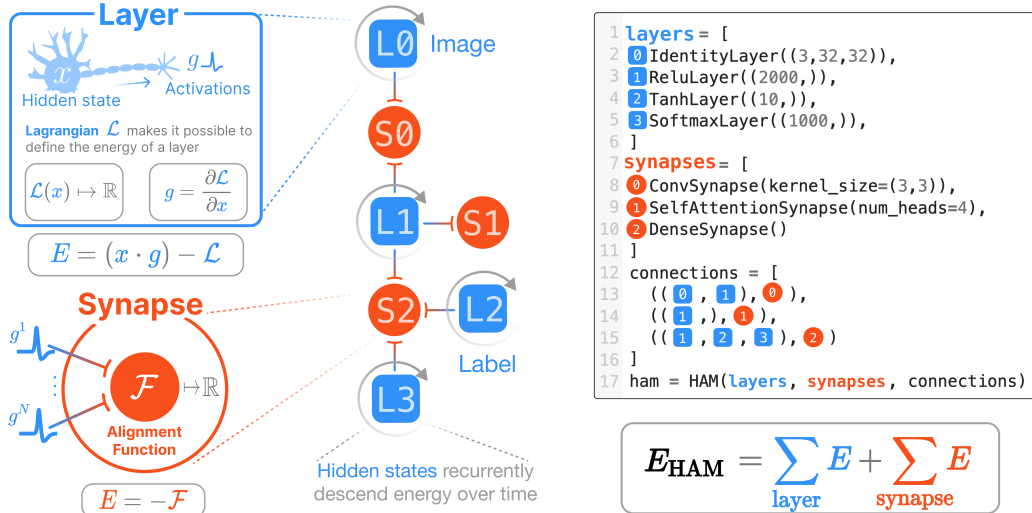


Figure 1: **Left:** the fundamental building blocks of our abstraction. Each *Layer* is given a Lagrangian function \mathcal{L} that fully defines both its activations g and energy. Each *Synapse* has energy dependent only on a learnable alignment function that converts the activations of connected layers into a scalar. **Center:** an example HAM composed of four layers and three synapses. The system’s energy function is the sum of the energies of its components. Layer states are recurrent in time. We use 0-based indexing to align the graphical representation with the matching Python code. **Right:** pseudocode illustrating how illustrating how a few lines of code can build a HAM. Each layer (lines 2-5) is assigned a Lagrangian (denoted by its familiar activation names e.g., identity, relu, tanh, softmax) and shape. Synapses (lines 8-10) are modeled after common network operations (e.g., convolution, attention, dense), which are all implemented in HAMUX and whose parameter shapes are fully defined by the connected layers. The hypergraph definition (lines 13-15) assembles the graph, one line per synapse.

38 the **Hierarchical Associative Memory** (HAM) [11]. The conventional HN is then a special limiting
 39 case of the HAM with two layers and one synapse.

40 The generalizability of the HN to the HAM is not common knowledge to the research community,
 41 and as such the behavior of the HAM as a “Hierarchical Hopfield Network” has not been well
 42 characterized or understood. In this work we propose a powerful abstraction and accompanying
 43 software framework called HAMUX that removes critical barriers that stand in the way of applying
 44 these networks at scale:

45 **Barrier 1: There is no standardized terminology for the energy fundamentals of a HAM.**
 46 Our framework proposes an abstraction that fully captures the behavior of any HN while being
 47 *modular* (i.e., it is easy to assemble deep HAMs connected to any number of signals), *generalizable*
 48 (i.e., it is possible to quickly propose novel operations and activation functions that satisfy energy
 49 constraints), and reminiscent of the *biological* inspiration for the original HN.

50 **Barrier 2: The energy of a HAM grows increasingly complex with the number of components and connections.**
 51 Our framework uses modern AI tooling (JAX [12]) to automatically calculate the gradient of the energy
 52 function for any given state and parameters, removing the need to manually derive complicated update
 53 rules.

