

# Biologically Realistic Computational Primitives of Neocortex Implemented on Neuromorphic Hardware Improve Vision Transformer Performance

Asim Iqbal<sup>\*1</sup>, Hassan Mahmood<sup>1</sup>, Greg J. Stuart<sup>2,3</sup>, Gord Fishell<sup>\*4,5</sup>, and Suraj Honnuraiah<sup>†\*2,4,5,6</sup>

<sup>1</sup>Tibbling Technologies, Redmond, WA, USA

<sup>2</sup>Eccles Institute of Neuroscience, John Curtin School of Medical Research, The Australian National University, Canberra, AU

<sup>3</sup>Department of Physiology, Monash University, Melbourne, Australia

<sup>4</sup>Harvard Medical School, Department of Neurobiology, Boston, MA, USA

<sup>5</sup>Stanley Center for Psychiatric Research, Broad Institute, Cambridge, MA, USA

<sup>6</sup>Institute of Neuroinformatics (INI), ETH Zurich and University of Zurich, Zurich, Switzerland

<sup>†</sup>Senior author and Lead contact

<sup>\*</sup>Corresponding authors: Asim Iqbal (asim@tibbtech.com), Gord Fishell (gordon.fishell@hms.harvard.edu) and Suraj Honnuraiah (suraj.honnuraiah@hms.harvard.edu)

## ORCID iDs:

Asim Iqbal: 0000-0003-2174-4554

Greg J. Stuart: 0000-0001-9395-2219

Gord Fishell: 0000-0002-9640-9278

Suraj Honnuraiah: 0009-0007-2867-2211

October 6, 2024

Understanding the computational principles of the brain and replicating them on neuromorphic hardware and modern deep learning architectures is crucial for advancing neuro-inspired AI (NeuroAI). Here, we develop an experimentally-constrained biophysical network model of neocortical circuit motifs, focusing on layers 2-3 of the primary visual cortex (V1). We investigate the role of four major cortical interneuron classes in a competitive-cooperative computational primitive and validate these circuit motifs implemented soft winner-take-all (sWTA) computation for gain modulation, signal restoration, and context-dependent multistability. Using a novel parameter mapping technique, we configured IBM's TrueNorth (TN) chip to implement sWTA computations, mirroring biological neural dynamics. Retrospectively, we observed a strong correspondence between the biophysical model and the TN hardware parameters, particularly in the roles of four key inhibitory neuron classes: Parvalbumin (feedforward inhibition), Somatostatin (feedback inhibition), VIP (disinhibition), and LAMP5 (gain nor-

malization). Moreover, the sparse coupling of this sWTA motif was also able to simulate a two-state neural state machine on the TN chip, replicating working memory dynamics essential for cognitive tasks. Additionally, integrating the sWTA computation as a pre-processing layer in the Vision Transformer (ViT) enhanced its performance on the MNIST digit classification task, demonstrating improved generalization to previously unseen data and suggesting a mechanism akin to zero-shot learning. Our approach provides a framework for translating brain-inspired computations to neuromorphic hardware, with potential applications on platforms like Intel's Loihi2 and IBM's Northpole. By integrating biophysically accurate models with neuromorphic hardware and advanced machine learning techniques, we offer a comprehensive roadmap for embedding neural computation into NeuroAI systems.

**Keywords:** Biophysics of neocortical computation, IBM's TrueNorth neuromorphic chip, Winner-take-all, Brain-inspired computing, Vision Transformers, Domain generalization, and NeuroAI.

## Introduction

Recent advances in machine learning and computational neuroscience have significantly accelerated the progress toward the development of synthetic cognitive agents with artificial general intelligence (AGI). Vision transformers (Dosovitskiy et al. [2020]) and natural language models (Shanahan et al. [2023]) have achieved notable success in image recognition and natural language processing. However, despite surpassing human performance in specific tasks like chess (Campbell et al. [2002]) and Go (Silver et al. [2017]), AI systems still encounter significant challenges when learning in novel environments. These systems require substantially more computational resources and annotated data than biological brains. This disparity may also arise from fundamental differences in how artificial and biological neural networks process information. In this work, we investigate the potential of reverse-engineering the brain's computational principles and integrating them into AI systems. This exploration aligns with the core tenet of the NeuroAI approach (Zador et al. [2023]), aiming to bridge the existing gap between artificial and biological intelligence.

The execution of cognitive behavior in the brain relies on the ability to select actions based on external stimuli and context (Dayan [2008]). In animals, the learning of state-dependent sensorimotor mappings (Asaad et al. [2000], Banerjee et al. [2020], Xu et al. [2022], Condylis et al. [2020]) is primarily mediated by the neocortex, which facilitates cognition through computations enabled through its modular, laminar microcircuits. These microcircuits consist of excitatory and inhibitory neurons, including four major inhibitory classes — parvalbumin (PV), somatostatin (SST), vasoactive intestinal peptide (VIP), and Lamp5 (Rudy et al. [2011], Tremblay et al. [2016]). These interneurons play a crucial role in regulating state-dependent computations, performing tasks such as arithmetic, logical operations, timing, and gain modulation (Fishell and Kepcs [2020], Kepcs and Fishell [2014], Ferguson and Cardin [2020], Niell and Scanziani [2021]). Importantly, these four inhibitory neuron classes are conserved across cortical regions and species (Pfeffer et al. [2013], Campagnola et al. [2022]), indicating that their computational logic is generalizable for diverse high-order tasks, including motor execution and working memory.

Several key candidate computational princi-

ples have been proposed to elucidate neocortical function, including normalization (Carandini and Heeger [2012]), dynamic field theory (Schöner and Spencer [2016]), attractor networks (Vyas et al. [2020]), predictive coding (Keller and Mrsic-Flogel [2018]), Bayesian inference (Bastos et al. [2012]) and winner-take-all (WTA) computations (Douglas and Martin [2007]). However, direct evidence at the level of microcircuit and biological hardware implementation remains limited. Among these, the WTA mechanism is amenable to neocortical architecture and combines key elements of these various computational approaches. By employing competitive-cooperative dynamics, the WTA mechanism facilitates selective amplification and noise minimization, thus enhancing signal restoration (Douglas and Martin [2007]). These characteristics resemble signal processing in both primary sensory and motor cortices, where superficial pyramidal neurons receive sparse and weak thalamic inputs that require amplification to extract relevant information (Balcioglu et al. [2023], Lien and Scanziani [2018], Bopp et al. [2017], Binzegger et al. [2004]). Consequently, the WTA mechanism may be a fundamental computational strategy employed by cortical circuits and represent a ubiquitous computational strategy implemented by the neocortex.

In addition, WTA models show considerable promise for neuromorphic hardware (Mead, 1990; 2023), especially in energy-efficient, real-time processing (Chicca et al. [2014], Qiao et al. [2015], Indiveri and Sandamirskaya [2019]). To execute such computations in silico, IBM's TrueNorth (TN) chip offers a tractable platform for integrating brain-inspired principles. For example, it features a reconfigurable, asynchronous, multi-core digital architecture optimized for real-time, ultra-low-power, event-driven processing with physical neurons (Modha et al. [2023], Merolla et al. [2014], Neckar et al. [2018]). As a result, TN is especially well-suited for implementing brain-like computations (Indiveri and Sandamirskaya [2019]). While prior research has focused on leveraging statistical relationships among neuronal populations to emulate biological circuits on TrueNorth hardware (Imam [2021]), developing generalizable techniques for integrating diverse biophysical and theoretical models remains an open challenge. In this work, we demonstrate that by employing biophysically realistic computational principles, parameters of IBM TrueNorth - such as thresholds, leak rates, and crossbar weights, can be modeled to reflect those found in cortical

microcircuits. By utilizing this approach, TrueNorth (TN) hardware can be programmed to perform computations analogous to those observed in simplified, biologically realistic V1 cortical circuit motifs, that may potentially underlie key V1 functions such as orientation and direction tuning (Rossi et al. [2020], Niell and Stryker [2008], Douglas and Martin [2007], Hubel and Wiesel [1962]).

Our goal here was not to build an exhaustive model of V1, such as that outlined in Billeh et al. [2020], but rather to design a simplified, generalizable circuit motif that validates the core computational principles utilized in cortical processing. The retrospective analysis confirmed that the optimal parameters for configuring TN hardware to display sWTA dynamics closely aligned with the primary functions of different interneuron classes. Furthermore, our findings demonstrated that hardware-optimized abstractions could effectively replicate biological circuits. Finally, to test the functionality of this approach, we investigated whether integrating this hardware-constrained sWTA computation could be utilized to implement a neural state machine for working memory or to enhance the performance of state-of-the-art deep learning models such as Vision Transformers (ViT). For the former, we successfully achieved persistent activity in the TN hardware by leveraging sparsely coupled sWTA motifs, a critical requirement for instituting working memory. Additionally, when this approach was applied as a pre-processing layer into the ViT architecture, we observed a substantial increase in classification accuracy for previously unseen test data. Together these results suggest that adapting biophysical principles to neuromorphic chips may offer a promising pathway for NeuroAI performance.

## Results

Our objective is to extract general principles of neocortical function, such as soft Winner-Take-All (sWTA), and implement them efficiently on neuromorphic hardware. By leveraging the hardware's parametric constraints, we aim to apply these simplified neocortical computations to enhance working memory capabilities. This approach not only mimics brain-like processing but also has the potential to improve AI models' performance across various machine learning tasks, bridging the gap between neuroscience and artificial intelligence.

## Biophysical model implementation of neocortical circuit motifs

Sensory information processing in primary sensory cortices, such as V1 relies on pyramidal neurons integrating bottom-up signals with top-down feedback from higher-order visual areas. Key components in this process include recurrent excitation, feedforward and feedback inhibition, disinhibition, and divisive normalization. These functions are primarily mediated by parvalbumin (PV), somatostatin (SST), vasoactive intestinal peptide (VIP), and lysosomal-associated membrane protein 5 (LAMP5) interneurons, respectively, working together to selectively amplify thalamic inputs (Reinhold et al. [2015], Reinhold et al. [2015], Pfeffer et al. [2013], Lien and Scanziani [2013]). Mouse V1 exhibits strong recurrent connections among layer 2/3 (L2/3) pyramidal neurons (Ko et al. [2011], Harris and Mrsic-Flogel [2013], Rossi et al. [2020]), with PV interneurons providing local feedforward inhibition and SST interneurons delivering global feedback inhibition targeting axon initial segments and L2/3 dendrites (Schneider-Mizell et al. [2021], Atallah et al. [2012], Naka et al. [2019], Adesnik and Scanziani [2010]). VIP interneurons modulate feedback inhibition (Karnani et al. [2016], Pfeffer et al. [2013]), while LAMP5 cells regulate top-down and bottom-up signals through normalization (Ibrahim et al. [2021], Huang et al. [2023], Malina et al. [2021], Hartung et al. [2024]).

Consistent with earlier studies, our experimental data confirmed that L2/3 pyramidal neurons receive strong inhibition during bottom-up sensory input stimulation in V1 (Supplementary Fig. 1A-C). To dissect the specific contributions of PV and SST interneurons, we employed optogenetics in PV-Cre and SST-Cre mice and analyzed how their activation modulated the current-frequency (f/I) response of L2/3 pyramidal neurons (Supplementary Fig. 1D-F). Additionally, top-down auditory cortex inputs, which primarily convey contextual information, were found to target Lamp5 expressing neurogliaform interneurons (Supplementary Fig. 1G-J). Though not explicitly tested, we modeled Lamp5-mediated global inhibition via volumetric transmission (Ibrahim et al. [2021], Huang et al. [2023]) affecting L2/3 pyramidal neurons. Finally, VIP interneurons, although not directly included, likely modulated SST-mediated inhibition through disinhibition (Karnani et al. [2016], Pfeffer et al. [2013]).

Next, we built a biophysically detailed network model in the NEURON simulation environment to validate whether the simplified cortical circuit motifs obtained from V1 indeed implemented sWTA computations. To do this, we incorporated excitatory and inhibitory cell types with diverse spiking patterns in a conductance-based Hodgkin-Huxley network model (Supplementary Fig. 2A-D). Synaptic parameters were constrained by our experimental data (Supplementary Fig. 1F). Poisson-modulated excitatory synaptic inputs were used to assess the input-output (IO) function of L2/3 pyramidal neurons under PV and SST inhibition (Supplementary Fig. 2A). Synaptic weights were tuned to match experimentally observed inhibitory postsynaptic potentials, adjusting the pyramidal neuron f/I curve (Supplementary Fig. 2D-F). SST inhibition primarily influenced the slope of the IO function, while PV inhibition altered the offset (Supplementary Fig. 2E-F). Lamp5-mediated volumetric inhibition was represented as non-specific inhibition across pyramidal neuron dendrites, achieved by reducing both SST and recurrent excitatory weights. Notably, although VIP inhibition was not explicitly modeled, its effects were captured by reducing SST weights.

Based on experimental data, we developed a generalized cortical microcircuit model comprising 10 pyramidal neurons, 10 PV interneurons, 1 SST interneuron, and 1 Lamp5 interneuron, with their biophysical properties constrained by our mouse V1 physiology data (Figure 1A). We examined the computational behavior of this model under Poisson-modulated excitatory synaptic inputs mimicking thalamic activity, where stronger inputs were selectively amplified, and weaker inputs suppressed (Figure 1B). Stronger thalamic inputs represent a tuned orientation or direction information carried to L2/3 pyramidal neurons in V1. The soft winner-take-all (sWTA) mechanism enables neural circuits to prioritize the strongest input by balancing competitive and cooperative interactions. Non-linear amplification of stronger inputs, coupled with suppression of weaker ones, forms the core of sWTA computations. Consistent with previous findings (Douglas et al. [1995], Somers et al. [1995], Reinhold et al. [2015]), our model reproduces these dynamics through recurrent excitation and lateral feedback inhibition (Figure 1B-C). Using this model, we assessed how various interneuron populations contribute to sWTA computations, focusing on their role in enhancing responses to strongly

stimulated neurons (Pyr 4) while suppressing weaker responses (Pyr 1-3, 5-7; Figure 1C). By modulating the conductances of PV, SST, VIP, and Lamp5 interneurons within physiological ranges, we quantified their influence on pyramidal neuron dynamics (Figure 1D-E). PV and SST inhibition independently shaped the width and gain of the sWTA function, while Lamp5-mediated inhibition primarily adjusted gain (Figure 1F). VIP-mediated disinhibition was examined by reducing SST synaptic weights.

We computed a selectivity index to evaluate the network's ability to suppress weaker inputs while amplifying stronger ones. This index, calculated as the difference in pyramidal neuron responses to closely tuned inputs, revealed that recurrent excitation, along with PV and SST inhibition, is crucial for maintaining high selectivity (Supplementary Fig. 2G-I). Lamp5 inhibition modulates gain, preserving tuning specificity while enabling flexible responses to factors like attention and locomotion (Ferguson and Cardin [2020], Bugeon et al. [2022]). Our model replicates key features of cortical circuits, showing both linear gain scaling and non-linear selectivity. These dynamics are governed by the balance between excitatory and inhibitory synaptic weights and the feedback inhibition threshold, modulated by network activity. Deviations from optimal weights diminished sparsity by amplifying secondary inputs (Figure 1G-I). Linear gain modulation enhanced weak thalamic inputs (Oldenburg et al. [2024], Sievers et al. [2024], Lien and Scanziani [2018], Lien and Scanziani [2013]), while non-linear computations such as signal restoration for sharply tuned, noise-embedded inputs (Figure 1G-H) — aligned with *in vivo* observations. Furthermore, the model captures hysteresis and multistability, enabling the circuit to amplify relevant inputs based on initial conditions or contextual cues (Figure 1I).

In summary, our model provides a biophysical foundation for a simplified and generalized computational mechanism, such as soft Winner-Take-All (sWTA), which may represent a universal computation in sensory cortices (Douglas and Martin [2007], Niell and Stryker [2008]), that might contribute to orientation and direction tuning in the visual cortex (Hubel and Wiesel [1962]), angular whisker tuning in the barrel cortex (Lavzin et al. [2012]), and frequency tuning in the auditory cortex (Kato et al. [2017]).

## Mapping neocortical algorithms onto IBM TrueNorth neuromorphic hardware

A few years ago, IBM released their neuromorphic TrueNorth (TN) chip, offering a reconfigurable, asynchronous, multi-core digital architecture ideal for implementing brain-inspired computations. We aimed to program the TN chip to implement a simplified sWTA computational primitive, inspired by the neocortex. A key challenge was that TN's neural dynamics were governed by parameters such as thresholds, leaks, and crossbar weight that did not directly align with biophysical or artificial neural network models. Notably, these strongly resemble gain modulation regulated by the four interneuron types considered in the biophysical modeling described previously.

To implement an sWTA computation, we developed an automated gain-matching technique to match TN network dynamics to the biophysical model, enabling accurate parameter mapping (Appendix A1). Initially, we created an abstract rate-based model that mimicked the input-output (IO) function of the biophysical neurons. We then derived constraints to map these dynamics onto the TN hardware, producing a linear threshold response that closely approximated the physiological behavior of cortical neurons (Figure 2A-C). Inputs to the TN network were generated by configuring the on-chip neurons within neurosynaptic cores to produce a range of frequencies (Figure 2D-E). This allowed us to match the TN network's IO gain to the abstract model (Figure 2F). Using contraction theory (Rutishauser and Douglas [2009]; Appendix A2), we derived optimal TN parameters, enabling us to implement all sWTA operations, as we performed in our biophysical analysis of V1 processing, including signal restoration, hysteresis, and multi-stability (Figure 2G-I).

After programming the TN chip to perform sWTA computations, we retrospectively compared TN parameters such as thresholds and leaks with those in our biophysical model. The optimized TN parameters closely aligned with the functions of the excitatory-inhibitory balance observed in the biophysical model (Supplementary Fig. 2E-F). Specifically, TN parameters such as threshold and leak mirrored the roles of PV and Lamp5-mediated inhibition in the biophysical model (Supplementary Fig. 3A-C). Moreover, the TN neurons replicated the recurrent excitation and global feedback inhibition

motifs found in cortical circuits, mirroring the effects of SST interneurons on pyramidal neuron IO functions, as well as the role of VIP interneurons in disinhibiting this population (Figure 2E). Notably, the parameters derived using our gain-matching technique closely resembled those observed in experimentally constrained models of neocortical circuits across all conditions, highlighting the fidelity of TN hardware in simulating cortical computations (Supplementary Fig. 3D-F).

While our initial efforts efficiently mapped the sWTA computations in rate mode, we extended the method to implement population-level sWTA networks in spiking mode configuration (Figure 2J). This configuration, tested on a population of 10 excitatory neurons gated by shared inhibition, allowed us to translate rate-based computations into spiking dynamics, better reflecting neocortical organization. Under optimal conditions, the TN network's dynamic range during sWTA closely matched the firing rates of cortical neurons in V1 (Figure 2K-L). We evaluated whether the TN network in spiking mode could implement sWTA computations under noisy conditions, simulating biological variability (Figure 2K-L). The noise was controlled using a parameter called threshold mask noise (TMN), which emulated spontaneous cortical activity. At TMN values up to 10, synaptic and spiking variability resulted in stable network dynamics that supported sWTA operations (Supplementary Fig. 3G-H). This demonstrated that TN spiking networks remain stable in noisy environments and avoid synchrony driven by inhibition. Without inhibition, excitatory activity would increase exponentially; however, when inhibitory neurons are activated, their gain is tuned to stabilize excitatory activity. This balance between positive and negative feedback forms an attractor state.

We next aimed to demonstrate a practical use case for sWTA circuit motifs in hardware applications, specifically by implementing a neural state machine (NSM). We hypothesized that NSMs could serve as foundational elements for complex cognitive tasks in robotics, and will benefit from the energy-efficient framework of sWTA networks. A key aspect of cognitive function is working memory, which allows for the retention of cue information even in the absence of stimuli, enabling appropriate action selection based on environmental cues. NSMs with working memory encode stimuli as distinct states, transitioning between them to support

context-dependent tasks (e.g., Fuster and Alexander [1971], Wang [2001], Harvey et al. [2012]). Previous studies have shown that sparsely coupled sWTA motifs can sustain persistent activity (Neftci et al. [2013], Rutishauser and Douglas [2009]). Given the conserved use of circuit motifs across cortical areas, we hypothesized that this sWTA architecture would efficiently support working memory. To implement stable attractor states, TN hardware parameters were tuned to balance positive and negative feedback as outlined in Neftci et al. [2013] to avoid inhibition-mediated synchrony, which is critical for maintaining persistent activity in an NSM using spiking dynamics. We tested whether these motifs could facilitate action selection in response to environmental cues while retaining information in their absence. Our results showed that TN neurosynaptic cores, initially designed for sensory sWTA, could be effectively repurposed for NSM implementation. Using sparsely coupled sWTA motifs, we achieved persistent activity in TN hardware, with time constants that aligned closely with experimental data.

We then evaluated the ability of this architecture to implement a two-state NSM. Transitions between states S1 and S2 were driven by input signals ( $X$ ,  $Y$ ) and regulated by pointer neurons (P12, P21) within the sparsely coupled sWTA motif (Supplementary Figure 4A-C). This configuration generated stable sWTA dynamics and persistent attractor states (Neftci et al. [2013]; Appendix A2). Consistent with previous findings, gamma coupling through bidirectional excitatory weights in TN hardware sustained persistent activity even without external input. When noise was introduced, it disrupted synchronous firing but optimizing gamma coupling maintained stable persistence (Supplementary Fig. 4D). By fine-tuning the coupling strength, we identified the minimal gamma required for maintaining persistent activity across varying noise levels. Striking this balance was critical for stability, especially when environmental cues were unreliable. We further confirmed the stability of attractor states by removing one transition input, demonstrating that the circuit continued to sustain activity (Supplementary Fig. 4E).

In summary, we developed a two-state NSM where transitions were governed by input signals and the current state (Supplementary Fig. 4F). The sWTA dynamics facilitated smooth state transitions within a finite state automaton (FSA)

framework. Efficiency analysis revealed that both time and energy in TN hardware scaled linearly with the number of states and computational load, contrasting with the quadratic scaling observed in Compass simulations. Notably, runtime on TN hardware was independent of firing rates, synapse activity, and neuron counts, with asynchronous state updates (Supplementary Fig. 4G-H; Appendix A3). This underscores the efficiency of neuromorphic hardware in implementing NSMs, supporting higher cognitive functions.

## Neocortex-inspired WTA implementation for Artificial Intelligence applications

**Performance boost in Image Classification:** Finally, we tested whether incorporating a pre-processing WTA layer into deep learning models, such as Vision Transformers (ViTs), could enhance performance on real-world vision tasks. Specifically, we explored the role of sWTA computations in spatial feature extraction for object classification tasks. We developed a novel neural layer inspired by the hardware-constrained sWTA motif and integrated it into the conventional ViT architecture to assess its impact on classifying unseen digit datasets. This approach (Appendix A4) leveraged sWTA as a pre-processing layer to reduce redundancies and enhance contrast in visual inputs. A sliding window-based computation (Figure 3A) was employed for feature amplification, minimizing domain shifts (Figure 3B). This allowed parameters extracted by the TN hardware constraints to execute sWTA computations using recurrent excitation and lateral inhibition across pixels. In this setup, the patch with the highest variance, or "winner patch," received the highest normalized value, while other patches were scaled accordingly. The selection of the "salient" patch size was optimized to maintain stable circuit dynamics in line with TN hardware parameters.

We next evaluated domain generalization to assess the ability of the sWTA model to adapt to unseen data distributions — an ideal test for mimicking the brain's ability to generalize across diverse sensory inputs and maintain robust performance in novel environments. We trained the ViT model, with and without the WTA layer (Figure 3C), on a single source domain, and then tested its performance on unseen target domains. Significant improvements were observed across all source/target combinations (Figure 3D). Beyond ViT, similar

**Table 1:** Comparison of classification accuracy between WTA-based DNN architectures (Vision Transformer, EfficientNet, CapsuleNet, MobileNet, and ResNet) and source-only models trained on the MNIST (M), MNIST-M (MM), SVHN (S), and USPS (U) datasets with respective combinations as highlighted on top of each column. The top panel shows the performance of the models without adding the WTA layer whereas the bottom panel shows the performance boost by adding the WTA layer to the network architectures. The models are tested on completely ‘unseen’ target datasets. (**bold-red** indicates the best and **bold-black** indicates the 2nd best)

Source-only Models	M→U	U→M	S→M	M→S	M→MM	MM→M
ViT	75.0	72.0	66.6	22.0	42.0	98.3
EfficientNet	77.9	50.7	61.4	18.6	18.8	95.0
CapsuleNet	<b>96.4</b>	<b>87.2</b>	58.1	11.8	22.5	<b>98.4</b>
MobileNet	84.4	60.0	72.2	22.4	33.9	97.3
ResNet	82.5	58.5	63.4	27.2	38.2	97.4
ViT+WTA	84.26	78.0	<b>73.6</b>	<b>52.7</b>	70.0	98.1
EfficientNet+WTA	83.5	74.2	69.1	19.6	48.8	96.3
CapsuleNet+WTA	<b>94.1</b>	<b>87.8</b>	<b>75.9</b>	32.1	57.2	<b>98.6</b>
MobileNet+WTA	82.5	70.9	73.5	<b>40.5</b>	<b>73.4</b>	97.8
ResNet+WTA	82.8	66.0	71.2	27.7	<b>70.2</b>	97.8

performance gains were observed in EfficientNet (Tan and Le [2019]), CapsuleNet (Sabour et al. [2018]), MobileNet-V2 (Sandler et al. [2018]), and ResNet-50 (He et al. [2016]), where the WTA layer enhanced the models’ ability to learn generalizable features, improving domain shift robustness in object recognition tasks for MNIST and other digit datasets (Table 1). Figure 3B illustrates how the WTA layer reduces domain shift, showing high similarity across sample images post-processing. These results were achieved without intensity-based augmentation, using only geometric augmentations. EfficientNet, MobileNet, and ResNet-50 were initialized with pre-trained ImageNet weights (Russakovsky et al. [2015]). Supplementary Figures 5A-B present train/loss curve examples for ViT and CapsuleNet, while Supplementary Table 2 details model architectures and training settings. Table 1 summarizes the results, with Supplementary Table 1 showing performance improvements compared to state-of-the-art models. We also compared the WTA layer’s ability to minimize domain shift against traditional pre-processing techniques such as Local Response Normalization (LRN) Krizhevsky et al. [2017], Local Contrast Normalization (LCN) Placidi and Polzinelli [2021] and Z-score normalization. Supplementary Figure 6A provides qualitative comparisons, while Supplementary Figure 6B visualizes UMAP embeddings of MNIST and

MNIST-M datasets post-normalization. Our WTA implementation significantly minimized domain shift (Supp Fig. 6B), leading to a marked improvement in classification accuracy (0.7) on unseen test data compared to baseline techniques (Supp Fig. 6C).

**Performance boost in Image Segmentation:** Lastly, we evaluated our approach in a deep learning model for the challenging task of natural image segmentation. Similar to the results seen in image classification, incorporating the sWTA layer into the RefineNet architecture Lin et al. [2017]) with ResNet-101 significantly improved performance in semantic segmentation. This underscores the broad applicability of our approach across various vision tasks. The model was trained on the Cityscapes dataset (Cordts et al. [2016]), consisting of 2,975 daytime driving images, and tested on 50 coarsely annotated Nighttime Driving dataset (Dai and Van Gool [2018]) (Figure 3E). We employed the same training setup, using a dynamic learning rate of 0.1 with stochastic gradient descent (SGD) on an RTX A6000 GPU and a batch size of 6. Performance was measured using mean Intersection over Union (mIoU).

Notably, adding the WTA layer to RefineNet achieved performance on par with nighttime driving data using only source-trained models (Figure 3F).

**Table 2:** Segmentation performance on Nighttime Driving (*Dai and Van Gool [2018]*), reported as mIoU scores. WTA-RefineNet outperforms other methods trained only on daytime data and has a competitive performance to methods also using nighttime images

Method	Nighttime Driving (mIOU)
RefineNet Lin et al. [2017]	34.1
W-RefineNet Lin et al. [2017] (pre-trained RGB weights of Resnet on ImageNet)	34.6
RefineNet-AdaBN Lin et al. [2017]	36.3
WTA-RefineNet (pre-trained RGB weights of Resnet on ImageNet)	<b>36.5</b>

Table 2 highlights the performance improvements in RefineNet for natural image segmentation tasks, with and without the sWTA layer, demonstrating its robustness in handling domain shifts. Supplementary Fig. 7 shows qualitative results for the nighttime dataset, comparing performance with and without the WTA layer.

In conclusion, implementing sWTA motifs as a layer in various deep learning architectures substantially improves their performance. Using the competitive-cooperative dynamics of biologically inspired WTA mechanisms as a preprocessing layer creates a synergy between biological principles and artificial neural networks. Our WTA-inspired layer enhances performance in real-world tasks like image classification and segmentation by acting as an adaptive filter that sharpens focus on the most relevant features while reducing noise and irrelevant information.

## Discussion

The overall goal of this study was threefold: First, to validate that computational principles, such as winner-take-all (WTA), are implemented by neocortical circuit motifs by delineating the contribution of four major cardinal interneuron classes in biophysical models. Second, to emulate WTA computational primitives and extend them to construct neural state machines on IBM’s TrueNorth (TN) neuromorphic hardware. Third, to integrate hardware-derived parametric constraints into deep learning models, such as Vision Transformers, demonstrating that biological principles can enhance performance and reduce training sessions through zero-shot learning. Figure 4 shows the block diagram architecture of our proposed framework.

Our gain-matching technique allows for the emulation of any neural network architecture and computational principle on neuromorphic hardware, similar to previous studies on analog systems (Neftci et al. [2011]). This work builds on prior neural computation research that has applied sWTA networks to object recognition (Yuille and Grzywacz [1988]; Riesenhuber and Poggio [1999]; Erlhagen and Schöner [2002]), attention (Itti et al. [1998]; Deco and Rolls [2005]), orientation selectivity (Ben-Yishai et al. [1995]; Somers et al. [1995]), decision making (Amari and Arbib [1977]) and sparse coding (Rozell et al. [2008]). Notably, our proposed framework (Figure 4) can be easily adapted for next-generation neuromorphic platforms like Braindrop (Neckar et al. [2018]), IBM’s Northpole (Modha et al. [2023]), and Intel’s Lohi2 (Davies et al. [2021]). Facilitating a direct translation of computational insights from theoretical neurobiology to hardware and AI systems.