54 **Barrier 3: It can be challenging to conceptualize HAM architectures as modern machine learning pipelines.**
 55 It requires the shift of perspective from the standard “classification” setup (take an input, pass it
 56 through the network, and return an output), to an “association” setup, where labels are just another
 57 attribute of a data point (data and labels are equally important and treated the same). With the
 58 examples released with HAMUX, we show how one can apply HAMs to conventional “classification”
 59 pipelines while reusing existing tools in the Deep Learning ecosystem to create custom energy
 60 components.

61 Our framework, though still in its infancy, aims to integrate HAMs into modern Deep Learning. All
 62 experiments in this paper were conducted using our framework.

63 2 The Abstraction

64 Hopfield Networks are recurrent networks whose behavior at time t is completely defined as a
 65 function of the neuron states. Our abstraction pivots on understanding how the individual ener-
 66 gies of *neuron layers* and *synapses* (i.e., the “energy building blocks” of our abstraction) oper-
 67 ate within the constraints of a *hypergraph*. Concretely, our HAM framework with N layers and
 68 K synapses is fully defined by a list of all **neuron layers** $\mathcal{X} = \{\mathcal{X}^1, \mathcal{X}^2, \dots, \mathcal{X}^N\}$, a list of all
 69 **synapses** $\mathcal{S} = \{\mathcal{S}^1, \mathcal{S}^2, \dots, \mathcal{S}^K\}$, and a list of connections specifying the **connection hypergraph**
 70 $\mathcal{G} = \{\mathcal{G}^1, \mathcal{G}^2, \dots, \mathcal{G}^K\}$. We emphasize that a HAM uses a hypergraph and not a normal graph: i.e.,
 71 a synapse (edge) can operate on an arbitrary number of layers (nodes). In practice, we represent \mathcal{G}^κ
 72 as the collection of integers $(\{\alpha, \beta, \dots\}, \kappa)$ to specify that synapse \mathcal{S}^κ operates on the activations of
 73 layers $\{\mathcal{X}^\alpha, \mathcal{X}^\beta, \dots\}$. Figure 1 summarizes our abstraction. The system’s total energy is defined as
 74 a sum of the energies of its components.

$$E_{\text{total}} = \sum_{\alpha=1}^N E_{\text{layer}}^\alpha + \sum_{\kappa=1}^K E_{\text{synapse}}^\kappa \quad (1)$$

75 2.1 Neuron Layers

76 A neuron layer \mathcal{X}^α is an assembly of D^α neurons each with a *hidden state* x_i^α governed by a
 77 *Lagrangian function* $\mathcal{L}^\alpha : \mathbb{R}^{D^\alpha} \mapsto \mathbb{R}$. All neurons share a scalar *time constant* τ^α that governs how
 78 quickly x_i^α will evolve in time. Each neuron can additionally have a *resting state* I_i^α that we call the
 79 bias in conventional deep learning. The choice of the Lagrangian \mathcal{L}^α fully defines the activation or
 80 the *gain function* $g^\alpha : \mathbb{R}^{D^\alpha} \mapsto \mathbb{R}^{D^\alpha}$ for a given neuron layer as $g^\alpha := \frac{\partial \mathcal{L}^\alpha}{\partial x^\alpha}$.

81 The most important rule of the neuron layer is the following: *the hidden state x_i^α is completely*
 82 *invisible to the rest of the network at any point in time*. The ONLY way a neuron can influence the
 83 rest of the network is via its activation g_i^α (i.e., all synapses must operate on activations).

84 The energy of neuron layer \mathcal{X}^α can be computed from its state x^α as follows:

$$E_{\text{layer}}^\alpha = \sum_i (x_i^\alpha - I_i^\alpha) g_i^\alpha - \mathcal{L}_\alpha(x^\alpha) \quad (2)$$

85 In practice, x_i^α need not be a scalar. For example, a layer of shape $\mathbb{R}^{D^\alpha \times H \times W}$ has D^α neurons
 86 whose states x_i^α and activations g_i^α are image patches in $\mathbb{R}^{H \times W}$.