Building on prior studies, we demonstrate that generic circuit motifs used for sensory computations can be adapted to implement working memory for cognitive decision-making tasks on TN neuromorphic hardware (Rutishauser and Douglas [2009], Neftci et al. [2013]). This multiplexing of sensory and decision-making computations is both hardware-friendly and efficient, with computation time and energy scaling linearly with the number of states or computational load. Our approach provides a scalable solution for mapping large state machines, outperforming current Long Short-Term Memory (LSTM) models. While the components in LSTM architectures scale quadratically with the number of states, our method scales linearly, offering a more efficient alternative. Future research could investigate the role of specific interneurons and circuit

motifs beyond sensory regions, potentially revealing principles for context-dependent decision-making.

Our framework can also be extended to mimic dendritic computation in cortical neurons and implement dendrocentric learning rules (Boahen [2022]), such as BCM plasticity (Bienenstock et al. [1982]), spike-timing dependent plasticity (STDP; Bi and Poo [1998]), and non-Hebbian behavioral time-scale plasticity (BTSP; Bittner et al. [2017]). These principles could inspire better hardware design for constructing self-learning AI cognitive agents with behavioral flexibility akin to that of animals. In this study, we focused on implementing aspects of the primary visual cortex within Vision Transformers and deep learning models. Future work could extend this framework to other cortical regions involved in sensory processing and decision-making. Together by implementing realistic biophysical models to neuromorphic systems we are able to improve AI networks, providing a key step towards fostering NeuroAI development.

## Acknowledgments

We thank Ehsan Arabzadeh, William Connelly, and Giacomo Indiveri for the fruitful discussion. We thank Jan Drugowitsch, Christopher Harvey, Debanjan Dasgupta, Alessandro Galloni, and Samuel Gershman for their helpful comments and constructive criticism of the manuscript. We thank Dharmendra Modha, Hayley Hu, and Ben Shaw for assistance and guidance on the TrueNorth hardware. The neuromorphic implementation was performed at the Institute of Neuroinformatics (INI), UZH/ETH Zurich as a part of the IBM-INI collaboration. We thank the organizers of the Telluride workshop and IBM Bootcamp. We acknowledge the generous support of high-computing GPUs from Tibbling Technologies for running experiments to train and test the WTA implementation in deep learning architectures.

## Author Contributions

A.I., G.F., and S.H. designed research and wrote the manuscript; S.H. performed *in-vitro* experiments, biophysical modeling, and neuromorphic hardware implementation; A.I. and S.H. conceived mapping of neocortical computation in deep learning architectures. A.I. and H.M. performed mapping of WTA in deep learning architectures and ran experiments for image classification and segmentation tasks; G.J.S

supervised *in-vitro* experiments; G.F supervised biophysical modeling and contributed to the organization of the manuscript.

## Competing Interest Statement

The authors have declared no competing interests.

## Methods

**Animals:** All experimental procedures were approved by and conducted in accordance with Harvard Medical School and The Australian National University Institutional Animal Care and Ethics Committee.

**Viral injections and Whole-cell patch-clamp recordings:** For labeling bottom-up sensory input and top-down contextual inputs, AAV1-hSyn-hChR2(H134R)-EYFP-WPRE-hGH (ChR2) was injected in either the contralateral visual cortex (or visual thalamus, dLGN) and primary auditory cortex respectively (Honnuraiah et al. [2024]; Godenzini et al. [2021]). Ipsilateral eye input is stimulated by contralateral V1, and contralateral eye input is stimulated by dLGN stimulation (Honnuraiah et al. [2024]). For PV and SST inhibition experiments, Cre-dependent ChR2 (AAV1-EF1a-DIO-hChR2(H134R)-EYFP-WPRE-hGH) was injected in the binocular visual cortex of the transgenic mice expressing Cre in either PV or SST. Three to four weeks after viral injection mice were deeply anesthetized with isoflurane (3% in oxygen) and immediately decapitated. Slice preparation protocols and the experimental recordings are explained in detail in this study (Honnuraiah et al. [2024]). All recordings were made in the current-clamp using a current clamp BVC-700A amplifier (Dagan Instruments, USA). Data was filtered at 10 kHz and acquired at 50 kHz by a Macintosh computer running Axograph X acquisition software (Axograph Scientific, Sydney, Australia) using an ITC-18 interface (Instrutech/HEKA, Germany). Hyperpolarizing and depolarizing current steps (200 pA to +600 pA; intervals of 50 pA) were applied via the somatic recording pipette to characterize the passive and active properties of neurons. Brain slices were bathed in gabazine (10  $\mu$ M) to block inhibition mediated by GABA-A receptors. Other pharmacological agents used in these experiments included tetrodotoxin (TTX; 1  $\mu$ M) and 4-aminopyridine (4-AP; 100  $\mu$ M), as

noted in the Results. For photo-stimulation of ChR2-expressing neurons and axon terminals, a 470 nm LED (Thorlabs) was mounted on the epi-fluorescent port of the microscope (Olympus BX50) allowing wide-field illumination through the microscope objective. The timing, duration, and strength of LED illumination were controlled by the data acquisition software (Axograph).

**Computational Modeling:** We conducted biophysical modeling using the NEURON 8.2 simulation environment (Carnevale and Hines [2006], Hines and Carnevale [1997]), with an integration time constant of 25  $\mu$ s. The active and passive properties of the model were optimized to match the experimental recordings (Supp Fig 1). We set the passive parameters as follows: Internal/axial resistance ( $R_i/R_a$ ) to 150  $\Omega\text{cm}$ , membrane resistance ( $R_m$ ) to 30  $K\Omega\text{cm}^2$ , capacitance ( $C_m$ ) to 1  $\mu\text{F}/\text{cm}^2$  and resting membrane potential ( $V_m$ ) to -75 mV. All neurons were simplified and implemented as a “ball and stick” model consisting of a somatic compartment (dimensions: Length=50  $\mu\text{m}$ ; diameter=50  $\mu\text{m}$ ) and a single dendritic compartment (dimensions: Length=100  $\mu\text{m}$ ; diameter=1  $\mu\text{m}$ ). Dendritic compartments were passive and were not adjusted for spines in the interneuron population but were adjusted for spines in pyramidal neurons by scaling the  $C_m$  by 2 and  $R_m$  by 0.5. Active conductances were included in the somatic compartment to mimic the regular firing pattern of pyramidal neurons, fast-spiking pattern of PV, burst-spiking for SST, and delayed spiking from the Lamp5. Active ion-channel distribution and its conductance values are obtained from our previous study (Soldado-Magraner et al. [2020]). A synapse was modeled as a co-localized combination of NMDA and AMPA receptor currents. A default value of NMDAR:AMPAR ratio was set at 1.5. All the values related to synaptic parameters were obtained from our previous study (Honnuraiah and Narayanan [2013], Testa-Silva et al. [2022]).

**Rate-based abstract neural network model:** We developed a simplified abstract model to reduce computational demands and extract the principles from detailed biophysical network models. We have used a rate-based approach to model neuronal activity. We approximate the neuronal activation by a linear-threshold function that describes the output action potential discharge rate of the neuron as a function of its input Bauer et al. [2014]. This type of neuronal activation function is a good approximation

to experimental and biophysical observations of the frequency of action potential discharge to synaptic or current inputs. The change in activity of a neuron is modeled as the summation of synaptic input with a decay of the current activity. The dynamics of the activity of the rate-based neurons implemented are given below:

$$\eta \frac{\partial exct}{\partial t} = -exct + \alpha[exct]^+ - \beta_1[inhb]^+ + I \quad (1)$$

$$\eta \frac{\partial inhb}{\partial t} = -inhb + \beta_2[exct]^+ \quad (2)$$

$$exct^+ = \max(exct, 0) \quad (3)$$

$$inhb^+ = \max(inhb, 0) \quad (4)$$

Where,  $[exct]^+$  is excitatory neuron activity,  $[inhb]^+$  is inhibitory neuron activity, and  $\tau$  is the neuronal time constant,  $\alpha$ ,  $\beta_1$ , and  $\beta_2$  are synaptic weights (Figure 2B).

The implementation of the computational primitives obtained from the biophysical models in the rate-based abstract models was crucial. This is because it provided analytically tractable solutions for the dynamics of neural activity. The analytical solution was later used to derive constraints for the TN hardware parameters as explained in the Appendix A2.

**IBM TrueNorth hardware and Compass software emulator:** The IBM TN neuromorphic chip is composed of  $64 \times 64$  (4096) digital neurosynaptic cores tiled in a 2-D array, containing an aggregate of 1 million neurons and 256 million synapses. Each core implements 256 neurons single-compartment leaky-integrate-and-fire neurons which could be operated in either rate or spike mode configuration. Each core is supported by a  $256 \times 256$  crossbar synapse array, and communication circuits to transfer spike trains. The crossbar array is flexible and can be configured freely. Each row of the crossbar corresponds to an axon of the neuron represented by horizontal lines which could be driven by any on-chip neurons. Inputs to the cores are generated by configuring the on-chip neurons to generate various frequencies either in rate or spike mode. Each column corresponds to a dendrite of that particular neuron represented by a horizontal line. A connection between an axon and a dendrite is a synapse and is organized into a synaptic crossbar (Supplementary Figure 4A). A peripheral memory core is located at the intersection of each row and a column, and the

binary value stored in the core represents whether or not a connection exists between the particular axon-dendrite pair. Therefore, each neuron can be configured to receive up to 1024 synaptic inputs (through its dendrite) depending on the crossbar value and the activity of the axons. TN operates in a mixed asynchronous–synchronous approach. All the communication and control circuits operate in asynchronous design while computations are done in synchronous design. Since TN cores operate in parallel and are governed by spike events, it is natural to implement all the routing mechanisms asynchronously. All the core computations must finish with finish in the current tick which spans 1 ms. Compass is the software emulator to program and simulate the full 4096 neurosynaptic cores and the digital asynchronous–synchronous design ensures one-to-one compass to TN correspondence (Akopyan et al. [2015] Merolla et al. [2014]).

**Vision Transformer (ViT) architecture:** Inspired by the original transformer (Vaswani [2017]) architecture for Natural Language Processing, ViT (Dosovitskiy et al. [2020]) is a self-attention-based architecture. It works as follows: the input image is distributed into N (flattened 2D) patches (where we keep N=6 for the MNIST experiments) and linear embeddings of these patches are fed as an input to the encoder of the trained ViT. The image patches of digits are embedded as tokens. The encoder block has several multi-headed layers with self-attention along with a normalization layer at the start of each layer. Furthermore, a Multi-Layer Perceptron (MLP) with a single hidden layer is used as a classification head that predicts the object categories present in the input image.

**WTA implementation in ViT:** We present a Winner-Take-All (WTA) approximation as a neuro-inspired layer in the Vision Transformer (ViT) architecture. Inspired by distinctive properties of cortical circuits in the mammalian visual cortex, this layer captures the neural signal regulation characteristics. A defining feature of these neurons lies in their ability to intricately capture and encode the contrasting elements and structural nuances of visual stimuli. This capability is reflected in the variable neural firing sequence of the neurons that WTA emulates, aligning closely with the varying contrasts present in the stimuli. This results in a more refined and contextually aware representation, which is particularly beneficial in object classification contexts where adaptability and nuanced

understanding of visual stimuli are crucial. Mathematical implementation is described in Appendix A4.

**MobileNet with WTA:** Developed for embedded devices such as mobile phones, etc., MobileNet-v2 (Sandler et al. [2018]) is successfully reducing the number of parameters by depth-wise separable convolutions, while keeping the accuracy comparable to the state-of-the-art. We initialized the weights of MobileNet-V2, trained on ImageNet, and added the WTA layer to the model.

**EfficientNet with WTA:** EfficientNet-B0 (Tan and Le [2019]) is used with pre-trained weights for ImageNet. It is pre-trained to classify 1000 image classes and trained on more than a million images. We initialized the convolution layers with the pre-trained model weights. We added the WTA layer after the input layer into the model.

**ResNet with WTA:** Among different variants of ResNet, we select ResNet-50 (He et al. [2016]), which contains 50 neural network layers. The introduction of skip connections reduces the problem of vanishing gradient and also ensures that the higher layers do not perform any worse than the layers before by learning the identity function. Similar to the other models above, we are initializing the weights with the pre-trained model on ImageNet and including the WTA layer after the input layer in the model architecture.

**CapsuleNet with WTA:** Unlike other Convolutional Neural Network (CNN) architectures, CapsuleNet (Sabour et al. [2018]) applies pattern matching by decomposing the hierarchical representations of the input features. The eventual representation of this network is supposed to be invariant to the view-angle of the input samples. One of the major differences between a typical CNN and CapsuleNet is the output of the individual units in their architecture. While the output of a single neuron in a CNN is mostly a scalar value, it is a vector in the case of CapsuleNet. Similar to the ViT, we include the WTA layer as an initial layer in the CapsuleNet architecture to make it robust for domain adaptation tasks for digit datasets.

**RefineNet with WTA:** RefineNet (Lin et al. [2017]) is a versatile multi-path refinement network that leverages all the information gathered during the down-sampling process to facilitate high-resolution prediction through long-range residual

connections. This approach enables the deeper layers, which capture high-level semantic features, to be directly refined using fine-grained features from earlier convolutions. The individual components of RefineNet employ residual connections following the identity mapping principle, enabling efficient end-to-end training. In our experiments, we used pre-trained RGB weights of ResNet on ImageNet for training and testing RefineNet with and without adding a WTA layer.

**MNIST and digit datasets:** For the domain generalization task, we utilized a suite of digit datasets that included the MNIST, SVHN, USPS, and MNIST-M (LeCun et al. [2010], Netzer et al. [2011], Hull [1994], Ganin et al. [2016]). Each dataset was split into 70-20-10 train, val, and test splits.

**MNIST dataset**, introduced by LeCun et al. [2010], is one of the most widely used datasets for handwritten digit classification. It contains a total of 70,000 grayscale images of handwritten digits. Each image is of size 28x28 pixels, and the dataset has been instrumental in benchmarking various machine learning algorithms.

**SVHN (Street View House Numbers) dataset**, presented by Netzer et al. [2011] is a real-world image dataset obtained from house numbers in Google Street View images. It comprises over 600,000 digit images. Specifically, it contains 73,257 digits for training, 26,032 digits for testing, and an additional 531,131 somewhat less difficult samples that can be used as extra training data. This dataset challenges models with recognizing digits in more complex and varied scenarios compared to the controlled environment of MNIST.

**USPS (United States Postal Service) dataset**, introduced by Hull [1994] is another handwritten digit dataset used for text recognition research. It contains 9,298 16x16 grayscale images of handwritten digits. The dataset was derived from scanned mail and has been a staple in the handwritten digit recognition field.

**MNIST-M dataset**, presented in the work by Ganin et al. [2016], is a modified version of the original MNIST dataset. It was created by overlaying MNIST digits onto patches randomly extracted from color photos of the BSDS500 dataset (Arbelaez et al. [2010]), resulting in a blend of digits and colored backgrounds. The MNIST-M dataset contains 149,002 images. This combination introduces additional challenges due to the color and texture variations in the background, making

it a valuable dataset for studying domain adaptation.

For domain generalization tasks, these datasets are particularly valuable because they offer variations in terms of image quality, resolution, and real-world applicability. The diversity in these datasets, ranging from clean handwritten digits to digits in natural scenes, challenges models to generalize well across different domains. This makes them ideal benchmarks for evaluating the robustness and adaptability of machine learning algorithms, especially in scenarios where the training and test data distributions differ significantly.

**Natural image dataset:** To explore the effect of WTA for segmentation tasks on natural images, we select a cityscape Cordts et al. [2016] data for training and nighttime driving Dai and Van Gool [2018] dataset for testing. This evaluation aims to test the robustness of the model against day-to-night time domain shifts.

## Appendix

### 0.1 Emulating Cortical Neuron Physiology in TN Hardware

We have derived a relationship between the parameters of the COMPASS neurons in order to obtain the desired dynamic range, shown below:

$$\Delta_r = \frac{f_{out}^{max} - f_{out}^{min}}{f_{in}^{max} - f_{in}^{min}} \quad (5)$$

$$W_{syn}^1 = \frac{\Delta_r}{N_{syn}} \quad (6)$$

$$\|\tilde{W}_{ij}\|_{min}^{max} = \alpha \cdot W_{syn}^1 \quad (7)$$

$$\lambda_l = \frac{\|\tilde{W}_{ij}\|_{min}^{max} \cdot N_{syn}}{\Delta_r} \quad (8)$$

Here,

$\Delta_r$  = Dynamic range of the neuron,

$W_{syn}^1$  = Sensitivity of single weight,

$N_{syn}$  = Total number of synapses,

$\tilde{W}_{ij}$  = Crossbar synaptic weight,

$\|\tilde{W}_{ij}\|_{min}^{max}$  = Range of the crossbar weight,

$\lambda_l$  = Leak parameter of the TN neuron.

The threshold ( $\lambda_l$ ) of the COMPASS neuron is decided based on the desired dynamic range and the total number of synapses ( $N_{syn}$ ), such that the sensitivity of a single synapse is preserved while setting up the actual crossbar synaptic weights ( $\tilde{W}_{ij}$ ).

For example, let us assume that we want to set the parameters of a COMPASS neuron that receives 10 synaptic inputs and has a dynamic range of 1 ( $\Delta_r$ ). If all of the synaptic weights are equal, then each synapse will have an impact factor of 0.1 ( $W_{sym}^1$ ) and we want a range of 50 for the crossbar synaptic weight ( $\|\tilde{W}_{ij}\|_{min}^{max}$ ). Then, according to the above equations, the threshold should be set to 2.

$$\begin{aligned}\lambda_l &= \frac{\|\tilde{W}_{ij}\|_{min}^{max} \cdot N_{syn}}{\Delta_r} \\ &= \frac{10 \cdot 10}{50}\end{aligned}$$

Thus, by using this relation, we can set the parameters of the TN neuron to obtain any behavior within the physiological range that closely matches the cortical neurons. TN simulation results below verifying the above relation and to understand the role of synaptic weight on the transfer function.

We tuned the feedback and feedforward inhibition of the model to match the impact of PV and SST activation on the pyramidal neuron. The subtractive inhibition is obtained by tuning the threshold ( $\theta$ ) value. The divisive inhibition is implemented by tuning the crossbar synaptic weight ( $\tilde{W}_{ij}$ ) to negative values and leak parameter value ( $\lambda_l$ ). Subtractive inhibition is implemented by tuning the threshold value of the TN neuron. The TN simulation results were verified with a conductance-based model implemented in NEURON.

We have derived a relationship between the parameters of the TN neurons to incorporate biophysically plausible excitatory and inhibitory synaptic interaction as shown below:  $\Delta_r$  = Dynamic range of the Neuron.

$N_{syn}$  = Total number of synapses.

$N_{ext}$  = Excitatory synapses.

$N_{inh}$  = Inhibitory synapses.

$\tilde{W}_{ext}$  = Crossbar excitatory weight.

$\|\tilde{W}_{ext}\|_{min}^{max}$  = range of the excitatory weight.

$\tilde{W}_{inh}$  = Crossbar inhibitory weight.

$\|\tilde{W}_{inh}\|_{min}^{max}$  = range of the inhibitory weight.

$\lambda_l$  = leak of the TN neuron.

$$\|\lambda_l\|_{min}^{max} = \begin{cases} \|\tilde{W}_{ext}\|_{min}^{max} \left( \frac{N_{ext}}{N_{inh}} \right) - 1 & \text{if } N_{inh} \neq 0 \\ 2 \cdot (\tilde{W}_{ext}^{max} - 1) \left( \frac{\Delta_r}{N_{syn}} \right) & \text{if } N_{inh} = 0 \end{cases}$$

$$\text{where, } \|\tilde{W}_{ext}\|_{min}^{max} = \frac{\lambda_l \cdot \Delta_r}{N_{syn}}$$

Thus, by appropriately tuning the parameters of the TN neuron  $\{\alpha, \lambda_l, \tilde{W}_{ij}\}$  we can achieve the desired, biophysically realistic synaptic integration that closely matches the cortical neurons.

## 0.2 Automated Parameter Mapping to TN parameters

TN neurons are configured to operate as linear threshold units (in rate-based mode) as described in the previous section. Based on this linear operation, we can estimate the role of the self-excitatory feedback connection on the transfer function, according to the equation below:  $f_{out}$  = output firing rate;  $f_{in}$  = Input rate;  $\lambda_l$  = Leak of TN neuron.

$$\begin{aligned}f_{out} &= \frac{w}{\lambda_l} (f_{in} - \theta) + \frac{\alpha}{\lambda_l} \cdot f_{out} \\ f_{out} \left( 1 - \frac{\alpha}{\lambda_l} \right) &= \frac{w}{\lambda_l} \cdot f_{in} - \left( \frac{w}{\lambda_l} - \theta \right) \\ \text{gain} &= \frac{f_{out}}{f_{in}} = \frac{w}{\lambda_l - \alpha}\end{aligned}$$

Based on the linear operation, we can estimate the role of the recurrent excitatory and inhibitory feedback connection on the transfer function, according to the equation below:

$$f_{out} = \left( \frac{w}{\lambda_l} \right) \cdot f_{in} + \left( \frac{\alpha}{\lambda_l} \right) \cdot f_{out} - \left( \frac{\beta_1 \cdot \beta_2}{\lambda_l} \right) \cdot f_{out} \quad (9)$$

$$f_{out} \left( 1 - \frac{\alpha}{\lambda_l} + \frac{\beta_1 \cdot \beta_2}{\lambda_l} \right) = \frac{w}{\lambda_l} \cdot f_{in} \quad (10)$$

$$\text{gain} = \frac{f_{out}}{f_{in}} = \frac{w}{\lambda_l - \alpha + \beta_1 \cdot \beta_2} \quad (11)$$

The effect of leak and threshold on the transfer function is quantified and the relationship between the parameters is shown below:  $f_{out}$  = output firing rate;  $f_{in}$  = Input rate;  $\theta$  = Threshold,  $\lambda_l$  = Leak of TN neuron and  $\sigma$  = Sign of Leak.

$$\begin{aligned}f_{out} &= \left( \frac{w}{1 - \sigma \lambda} \right) \cdot f_{in} + \left( \frac{\alpha}{1 - \sigma \lambda} \right) \cdot f_{out} \\ &\quad - \left( \frac{\beta_1 \beta_2}{1 - \sigma \lambda} \right) \cdot f_{out} - (\Omega[\theta, w] \cdot \theta)\end{aligned} \quad (12)$$

$$\begin{aligned}f_{out} \cdot \left( 1 - \frac{\alpha}{(1 - (\sigma \lambda))} + \frac{\beta_1 \beta_2}{1 - \sigma \lambda} \right) &= \left( \frac{w}{1 - \sigma \lambda} \right) \cdot f_{in} - (\Omega[\lambda, w] \cdot \lambda)\end{aligned} \quad (13)$$

$$f_{\text{out}} = \left( \frac{w}{1 - (\alpha + \sigma\lambda) + \beta_1\beta_2} \right) \cdot f_{\text{in}} - \left( \frac{\Omega[\theta, w] \cdot \theta}{1 - (\alpha + \sigma\lambda) + \beta_1\beta_2} \right) \quad (14)$$

$$\text{gain} = \left( \frac{w}{1 - (\alpha + \sigma\lambda) + \beta_1\beta_2} \right) \quad (15)$$

$$\Omega[\theta, w] = \begin{cases} \frac{f_{\text{out}, \theta=0}^{\max} - f_{\text{out}, \theta=1}^{\max}}{\text{gain}} & \text{if } \theta < \text{int}\left(\frac{w}{2}\right) \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

To achieve winner-take-all behavior, the parameters have to satisfy certain constraints imposed by Contraction analysis shown below:  $\alpha_m$ ;  $\beta_{1m}$ ;  $\beta_{2m}$  = parameters in programming platform.  $\alpha_c$ ;  $\beta_{1c}$ ;  $\beta_{2c}$  = parameters in TN. For the programming model, the parameters must satisfy the following criteria:

$$0 < \alpha_m < 2\sqrt{\beta_1^m \cdot \beta_2^m} \quad (17)$$

$$0 < \beta_1^m \quad (18)$$

$$0 < \beta_1^m \cdot \beta_2^m < 1 \quad (19)$$

The optimal solution in the programming environment's parametric space satisfying the above conditions is:

$$\alpha^m = 1.2, \quad \beta_1^m = -3, \quad \beta_2^m = 0.25 \quad (20)$$

Now, using these parameters we can obtain the corresponding Truenorth/Compass parameter values according to the equation we have derived that maps the parameters from the programming platform to TrueNorth/Compass space, shown below:

$$|\beta_1^c \cdot \beta_2^c - \alpha^c| = w_c (1 - \alpha_m + \beta_1^m \cdot \beta_2^m) - \theta^c \quad (21)$$

Substituting the values, we obtain the corresponding optimal TrueNorth/Compass parameters that satisfy the contraction analysis criteria in TrueNorth/Compass space:

$$\alpha^c = 22, \quad \beta_1^c = -7, \quad \beta_2^c = 1, \quad \theta^c = 20, \quad w_c = 10 \quad (22)$$

We plug in the above values in the TrueNorth/Compass circuit shown in Figure 2 and verify if the following WTA functional characteristics are satisfied:

1. Non-linear signal amplification (winner selection). (Validated in Figure 2G)
2. Robustness and signal restoration (broadly tuned inputs). (Figure 2H)

### 3. Dynamic switching and multi-stability. (Validated in Figure 2I)

Thus, by appropriately tuning the parameters of the TN neuron  $\{\alpha, \lambda_l, W_{ij}\}$  we can achieve the desired, biophysically realistic synaptic integration that closely matches the cortical neurons.

## 0.3 Hardware load analysis

Computation load =  $C(N) \times \text{numTicks}$

$$C(N) = I(N) + O(N)$$

$C(N)$  = Number of Connector pins

$I(N)$  = Number of Input pins

$O(N)$  = Number of Output pins

## 0.4 Mathematical formulation of WTA in ViT

To mathematically implement the WTA layer, we process an input image  $I \in \mathbb{R}^{W \times H \times C}$ , where  $W$ ,  $H$ , and  $C$  represent its width, height, and channel count, respectively. Our goal is to form a domain-independent representation, denoted as  $I^G$ . The image  $I$  is segmented into patches of size  $s$ , represented as  $P = \{p_{s1}, p_{s2}, \dots, p_{sn}\}$ , where each patch  $p$  of size  $s$  encircles a pixel  $k$  at coordinates  $i, j$ . For each patch, its mean  $\mu_{ps}$  and standard deviation  $\sigma_{ps}$  are calculated to construct  $I^G$ :

$$\sigma_{ps} = \left( \frac{1}{s^2} \sum_{i,j \in p,s} (k_{ij} - \mu_{ps})^2 \right)^{1/2}, \quad (23)$$

where:

$$\mu_{ps} = \frac{1}{s^2} \sum_{i,j \in p,s} k_{ij}$$

$$z = \max\{\sigma_{ps1}, \sigma_{ps2}, \dots, \sigma_{psn}\} \quad (24)$$

$$I^G = \frac{\sigma_{ps}}{z} \quad (25)$$

## References

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- Murray Shanahan, Kyle McDonell, and Laria Reynolds. Role play with large language models. *Nature*, 623(7987):493–498, 2023.
- Murray Campbell, A Joseph Hoane Jr, and Feng-hsiung Hsu. Deep blue. *Artificial intelligence*, 134(1-2):57–83, 2002.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- Anthony Zador, Sean Escola, Blake Richards, Bence Ölveczky, Yoshua Bengio, Kwabena Boahen, Matthew Botvinick, Dmitri Chklovskii, Anne Churchland, Claudia Clopath, et al. Catalyzing next-generation artificial intelligence through neuroai. *Nature communications*, 14(1):1597, 2023.
- Peter Dayan. Simple substrates for complex cognition. *Frontiers in neuroscience*, 2:411, 2008.
- Wael F Asaad, Gregor Rainer, and Earl K Miller. Task-specific neural activity in the primate pre-frontal cortex. *Journal of neurophysiology*, 84(1):451–459, 2000.
- Abhishek Banerjee, Giuseppe Parente, Jasper Teutsch, Christopher Lewis, Fabian F Voigt, and Fritjof Helmchen. Value-guided remapping of sensory cortex by lateral orbitofrontal cortex. *Nature*, 585(7824):245–250, 2020.
- Duo Xu, Mingyuan Dong, Yuxi Chen, Angel M Delgado, Natasha C Hughes, Linghua Zhang, and Daniel H O'Connor. Cortical processing of flexible and context-dependent sensorimotor sequences. *Nature*, 603(7901):464–469, 2022.
- Cameron Condylis, Eric Lowet, Jianguang Ni, Karina Bistrong, Timothy Ouellette, Nathaniel Josephs, and Jerry L Chen. Context-dependent sensory processing across primary and secondary somatosensory cortex. *Neuron*, 106(3):515–525, 2020.
- Bernardo Rudy, Gordon Fishell, Soohyun Lee, and Jens Hjerling-Leffler. Three groups of interneurons account for nearly 100% of neocortical gabaergic neurons. *Developmental neurobiology*, 71(1):45–61, 2011.
- Robin Tremblay, Soohyun Lee, and Bernardo Rudy. Gabaergic interneurons in the neocortex: from cellular properties to circuits. *Neuron*, 91(2):260–292, 2016.
- Gord Fishell and Adam Kepecs. Interneuron types as attractors and controllers. *Annual review of neuroscience*, 43:1–30, 2020.
- Adam Kepecs and Gordon Fishell. Interneuron cell types are fit to function. *Nature*, 505(7483):318–326, 2014.
- Katie A Ferguson and Jessica A Cardin. Mechanisms underlying gain modulation in the cortex. *Nature Reviews Neuroscience*, 21(2):80–92, 2020.
- Cristopher M Niell and Massimo Scanziani. How cortical circuits implement cortical computations: Mouse visual cortex as a model. *Annual Review of Neuroscience*, 44:381–404, 2021.
- Carsten K Pfeffer, Mingshan Xue, Miao He, Z Josh Huang, and Massimo Scanziani. Inhibition of inhibition in visual cortex: the logic of connections between molecularly distinct interneurons. *Nature neuroscience*, 16(8):1068–1076, 2013.
- Luke Campagnola, Stephanie C Seeman, Thomas Chartrand, Lisa Kim, Alex Hoggarth, Clare Gamlin, Shinya Ito, Jessica Trinh, Pasha Davoudian, Cristina Radaelli, et al. Local connectivity and synaptic dynamics in mouse and human neocortex. *Science*, 375(6585):eabj5861, 2022.
- Matteo Carandini and David J Heeger. Normalization as a canonical neural computation. *Nature reviews neuroscience*, 13(1):51–62, 2012.
- Gregor Schöner and John P Spencer. *Dynamic thinking: A primer on dynamic field theory*. Oxford University Press, 2016.
- Saurabh Vyas, Matthew D Golub, David Sussillo, and Krishna V Shenoy. Computation through neural population dynamics. *Annual review of neuroscience*, 43:249–275, 2020.
- Georg B Keller and Thomas D Mrsic-Flogel. Predictive processing: a canonical cortical computation. *Neuron*, 100(2):424–435, 2018.
- Andre M Bastos, W Martin Usrey, Rick A Adams, George R Mangun, Pascal Fries, and Karl J Friston. Canonical microcircuits for predictive coding. *Neuron*, 76(4):695–711, 2012.
- Rodney J Douglas and Kevan AC Martin. Recurrent neuronal circuits in the neocortex. *Current biology*, 17(13):R496–R500, 2007.