87 2.2 Synapses

88 A *synapse* \mathcal{S}^κ is a parameterized *alignment function* \mathcal{F}^κ that transforms the activations $\{g^\alpha, g^\beta, \dots\}$
 89 of one or more layers $\{\mathcal{X}^\alpha, \mathcal{X}^\beta, \dots\}$ into a meaningful scalar that represents the alignment of those
 90 layers. This is a novel construction of our abstraction as no previous work has considered anything
 91 beyond pairwise synapses. A synapse’s energy is defined as the negative of its alignment function.

$$E_{\text{synapse}}^\kappa = -\mathcal{F}^\kappa, \quad \text{where } \mathcal{F}^\kappa(g^\alpha, g^\beta, \dots) \mapsto \mathbb{R}. \quad (3)$$

92 2.3 The Update Rule

93 The update rule for x_i^α is the direct consequence of differentiating the total energy in Eq. 1. Incoming
 94 signals $\frac{\partial \mathcal{F}^\kappa}{\partial g_i^\alpha}$ into x_i^α can only come from connected synapses. See the derivations in [7, 11].

$$\tau^\alpha \frac{dx_i^\alpha}{dt} = -\frac{\partial E_{\text{total}}}{\partial g_i^\alpha} = \sum_{\kappa=1}^K \frac{\partial \mathcal{F}^\kappa}{\partial g_i^\alpha} + I_i^\alpha - x_i^\alpha \quad (4)$$

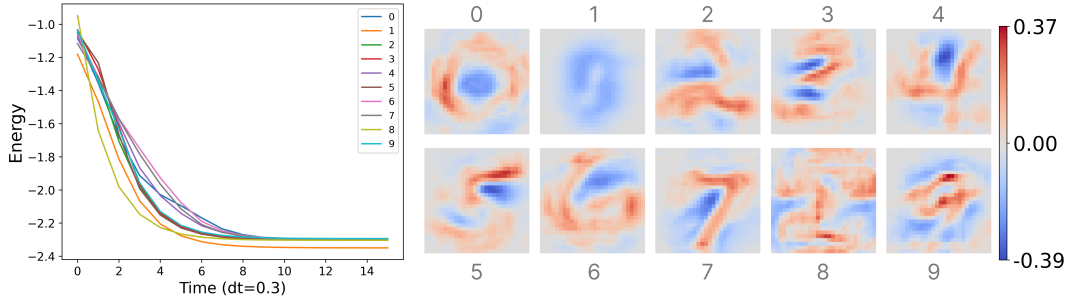


Figure 2: A HAM was trained on the MNIST classification task. After training, the image neurons of the trained HAM were initialized at random while the label neurons were clamped to one-hot encodings for each digit. The image neurons were then allowed to update their states so that the energy decreases in time (left) to the fixed point of the dynamics (in our experiments after 15 time steps). These fixed points are shown on the right as the $\tanh(\cdot)$ activations of the image layer. We use the HN (*SoftMax*) model from Table 1.

Table 1: Classification results reported as the best top-1 accuracy on the validation set at 100 (600) epochs.

Model	Top-1 Val %		# of Params (M)	
	MNIST	CIFAR10	MNIST	CIFAR10
HN (<i>ReLU</i>)	98.74 (99.03)	58.82 (66.23)	0.79	18.49
HN (<i>SoftMax</i>)	97.32 (97.99)	53.17 (58.89)	0.79	18.49
Conv HAM (<i>max pool</i>)	99.29 (99.51)	82.08 (86.39)	0.60	1.78

95 3 Classification and Observing the Dynamics

96 The problem of classification can be formulated in a HAM as an association problem: what one-hot
 97 encoded label is associated with a given collection of pixels? In this section we train several different
 98 architectures using HAMUX: the traditional HN in both the classical and modern paradigm where we
 99 concatenate flattened pixels to labels; and a novel HAM where we stack convolutional and pooling
 100 operations. All architectures behave under the global energy function and we do not use additional
 101 encoders, decoders, or classification heads. See Appendix A for technical details and discussion.
 102 The accuracy of our HAMs, shown in Table 1, stands in line with the accuracy of corresponding
 103 feedforward models on similar tasks.

104 Every system built with HAMUX is an associative memory. This means that training a HAM to
 105 “classify” images in one direction allows us to utilize the system in reverse — we can retrieve the
 106 memory most associated with a label by clamping the labels to a desired value over the dynamics.
 107 See Figure 2 for energy dynamics of a trained model with different clamped MNIST labels and
 108 Appendix B for details.

109 4 Conclusion

110 In this work, we have proposed a universal abstraction to describe the energy of general Hopfield
 111 Networks. Our abstraction proves particularly powerful for network design, allowing anyone to
 112 modularly construct deep HNs (HAMs) that can be applied to traditional machine learning tasks.
 113 At the same time, our abstraction generalizes the fundamental operations of synapses, giving us
 114 freedom to implement energy-constrained versions of convolutions, pooling, and even attention. We
 115 package our abstraction into a software framework called HAMUX that makes it trivial to implement
 116 complex HAMs. We believe that HAMUX gives HNs the representational power of modern Deep
 117 Learning by providing a framework that integrates modern AI tooling and advances into a regime
 118 that is governed by energy and whose cardinal function is association rather than prediction.

References

- 119
- 120 [1] J J Hopfield. Neural networks and physical systems with emergent collective computational
121 abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, April 1982.
- 122 [2] J. J. Hopfield. Neurons with graded response have collective computational properties like
123 those of two-state neurons. *Proceedings of the National Academy of Sciences of the United*
124 *States of America*, 81(10):3088–3092, May 1984.
- 125 [3] Dmitry Krotov and John J Hopfield. Dense associative memory for pattern recognition. *Ad-*
126 *vances in neural information processing systems*, 29, 2016.
- 127 [4] Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Upgang, and Franck Vermet. On a
128 Model of Associative Memory with Huge Storage Capacity. *Journal of Statistical Physics*,
129 168(2):288–299, July 2017.
- 130 [5] Dmitry Krotov and John Hopfield. Dense Associative Memory Is Robust to Adversarial Inputs.
131 *Neural Computation*, 30(12):3151–3167, December 2018.
- 132 [6] Hubert Ramsauer, Bernhard Schöfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas
133 Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield
134 networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- 135 [7] Dmitry Krotov and John J Hopfield. Large associative memory problem in neurobiology and
136 machine learning. In *International Conference on Learning Representations*, 2020.
- 137 [8] Michael Widrich, Bernhard Schöfl, Milena Pavlović, Hubert Ramsauer, Lukas Gruber, Markus
138 Holzleitner, Johannes Brandstetter, Geir Kjetil Sandve, Victor Greiff, Sepp Hochreiter, et al.
139 Modern hopfield networks and attention for immune repertoire classification. *Advances in*
140 *Neural Information Processing Systems*, 33:18832–18845, 2020.
- 141 [9] Qun Liu and Supratik Mukhopadhyay. Unsupervised learning using pretrained cnn and asso-
142 ciative memory bank. In *2018 International Joint Conference on Neural Networks (IJCNN)*,
143 pages 01–08. IEEE, 2018.
- 144 [10] Beren Millidge, Tommaso Salvatori, Yuhang Song, Thomas Lukasiewicz, and Rafal Bogacz.
145 Universal Hopfield Networks: A General Framework for Single-Shot Associative Memory
146 Models, June 2022.
- 147 [11] Dmitry Krotov. Hierarchical associative memory. *arXiv preprint arXiv:2107.06446*, 2021.
- 148 [12] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
149 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
150 Zhang. JAX: Composable transformations of Python+NumPy programs, 2018.
- 151 [13] Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter Buchlovsky,
152 David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Claudio Fantacci, Jonathan God-
153 win, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel, Shaobo Hou, Steven Kaptur-
154 owski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch, Lena Martens, Hamza
155 Merzic, Vladimir Mikulik, Tamara Norman, John Quan, George Papamakarios, Roman Ring,
156 Francisco Ruiz, Alvaro Sanchez, Rosalia Schneider, Eren Sezener, Stephen Spencer, Srivatsan
157 Srinivasan, Luyu Wang, Wojciech Stokowiec, and Fabio Viola. The DeepMind JAX Ecosys-
158 tem, 2020.
- 159 [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*
160 *preprint arXiv:1412.6980*, 2014.
- 161 [15] Ross Wightman. Pytorch image models. [https://github.com/rwightman/
162 pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 163 [16] Maxence Ernout, Julie Grollier, Damien Querlioz, Yoshua Bengio, and Benjamin Scellier.
164 Updates of equilibrium prop match gradients of backprop through time in an rnn with static
165 input. *Advances in neural information processing systems*, 32, 2019.
- 166 [17] Cristian Garcia. Treex, 2021.