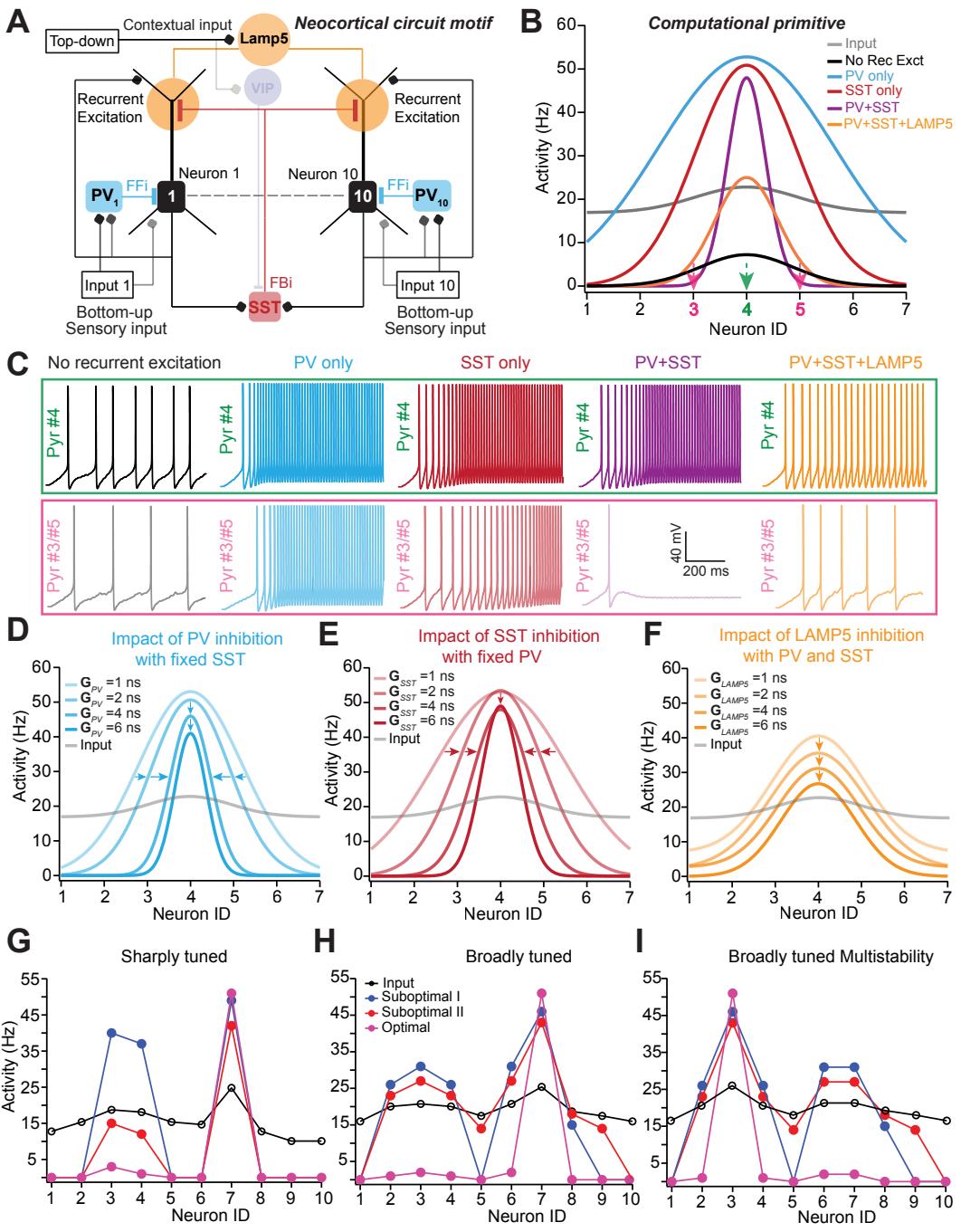
- Aygul Balcioglu, Rebecca Gillani, Michael Doron, Kendyll Burnell, Taeyun Ku, Alev Erisir, Kwanghun Chung, Idan Segev, and Elly Nedivi. Mapping thalamic innervation to individual l2/3 pyramidal neurons and modeling their ‘readout’ of visual input. *Nature Neuroscience*, 26(3):470–480, 2023.
- Anthony D Lien and Massimo Scanziani. Cortical direction selectivity emerges at convergence of thalamic synapses. *Nature*, 558(7708):80–86, 2018.
- Rita Bopp, Simone Holler-Rickauer, Kevan AC Martin, and Gregor FP Schuhknecht. An ultrastructural study of the thalamic input to layer 4 of primary motor and primary somatosensory cortex in the mouse. *Journal of Neuroscience*, 37(9):2435–2448, 2017.
- Tom Binzegger, Rodney J Douglas, and Kevan AC Martin. A quantitative map of the circuit of cat primary visual cortex. *Journal of Neuroscience*, 24(39):8441–8453, 2004.
- Elisabetta Chicca, Fabio Stefanini, Chiara Bartolozzi, and Giacomo Indiveri. Neuromorphic electronic circuits for building autonomous cognitive systems. *Proceedings of the IEEE*, 102(9):1367–1388, 2014.
- Ning Qiao, Hesham Mostafa, Federico Corradi, Marc Osswald, Fabio Stefanini, Dora Sumislawska, and Giacomo Indiveri. A reconfigurable on-line learning spiking neuromorphic processor comprising 256 neurons and 128k synapses. *Frontiers in neuroscience*, 9:141, 2015.
- Giacomo Indiveri and Yulia Sandamirskaya. The importance of space and time for signal processing in neuromorphic agents: the challenge of developing low-power, autonomous agents that interact with the environment. *IEEE Signal Processing Magazine*, 36(6):16–28, 2019.
- Dharmendra S Modha, Filipp Akopyan, Anthony Andreopoulos, Raghavendra Appuswamy, John V Arthur, Andrew S Cassidy, Pradeep Datta, Michael V DeBole, Steven K Esser, Carlos O Otero, Jun Sawada, Brian Taba, Arnon Amir, Dhiren Bablani, Paul J Carlson, Myron D Flickner, Ravi Gandhasri, Gilles J Garreau, Makoto Ito, Joseph L Klamo, Jason A Kusnitz, Janice L McClatchey, Jeffrey L McKinstry, Yasuhiro Nakamura, Tapan K Nayak, William P Risk, Kurt Schleupen, Brian Shaw, Jayadev Sivagnanam, David F Smith, Ivo Terrizzano, and Takashi Ueda. Neural inference at the frontier of energy, space, and time. *Science*, 112(658):22–25, 2023.
- Paul A Merolla, John V Arthur, Rodrigo Alvarez-Icaza, Andrew S Cassidy, Jun Sawada, Filipp Akopyan, Bryan L Jackson, Nabil Imam, Chen Guo, Yutaka Nakamura, et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345(6197):668–673, 2014.
- Alexander Neckar, Sam Fok, Ben V Benjamin, Terrence C Stewart, Nick N Oza, Aaron R Voelker, Chris Eliasmith, Rajit Manohar, and Kwabena Boahen. Braindrop: A mixed-signal neuromorphic architecture with a dynamical systems-based programming model. *Proceedings of the IEEE*, 107(1):144–164, 2018.
- Nabil Imam. Wiring up recurrent neural networks. *Nature Machine Intelligence*, 3(9):740–741, 2021.
- L Federico Rossi, Kenneth D Harris, and Matteo Carandini. Spatial connectivity matches direction selectivity in visual cortex. *Nature*, 588(7839):648–652, 2020.
- Cristopher M Niell and Michael P Stryker. Highly selective receptive fields in mouse visual cortex. *Journal of Neuroscience*, 28(30):7520–7536, 2008.
- David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *Journal of Physiology*, 160(1):106–154, 1962.
- Yazan N Billeh, Binghuang Cai, Sergey L Gratiy, Kael Dai, Ramakrishnan Iyer, Nathan W Gouwens, Reza Abbasi-Asl, Xiaoxuan Jia, Joshua H Siegle, Shawn R Olsen, et al. Systematic integration of structural and functional data into multi-scale models of mouse primary visual cortex. *Neuron*, 106(3):388–403, 2020.
- Kimberly Reinhold, Alexander D Lien, and Massimo Scanziani. Distinct recurrent versus afferent dynamics in cortical visual processing. *Nature Neuroscience*, 18(12):1789–1797, 2015.
- Anthony D Lien and Massimo Scanziani. Tuned thalamic excitation is amplified by visual cortical circuits. *Nature neuroscience*, 16(9):1315–1323, 2013.
- Ho Ko, Sonja B Hofer, Bruno Pichler, Katherine A Buchanan, P Jesper Sjöström, and Thomas D

- Mrsic-Flogel. Functional specificity of local synaptic connections in neocortical networks. *Nature*, 473(7345):87–91, 2011.
- Kenneth D Harris and Thomas D Mrsic-Flogel. Cortical connectivity and sensory coding. *Nature*, 503(7474):51–58, 2013.
- Casey M Schneider-Mizell, Agnes L Bodor, Forrest Collman, Derrick Brittain, Adam Bleckert, Sven Dorkenwald, Nicholas L Turner, Thomas Macrina, Kisuk Lee, Ran Lu, et al. Structure and function of axo-axonic inhibition. *Elife*, 10:e73783, 2021.
- Bassam V Atallah, William Bruns, Matteo Carandini, and Massimo Scanziani. Parvalbumin-expressing interneurons linearly transform cortical responses to visual stimuli. *Neuron*, 73(1):159–170, 2012.
- Alexander Naka, Julia Veit, Ben Shababo, Rebecca K Chance, Davide Risso, David Stafford, Benjamin Snyder, Andrew Egladyous, Desiree Chu, Savitha Sridharan, et al. Complementary networks of cortical somatostatin interneurons enforce layer specific control. *Elife*, 8:e43696, 2019.
- Hillel Adesnik and Massimo Scanziani. Lateral competition for cortical space by layer-specific horizontal circuits. *Nature*, 464(7292):1155–1160, 2010.
- Mahesh M Karnani, Jesse Jackson, Inbal Ayzenshtat, Ali Hamzehei Sichani, Kiarash Manoocheri, Sooyeon Kim, and Rafael Yuste. Opening holes in the blanket of inhibition: Localized lateral disinhibition by vip interneurons. *J Neurosci*, 36(13):3471–3480, 2016.
- Leena Ali Ibrahim, Shuhan Huang, Marian Fernandez-Otero, Mia Sherer, Yanjie Qiu, Spuriti Vemuri, Qing Xu, Robert Machold, Gabrielle Pouchelon, Bernardo Rudy, et al. Bottom-up inputs are required for establishment of top-down connectivity onto cortical layer 1 neurogliaform cells. *Neuron*, 109(21):3473–3485, 2021.
- Shuhan Huang, Sherry Jingjing Wu, Giulia Sansone, Leena Ali Ibrahim, and Gord Fishell. Layer 1 neocortex: Gating and integrating multidimensional signals. *Neuron*, 2023.
- Katayun Cohen-Kashi Malina, Emmanouil Tsivourakis, Dahlia Kushinsky, Daniella Apelblat, Stav Shtiglitz, Eran Zohar, Michael Sokoletsky, Gen-ichi Tasaka, Adi Mizrahi, Ilan Lampl, et al. Ndnf interneurons in layer 1 gain-modulate whole cortical columns according to an animal's behavioral state. *Neuron*, 109(13):2150–2164, 2021.
- Jan Hartung, Anna Schroeder, Rodrigo Alejandro Pérez Vázquez, Rogier B Poorthuis, and Johannes J Letzkus. Layer 1 ndnf interneurons are specialized top-down master regulators of cortical circuits. *Cell Reports*, 43(5), 2024.
- Rodney J Douglas, Christof Koch, Misha Mahowald, Kevan AC Martin, and Humbert H Suarez. Recurrent excitation in neocortical circuits. *Science*, 269(5226):981–985, 1995.
- David C Somers, Sacha B Nelson, and Mriganka Sur. An emergent model of orientation selectivity in cat visual cortical simple cells. *Journal of Neuroscience*, 15(8):5448–5465, 1995.
- Stephane Bugeon, Joshua Duffield, Mario Dipoppa, Anne Ritoux, Isabelle Pranker, Dimitris Nicoloutsopoulos, David Orme, Maxwell Shinn, Han Peng, Hamish Forrest, et al. A transcriptomic axis predicts state modulation of cortical interneurons. *Nature*, 607(7918):330–338, 2022.
- Ian Antón Oldenburg, William D Hendricks, Gregory Handy, Kiarash Shamardani, Hayley A Bounds, Brent Doiron, and Hillel Adesnik. The logic of recurrent circuits in the primary visual cortex. *Nature neuroscience*, pages 1–11, 2024.
- Meike Sievers, Alessandro Motta, Martin Schmidt, Yagmur Yener, Sahil Loomba, Kun Song, Johannes Bruett, and Moritz Helmstaedter. Connectomic reconstruction of a cortical column. *bioRxiv*, pages 2024–03, 2024.
- Maria Lavzin, Sophia Rapoport, Alon Polsky, Liora Garion, and Jackie Schiller. Nonlinear dendritic processing determines angular tuning of barrel cortex neurons in vivo. *Nature*, 490(7420):397–401, 2012.
- Hiroyuki K Kato, Samuel K Asinof, and Jeffry S Isaacson. Network-level control of frequency tuning in auditory cortex. *Neuron*, 95(2):412–423, 2017.
- Ueli Rutishauser and Rodney J Douglas. State-dependent computation using coupled recurrent networks. *Neural computation*, 21(2):478–509, 2009.
- Joaquin M Fuster and Garrett E Alexander. Neuron activity related to short-term memory. *Science*, 173(3997):652–654, 1971.

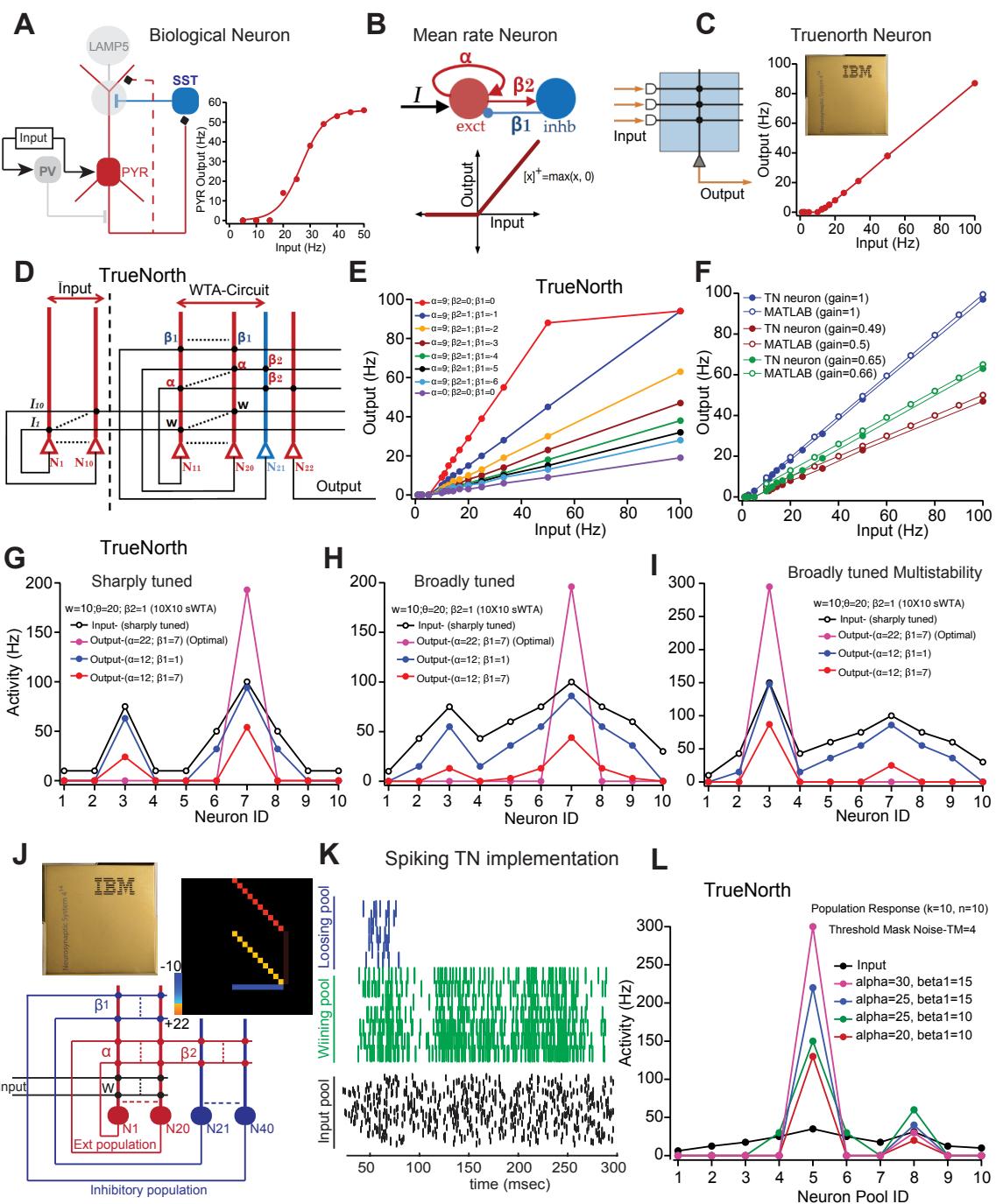
- Xiao-Jing Wang. Synaptic reverberation underlying mnemonic persistent activity. *Trends in neurosciences*, 24(8):455–463, 2001.
- Christopher D Harvey, Philip Coen, and David W Tank. Choice-specific sequences in parietal cortex during a virtual-navigation decision task. *Nature*, 484(7392):62–68, 2012.
- Emre Neftci, Jonathan Binas, Ueli Rutishauser, Elisabetta Chicca, Giacomo Indiveri, and Rodney J Douglas. Synthesizing cognition in neuromorphic electronic systems. *Proceedings of the National Academy of Sciences*, 110(37):E3468–E3476, 2013.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- Sara Sabour, Nicholas Frosst, and Geoffrey Hinton. Matrix capsules with em routing. In *6th international conference on learning representations, ICLR*, volume 115, 2018.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- Giuseppe Placidi and Matteo Polzinelli. Local contrast normalization to improve preprocessing in mri of the brain. In *International Conference on Bioengineering and Biomedical Signal and Image Processing*, pages 255–266. Springer, 2021.
- Dengxin Dai and Luc Van Gool. Dark model adaptation: Semantic image segmentation from daytime to nighttime. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 3819–3824. IEEE, 2018.
- Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1925–1934, 2017.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- Emre Neftci, Elisabetta Chicca, Giacomo Indiveri, and Rodney Douglas. A systematic method for configuring vlsi networks of spiking neurons. *Neural computation*, 23(10):2457–2497, 2011.
- Alan Yuille and Norbert Grzywacz. A computational theory for the perception of coherent visual motion. *Nature*, 333:71–74, 1988.
- Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex. *Nature Neuroscience*, 2(11):1019–1025, 1999.
- Wolfram Erlhagen and Gregor Schöner. Dynamic field theory of movement preparation. *Psychological review*, 109(3):545, 2002.
- Laurent Itti, Ernst Niebur, and Christof Koch. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, 1998.
- Gustavo Deco and Edmund Rolls. Neurodynamics of biased competition and cooperation for attention: a model with spiking neurons. *Journal of Neurophysiology*, 94(4):295–313, 2005.
- Rani Ben-Yishai, R Lev Bar-Or, and Haim Sompolinsky. Theory of orientation tuning in visual cortex. *Proceedings of the National Academy of Sciences*, 92(9):3844–3848, 1995.
- Shun-Ichi Amari and Michael A Arbib. Competition and cooperation in neural nets. *Systems neuroscience*, pages 119–165, 1977.

- Christopher J Rozell, Don H Johnson, Richard G Baraniuk, and Bruno A Olshausen. Sparse coding via thresholding and local competition in neural circuits. *Neural computation*, 20(10):2526–2563, 2008.
- Mike Davies et al. Taking neuromorphic computing to the next level with loihi2. *Intel Labs' Loihi*, 2: 1–7, 2021.
- Kwabena Boahen. Dendrocentric learning for synthetic intelligence. *Nature*, 612(7938):43–50, 2022.
- Elie L Bienenstock, Leon N Cooper, and Paul W Munro. Theory for the development of neuron selectivity: orientation specificity and binocular interaction in visual cortex. *Journal of Neuroscience*, 2(1):32–48, 1982.
- Guo-Qiang Bi and Mu-ming Poo. Synaptic modifications in cultured hippocampal neurons: dependence on spike timing, synaptic strength, and postsynaptic cell type. *J Neurosci*, 18(24):10464–10472, 1998.
- Katie C Bittner, Aaron D Milstein, Christine Grienberger, Sandro Romani, and Jeffrey C Magee. Behavioral time scale synaptic plasticity underlies ca1 place fields. *Science*, 357(6355):1033–1036, 2017.
- Suraj Honnuraiah, Helena H-Y Huang, William J Ryan, Robin Broersen, William M Connelly, and Greg Stuart. Cellular and circuit mechanisms underlying binocular vision. *bioRxiv*, pages 2024–03, 2024.
- L Godenzini, D Alwis, Robertas Guzulaitis, Suraj Honnuraiah, GJ Stuart, and LM Palmer. Auditory input enhances somatosensory encoding and tactile goal-directed behavior. *Nature Communications*, 12(1):4509, 2021.
- Nicholas T Carnevale and Michael L Hines. *The NEURON book*. Cambridge University Press, 2006.
- Michael L Hines and Nicholas T Carnevale. The neuron simulation environment. *Neural computation*, 9(6):1179–1209, 1997.
- Saray Soldado-Magraner, Federico Brandalise, Suraj Honnuraiah, Michael Pfeiffer, Marie Moulinier, Urs Gerber, and Rodney Douglas. Conditioning by subthreshold synaptic input changes the intrinsic firing pattern of ca3 hippocampal neurons. *Journal of neurophysiology*, 123(1):90–106, 2020.
- Suraj Honnuraiah and Rishikesh Narayanan. A calcium-dependent plasticity rule for hcn channels maintains activity homeostasis and stable synaptic learning. *PloS one*, 8(2):e55590, 2013.
- Guilherme Testa-Silva, Marius Rosier, Suraj Honnuraiah, Robertas Guzulaitis, Ana Morello Megias, Chris French, James King, Katharine Drummond, Lucy M Palmer, and Greg J Stuart. High synaptic threshold for dendritic nmda spike generation in human layer 2/3 pyramidal neurons. *Cell reports*, 41(11), 2022.
- Roman Bauer, Frédéric Zubler, Sabina Pfister, Andreas Hauri, Michael Pfeiffer, Dylan R Muir, and Rodney J Douglas. Developmental self-construction and-configuration of functional neocortical neuronal networks. *PLoS computational biology*, 10(12):e1003994, 2014.
- Filipp Akopyan, Jun Sawada, Andrew Cassidy, Rodrigo Alvarez-Icaza, John Arthur, Paul Merolla, Nabil Imam, Yutaka Nakamura, Pallab Datta, Gi-Joon Nam, et al. Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neuromimetic chip. *IEEE transactions on computer-aided design of integrated circuits and systems*, 34(10):1537–1557, 2015.
- Ashish Vaswani. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- Yann LeCun, Corinna Cortes, Chris Burges, et al. Mnist handwritten digit database, 2010.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 7. Granada, Spain, 2011.
- Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Fleuret, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of machine learning research*, 17(59):1–35, 2016.
- Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *IEEE transactions*

*on pattern analysis and machine intelligence*, 33  
(5):898–916, 2010.

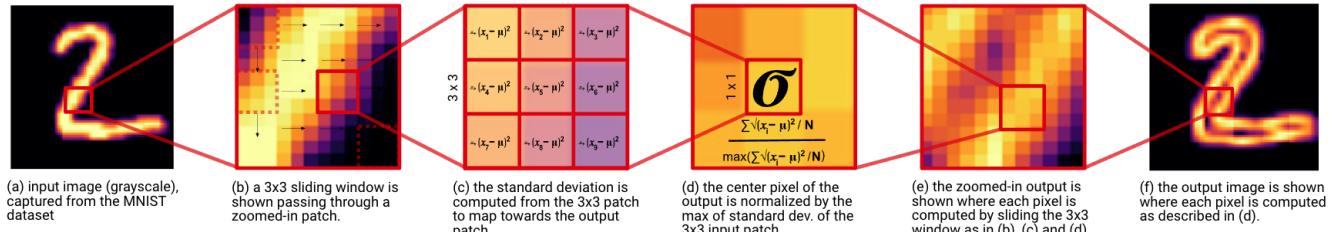


**Figure 1:** Biophysical implementation of experimentally validated neocortical circuit motifs using a conductance-based neural network model. **A)** Canonical neocortical circuit motif implemented using a conductance-based Hodgkin-Huxley neuron network model consisting of 10 pyramidal neurons, 10 PV interneurons providing feedforward and local feedback inhibition, and one common inhibitory neuron by LAMP5 for contextual modulation and one SST interneuron providing lateral inhibition. Sensory input is Poisson modulated excitatory synaptic inputs with varying frequencies. **B)** Computational primitives implemented by the neocortical circuit motifs performing non-linear input amplification, selective suppression, and soft-winner take-all (sWTA) computation. The contribution of various circuit elements like distinct interneurons: PV, SST, and LAMP5 mediated inhibition and recurrent excitation to sWTA computation is shown in various colors. Gaussian fits of the data points are used to extract quantifiable parameters such as mean and width. **C)** Voltage traces of representative pyramidal neurons receiving the highest (in green) and next to highest (in pink) input shown during various conditions: without recurrent excitation (black), with recurrent excitation and PV plus SST inhibition (magenta) along with LAMP5 inhibition (orange), and only SST (red) and only PV inhibition (cyan). **D-F)** Impact of modulating distinct inhibition on the computational primitives implemented by the circuit motifs in (A). Varying PV inhibitory conductance from 1-6 nS with a fixed SST inhibition of 1.5 nS along with recurrent excitation (D). Varying SST inhibitory conductance from 1-6 nS with a fixed PV inhibition of 1.5 nS along with recurrent excitation (E). Varying LAMP5 inhibitory conductance from 1-6 nS with an SST and PV inhibition along with recurrent excitation (F). **G-I)** Verification of sWTA properties such as nonlinear amplification and signal restoration for sharply (G) and broadly (H) tuned inputs. Multistability and signal invariance for broadly (I) tuned inputs for suboptimal (blue, red) and optimal (magenta) parameters are derived from our proposed method.

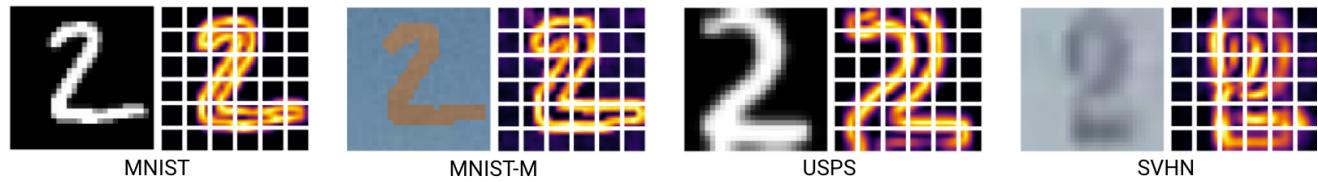


**Figure 2:** Configuring IBM TN neuromorphic hardware to implement neocortical circuit motifs and computational primitives. **A)** Left: Simplified biophysical circuit implementation (excitatory neurons in red; lateral inhibition in blue). Right: Transfer function of the L2/3 pyramidal neuron when stimulated with Poisson input. Stimulus frequency is varied between 5 to 50 Hz (in steps of 5 Hz) and the corresponding output response is plotted. **B)** Top: Abstract rate-encoding neuron model containing one excitatory (red) and inhibitory (blue) neuron with linear-threshold activation function (bottom). **C)** Left: Single TrueNorth neuron with crossbar weights along with input and output. Right: Transfer function of the TrueNorth neuron in a rate-encoding configuration similar to the transfer function of biological and abstract neuron models. **D)** Wiring diagram showing the connection configuration between 10 excitatory and 1 feedback inhibitory truenorth neurons for implementing the circuit motif shown in Figure 1A. **E)** Impact of recurrent excitation ( $\alpha$ ) and feedback inhibition ( $\beta_1$ ,  $\beta_2$ ) on the slope of the transfer function of the TrueNorth neuron shown in. **F)** Matching the slope/gain of the transfer function of the TN neuron with the abstract rate encoding neuron model implemented in software by tuning the TN parameters as described in Appendix A2, A3. **G-I)** Parameter tuning and verification of sWTA properties like nonlinearity amplification and signal restoration in TN simulation and TN hardware to sharply tuned inputs (G), signal invariance property of sWTA to broadly tuned inputs (H), and multistability (I). Output response for optimal parameters is shown in magenta, inputs are shown in black, and responses for suboptimal parameters are shown in blue and red. **J-L)** Population-level implementation of sWTA in the spiking-mode configuration of Truenorth neural network.