Table 2: TIMM Data Augmentation configuration for the experiments

	MNIST	CIFAR
re_prob	0.1	0.2
hflip	0.0	0.5
vflip	0.0	0.0
scale	(0.9, 1.0)	(0.2, 1.0)
color_jitter	0.4	0.5
auto_augment	None	None
ratio	(0.75., 1.33)	(0.75, 1.33)

167 A Classification: Training and Architectures

168 We now provide the technical details for the variety of architectures that were implemented using
 169 HAMUX in section 3. All models were trained for 600 epochs using the ADAM optimizer with
 170 weight decay (using optax defaults: b1=0.9, b2=0.99, weight decay=1e-4 [13]) at a constant learning
 171 rate of 0.001 and batch size of 400. For consistency, we do not optimize hyperparameters for any
 172 particular architecture. We train the networks using simple backpropagation through time using the
 173 ADAM optimizer [14]. We load MNIST and CIFAR images with minor data augmentations from
 174 the popular `timm` [15] library. All experiments were run on an A100 GPU.

175 We provide example code for MNIST architectures; CIFAR10 versions of the architecture are almost
 176 identical, though they require different layer shapes for all layers related to images and patches.
 177 Depending on the depth of the HAM, we vary the number of recurrent timesteps along with our step
 178 size Δt through time. This is to ensure that pixel information has enough time to propagate to a
 179 potentially distant label layer. All architectures and their training scripts will be released as demo
 180 models along with the framework.

181 A.1 Shallow Hopfield Network Configurations

182 We begin by implementing two different HNs where a visible layer, the concatenation of vectorized
 183 pixels and one-hot labels, is connected to a single hidden layer via a synaptic weight matrix. The
 184 difference in these two architectures is in the choice of activation function. The `relu` activation func-
 185 tion defines the model as a continuous, recurrently applied Classical Hopfield Network (CHN) [2],
 186 whereas `softmax` activation function defines a Dense Associative Memory (DAM), also known as
 187 the Modern Hopfield Network (MHN) [3, 4, 6, 7]. The latter has been shown to have a higher storage
 188 capacity than the former; however, it is our experience in writing this paper that they are also harder
 189 to train using backpropagation-through-time. For this reason, in our HN (*Softmax*) architecture we
 190 additionally normalize each *memory* in the weights to have L2-norm equal to 1. Specifically, given
 191 our synaptic matrix connecting 1000 hidden units to 794 visible units (784 pixels + 10 labels), we
 192 enforce that the matrix consists of 1000 unit vectors each of dimension 794. This promotes diversi-
 193 fying of the weights during training and stability of the dynamics during longer inference runs.