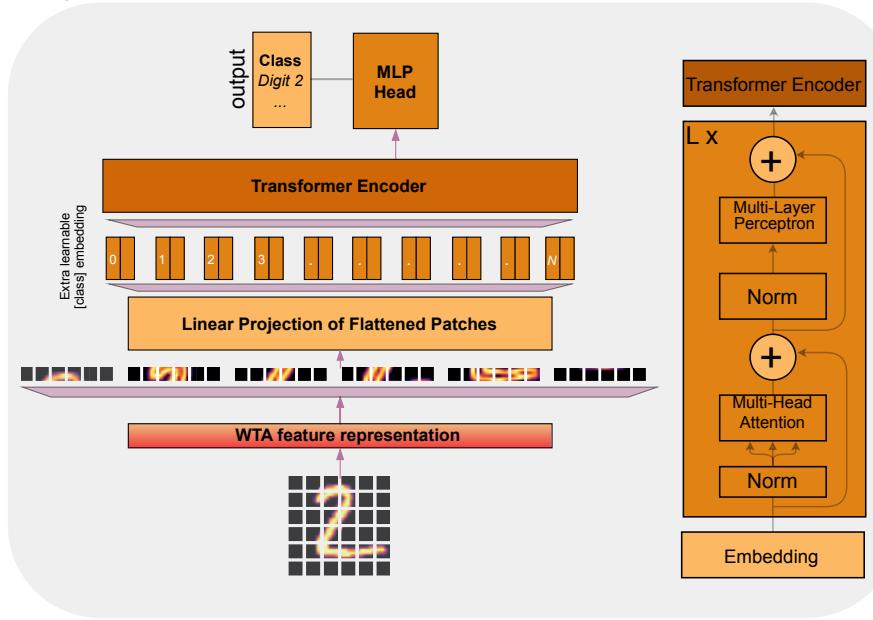
**A. Processing of WTA layer on a digit 2 MNIST image sample**



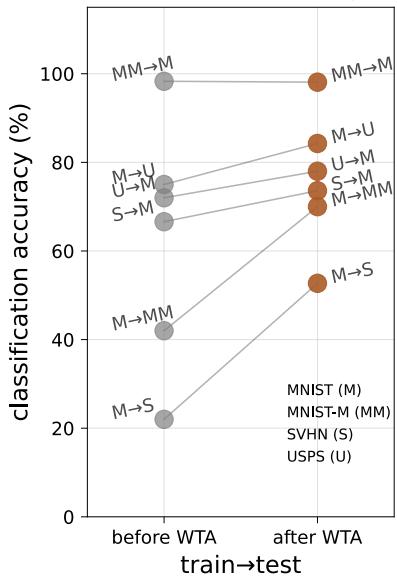
**B. WTA layer-based image patches for ViT on digit 2 image samples drawn from MNIST, MNIST-M, USPS and SVHN datasets**



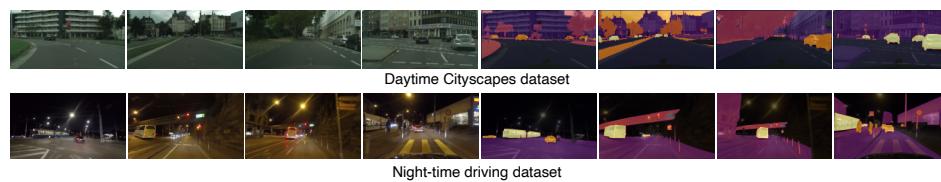
**C. Integration of WTA layer in ViT architecture**



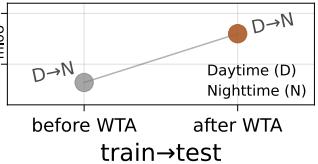
**D. Performance boost on unseen data Vision Transformer - (0-9 digits)**



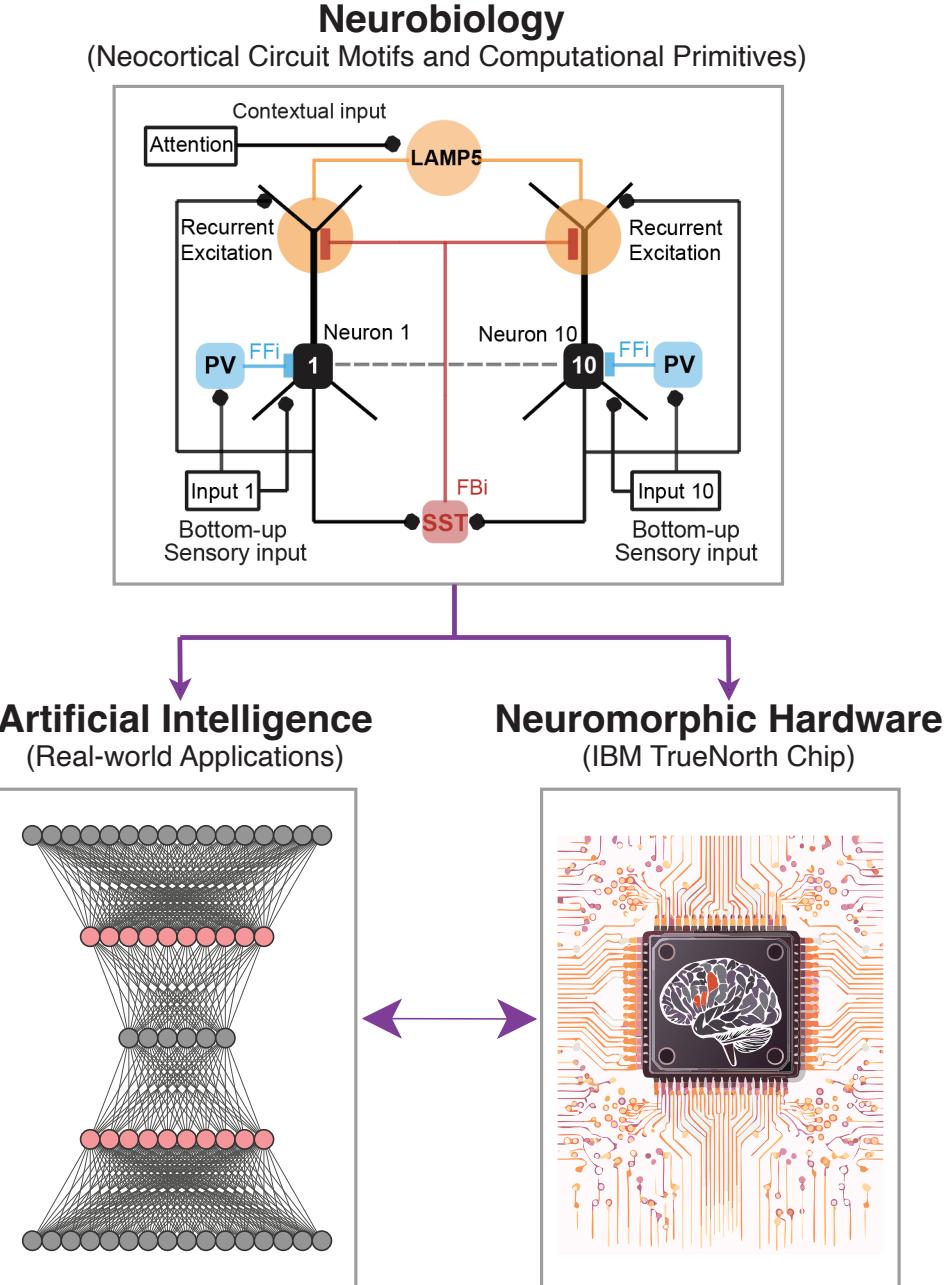
**E. Natural image segmentation samples from day and night-time driving dataset**



**F. Performance boost on unseen data RefineNet - (day/night images)**

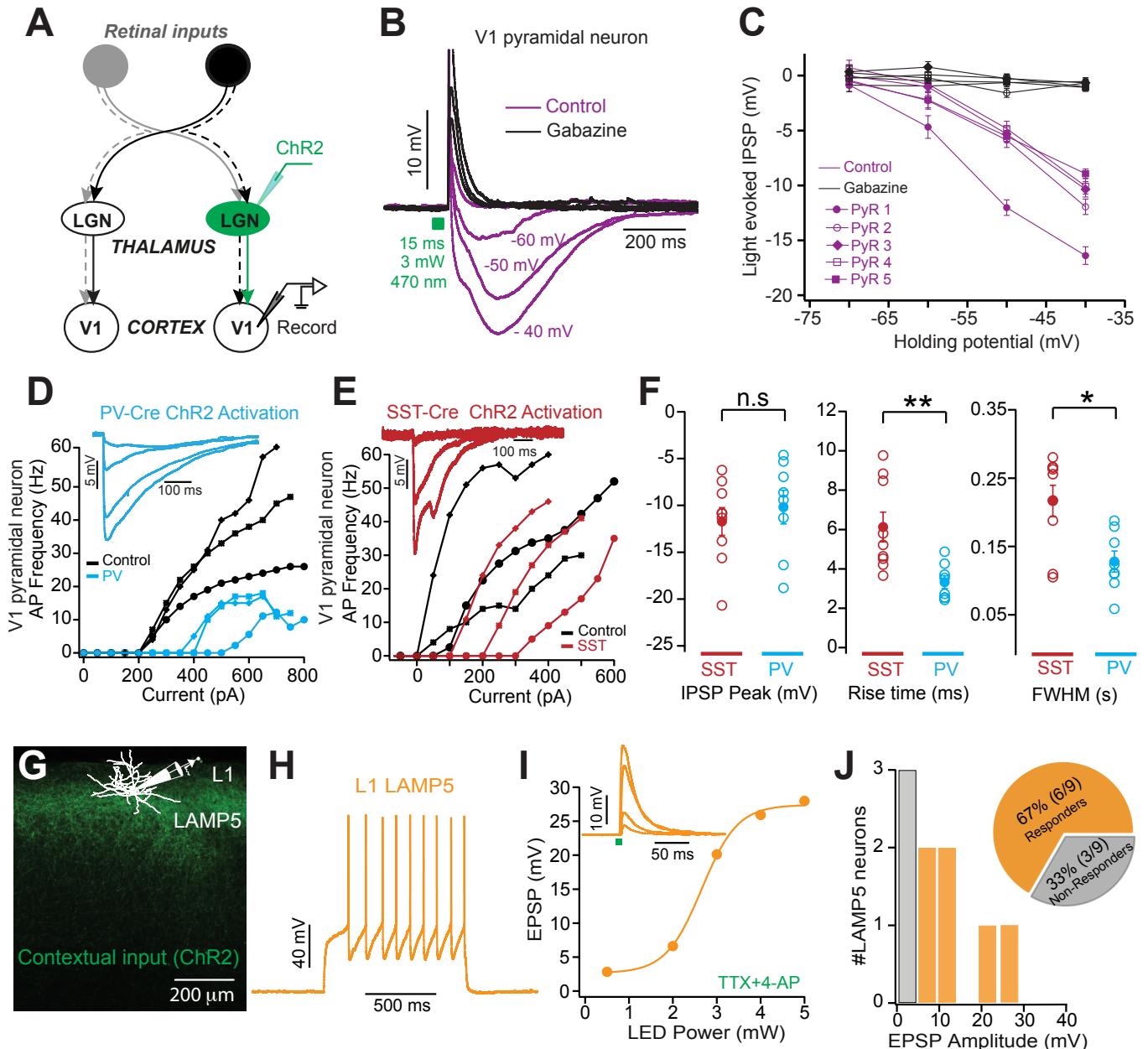


**Figure 3:** Neocortex-inspired algorithms for machine learning application and basis for zero-shot learning. **A)** The process involves a sliding window-based computation to enhance features, employing adjustable parameters that emulate the WTA implementation. This is accomplished through recurrent excitation and lateral inhibition acting across pixels, facilitating the feature enhancement mechanism. **B)** WTA layer-based image patches for ViT architecture for MNIST and digit datasets. **C-D)** Integration of WTA layer in Vision Transformer architecture for MNIST object classification task and results on training the model on source domain and testing on unseen target domains. **E)** Natural image segmentation samples from daytime and nighttime driving datasets. **F)** Performance improvement with and without adding a WTA layer in RefineNet for semantic segmentation task.