194 To implement this architecture using HAMUX we define an image layer with the `tanh` activa-
 195 tion function that is connected to a label layer with the `softmax` activation function. Function-
 196 ally, the states of these two layers are concatenated together and their union is referred to as
 197 a single visible layer. The second “layer” of this two-layer HN is hidden inside the synapse
 198 (`DenseMatrixSynapseWithHiddenLayer`). So why do we put this layer within the synapse? In
 199 our abstraction, the states of each layer must be propagated through time — that is, a layer’s state
 200 $x^\alpha(t)$ should depend on no other layer states at the same time t . The original definition of the HN
 201 defines the state of the hidden layer at time t as a function of the state of the visible layer at time t
 202 (which is equivalent to taking the limit as τ of this layer approaches 0). By including the Lagrangian
 203 function as a hidden layer within a synapse we propose a workaround that is mathematically consis-
 204 tent with the original definition of the HN (that is, states are independent of other states at time
 205 t).

```
layers = [  
    TanhLayer((28,28,1)),  
    SoftmaxLayer((10,)),  
]  
synapses = [  
    DenseMatrixSynapseWithHiddenLayer(1000, hidden_lagrangian=LRelu()),  
]  
connections = [  
    ((0, 1), 0),  
]  
ham = HAM(layers, synapses, connections)
```

Code 1: Shallow HN with ReLU, MNIST

206 **A.2 Convolutional and Pooling HAMS**

207 Dense matrix operations are only one kind of alignment function that describe layer-layer relationships. We introduce an architecture that additionally uses convolutions and pooling. We must first
208 calculate the shape of each layer that serves as the output of a convolution or pooling operation.
209 With this we can describe our 5-layer HAM as follows (layer 5 is hidden within our last synapse):
210

```

layers = [
  TanhLayer((28,28,1), tau=1.0),
  TanhLayer((7,7,64), tau=1.0),
  TanhLayer((2,2,128), tau=1.0),
  SoftmaxLayer((10,), tau=1.0),
]
synapses = [
  ConvSynapseWithPool(
    (4, 4),
    strides=(2, 2),
    padding=(2, 2),
    pool_window=(2, 2),
    pool_stride=(2, 2),
    pool_type="max",
  ),
  ConvSynapseWithPool(
    (3, 3),
    strides=(1, 1),
    padding=(0, 0),
    pool_window=(2, 2),
    pool_stride=(2, 2),
    pool_type="max",
  ),
  DenseMatrixSynapseWithHiddenLayer(1000, hidden_lagrangian=LRelu()),
]
connections = [
  ((0, 1), 0),
  ((1, 2), 1),
  ((2, 3), 2)
]

```

Code 2: Convolutional HAM with MaxPooling, MNIST

211 A.2.1 On the consequences of Max-Pooling

212 We make several architectural choices in this paper, not all of which have the most natural corre-
 213 spondence to energy. In particular, consider the max pooling operation which is implemented as part
 214 of our convolutional synapse. For each patch on which this operation is applied it will discard infor-
 215 mation from every element but one. This makes it unclear what the signal should be in the gradient
 216 of that synapse’s energy for all non-maximum elements in the patch. The JAX autograd system (and
 217 hence HAMUX) uses the definition of max-pooling consistent with that proposed by [16] where the
 218 gradients of all non-maximum elements of the patch are zero.

219 If we consider the operation of max-pooling to be a competitive operation within a single neuron
 220 where the “winner takes all,” we can easily consider a softer version of the maximum on each patch,
 221 e.g., “softmax pooling.” We leave this for future work.

222 A.3 Example forward pass

223 Our HAMs are all fully dynamic systems through time, so it can be unnatural to consider them as
 224 prediction engines. Here we provide the forward pass that we used for our classification pipelines.


```

import jax.numpy as jnp
import jax.tree_util as jtu

def simple_fwd(ham, x, depth, dt):
    """A simple version of the forward function for classification.

    Image layers are `layer[0]` and labels are `layer[-1]`
    """
    # Initialize hidden states given our data
    xs = ham.init_states(x.shape[0])
    xs[0] = jnp.array(x)

    # Masks allow us to clamp our visible data over time
    masks = jtu.tree_map(lambda x: jnp.ones_like(x), xs)
    masks[0] = jnp.zeros_like(masks[0])

    for i in range(depth):
        # Calculate update directions
        updates = ham.vupdates(xs)

        # Simple step down the energy function
        xs = ham.step(
            xs, updates, dt=dt, masks=masks
        )

    # Label layer has the softmax activation function
    # Use this to return probabilities
    return ham.layers[-1].g(xs[-1])

```

Code 3: An example of the forward function for classification

225 B Observing Dynamics

226 Using our trained HN-Softmax model from our classification experiments, we clamp the labels of
 227 our system to a desired value and let the dynamics of the system evolve around these signals. We
 228 display the images as our HAM sees them: as the tanh activations of our pixels (hence the negative
 229 values in Figure 2). To extract sharper memories (i.e., to prevent our system from choosing an
 230 incoherent superposition of memories at the limit of the dynamics), we decrease the temperature of
 231 the Softmax in our hidden layer by a factor of 10.

232 C About HAMUX: Software Details

233 The HAMUX software used for this paper is undergoing rapid API and tooling changes in an effort
 234 to incorporate more of the expected functionality present in a Machine Learning library. At the time
 235 of this writing, all components in HAMUX have been built using the excellent but unpopular Treex
 236 library [17], which we found to be the most robust and simple ML framework for building the unique
 237 constraints required HAMUX components. As part of the JAX ecosystem, any traditional tools for
 238 working with deep networks in JAX also work with HAMUX (e.g., it would be easy to consider
 239 alternative optimization procedures implemented in Optax [13] to descend the energy function).
 240 The software will be released prior to this workshop date.

241 D Examples of Lagrangian Functions

242 The choice of the Lagrangian function L_α for layer \mathcal{X}^α fully defines the activation or the *gain*
 243 *function* of that layer $g^\alpha : \mathbb{R}^{D_\alpha} \mapsto \mathbb{R}^{D_\alpha}$ as $g^\alpha := \frac{\partial L_\alpha}{\partial x^\alpha}$ for a given neuron layer. These functions
 244 introduce non-linearities into our HAM. Sometimes these functions operate elementwise (e.g., ReLU,
 245 GeLU) whereas other activation functions include normalization effects that scale an individual

246 neuron's state x_i^α given the states of other neurons x^α in the same layer (e.g., SoftMax, LayerNorm).
 247 For behaved dynamics, we must choose a Lagrangian that is both convex and differentiable.

248 The following common activation functions in use today have easy parallels in the Lagrangian world.

249 **ReLU** The Lagrangian of the ReLU is

$$\mathcal{L} = \frac{1}{2} \sum_i \max(x_i, 0)^2,$$

250 where

$$g_i = \frac{\partial \mathcal{L}}{\partial x_i} = \max(x_i, 0).$$

251 **SoftMax** The Lagrangian of the SoftMax is also known as the LogSumExp. It proves benefi-
 252 cial to consider the case of a softmax with a (potentially learnable) inverse temperature scalar
 253 β (e.g., Transformer attention defaults to $\frac{1}{\sqrt{D_{\text{key}}}}$). The Lagrangian of this operation can be
 254 expressed as

$$\mathcal{L} = \frac{1}{\beta} \log \sum_i \exp(\beta x_i),$$

255 where

$$g_i = \frac{\partial \mathcal{L}}{\partial x_i} = \frac{\exp(\beta x_i)}{\sum_j \exp(\beta x_j)}.$$

256 **Identity** The identity activation occurs with the Lagrangian

$$\mathcal{L} = \frac{1}{2} \sum_i x_i^2,$$

257 where

$$g_i = \frac{\partial \mathcal{L}}{\partial x_i} = x_i.$$

258 **LayerNorm** The LayerNorm activation modifies each neuron layer to have mean 0 and standard
 259 deviation 1, while optionally learning a scale γ on the standard deviation and a shift δ to the
 260 mean $\bar{x} = \frac{1}{D} \sum_{k=1}^D x_k$. The corresponding Lagrangian is

$$L = D\gamma \sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon} + \sum_j \delta_j x_j,$$

261 where

$$g_i = \gamma \frac{x_i - \bar{x}}{\sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon}} + \delta_i.$$