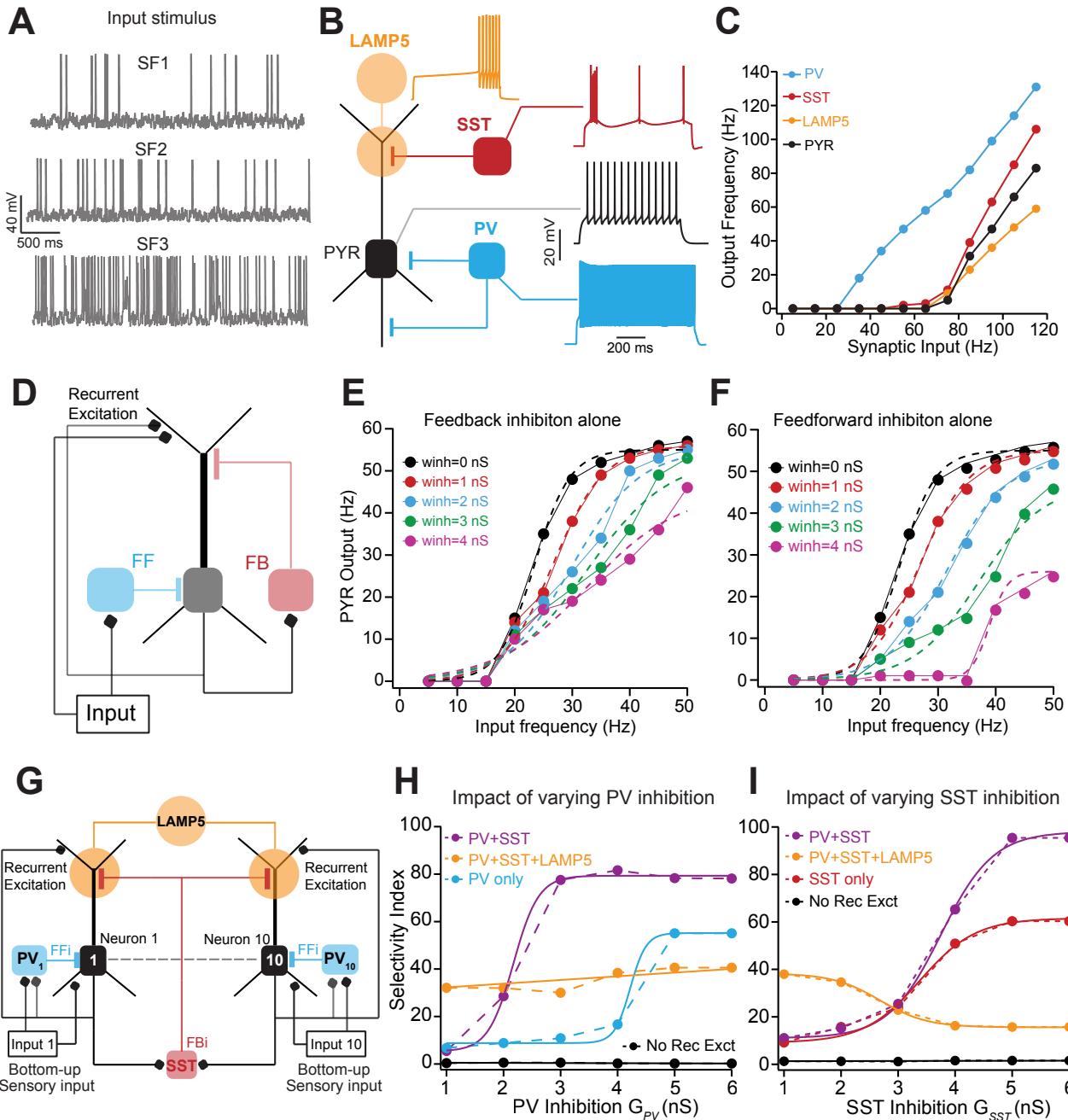
**Figure 4:** Block diagram architecture of our proposed framework. We draw inspiration from understanding the neurobiological systems of the brain and implement our findings into the neuromorphic hardware. In addition, we implement our findings as a pre-processing WTA layer in AI models and observe high-performance boosts in real-world object classification and natural image segmentation tasks.

## Experiment



**Supplementary Figure 1:** A) Schematic of the experimental arrangement showing ChR2 injections in LGN, and recordings in pyramidal neuron in V1. B) Average light-evoked synaptic responses (10 trials) in a representative layer 2/3 pyramidal neuron in V1 at the indicated holding potentials in control (magenta) and in the presence of GABAzine (black). C) Summary data showing IPSP amplitude versus holding potential in control (magenta) and in the presence of GABAzine (black) in different binocular layer 2/3 pyramidal neurons. D-E) Impact of PV and SST neuron activation on the f/I curves of L2/3 pyramidal neurons. Voltage responses of binocular pyramidal neurons at different resting membrane potentials during cre-dependent activation of PV (D) and SST (E) inhibitory neurons. F) Quantification of peak hyperpolarization (left panel), rise time (middle panel), and full-width half maximum (right panel) of the inhibition evoked by SOM and PV neuron activation in binocular pyramidal neurons at a resting membrane potential of -45 mV. G) Schematic showing contextual inputs from the Auditory cortex to layer 1 Lamp5 neurons in the somatosensory cortex (S1). H) Suprathreshold voltage response of an L1 Lamp5 interneuron to a somatic current injection of 200 pA. I) Amplitude of EPSPs in a layer 1 neuron in S1 versus LED power (470 nm, 2 ms) in the presence of TTX+4-AP. Inset: Example EPSPs during photo-activation of auditory cortex axons with increasing power (0.24 to 1 mW). J) Histogram of EPSP amplitude in layer 1 interneurons receiving (cyan) and not receiving (grey) A1 input. Inset: Pie chart of the distribution of responding (cyan) and non-responding (grey) layer 1 interneurons.

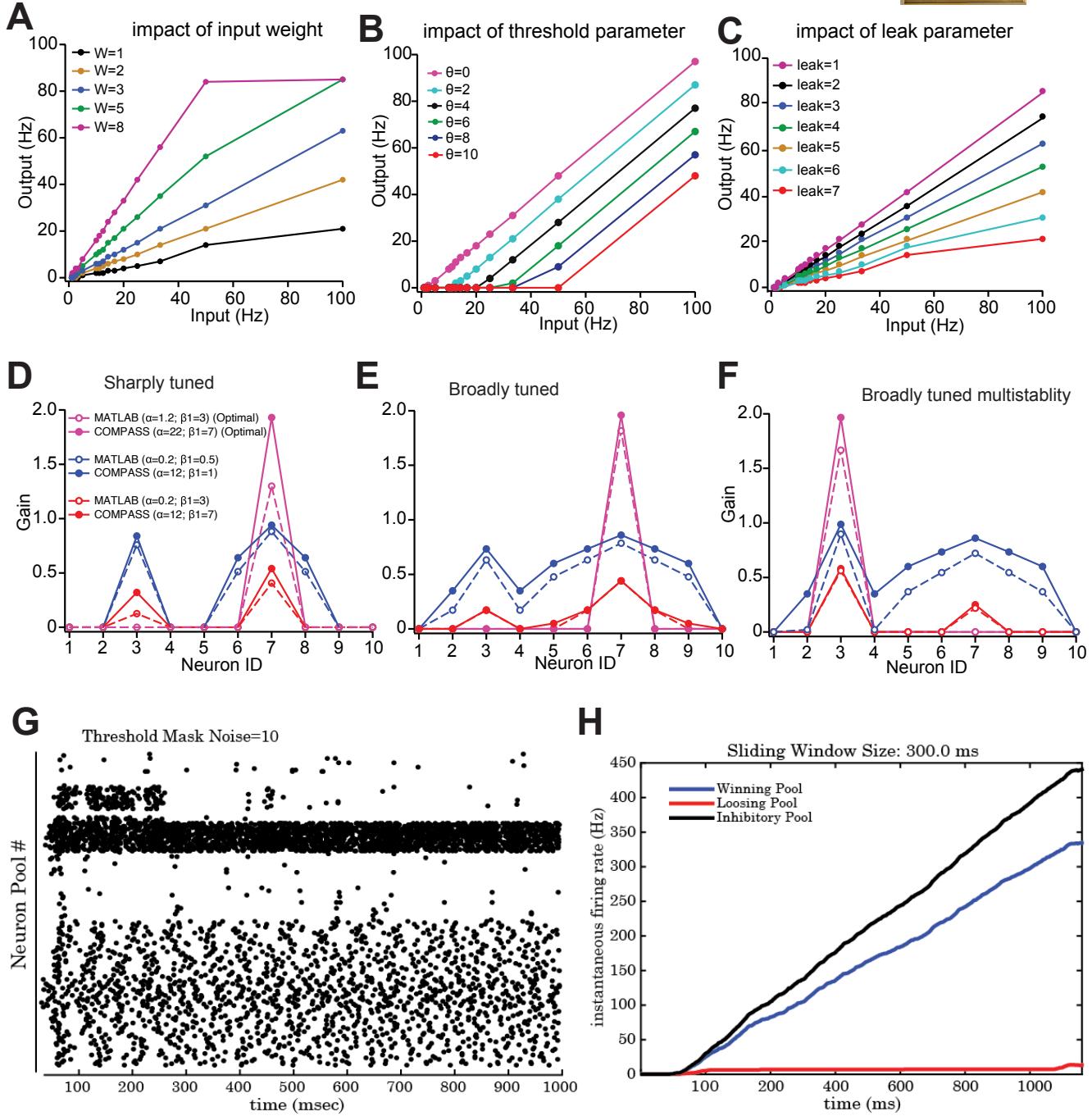
## Model



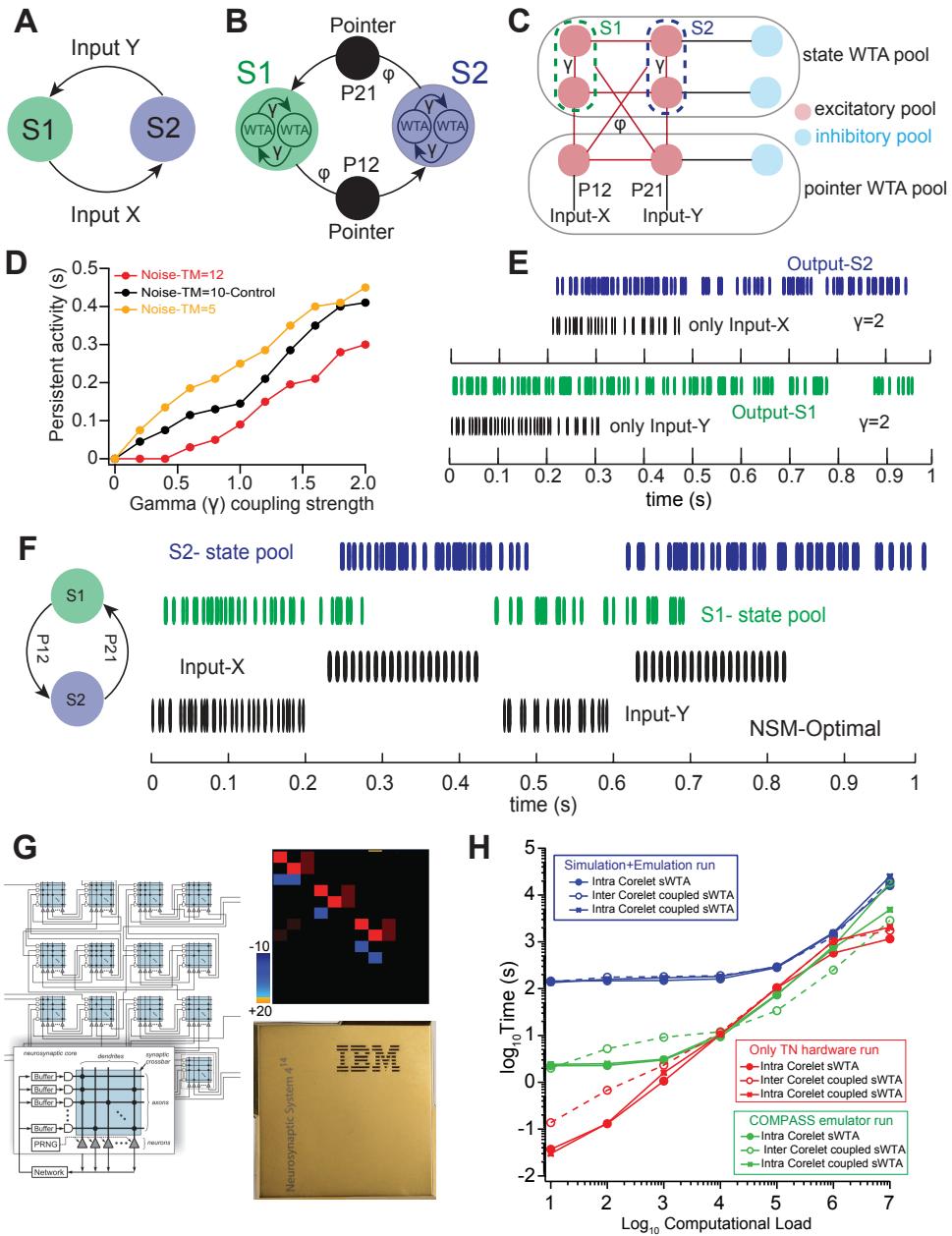
**Supplementary Figure 2:** A) Voltage traces depicting neuronal firing for Poisson-distributed synaptic stimulation at various stimulus frequencies (SF) varying from 5 to 100 Hz. B) Schematic showing inhibitory neurons preferred synaptic location onto pyramidal neurons and their action potential firing dynamics of pyramidal neuron (black) and distinct interneuron types: PV (blue), SST (red), Lamp5 (orange). C) Plot showing output frequency (FF) as a function of SF for individual neuron types as described in (B). D) Schematic showing simplified circuit organization of feedforward (FF) and feedback (FB) inhibition on pyramidal neurons. E-F) Quantification of Feedforward (E) and Feedback (F) inhibition impact on the pyramidal neuron's input-output function for different inhibitory synaptic conductance shown in various colors. Input to the network is Poisson-distributed synaptic stimulation at various SF as described in (A). G) Schematic showing the canonical cortical microcircuit motif. H-I) Impact of varying PV (H) and SST (I) inhibition on the selectivity index, which is calculated by subtracting the responses of pyramidal neurons to closely tuned inputs and multiplying the response difference with network gain. Selectivity index quantification as a function of inhibitory conductance strength for various scenarios such as in the presence of both PV, SST, and Lamp5 inhibition (orange), PV and SST only (magenta), PV only (cyan), SST only (red) and no recurrent excitation (black).



### TrueNorth neuron configuration

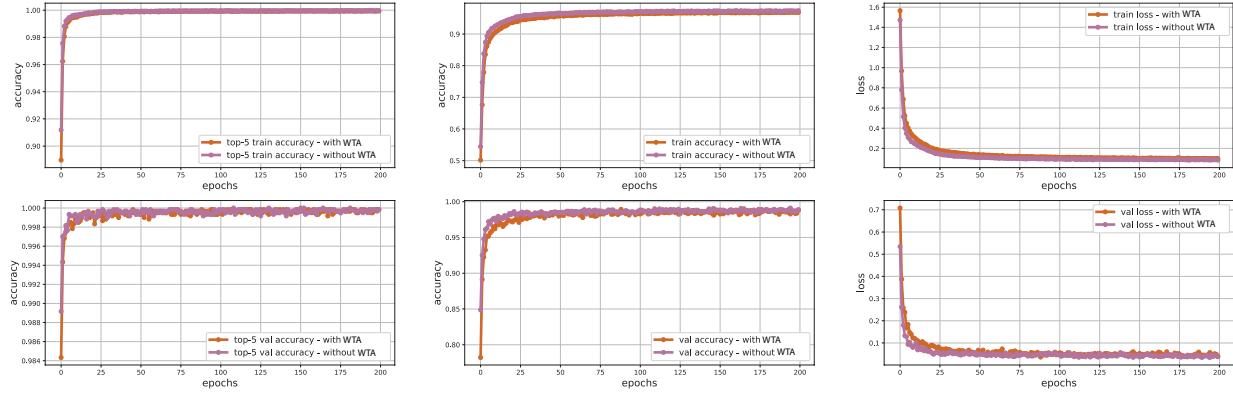


**Supplementary Figure 3:** A-C) Impact of excitatory synaptic weight (A), threshold parameter (B), and leak parameter (C) on TrueNorth neuron input-output function for various magnitudes of the tested parameter shown in the corresponding color. D-F) Comparison of network gain of the rate-based neural network model implemented in the programming platform and with the TN neural networks implemented in the Compass (TrueNorth's) emulator. For sharply tuned (D), broadly input with noise (E) and broadly tuned input with multistability (F) under optimal (magenta) and suboptimal (blue) parameter tuning. G-H) Impact of threshold mask noise parameter on WTA dynamics in TN. Spike raster plot of winning population showing WTA behavior for a threshold mask noise value of 10 (G). Population response dynamics of the winning and losing excitatory pool along with inhibitory population (H).

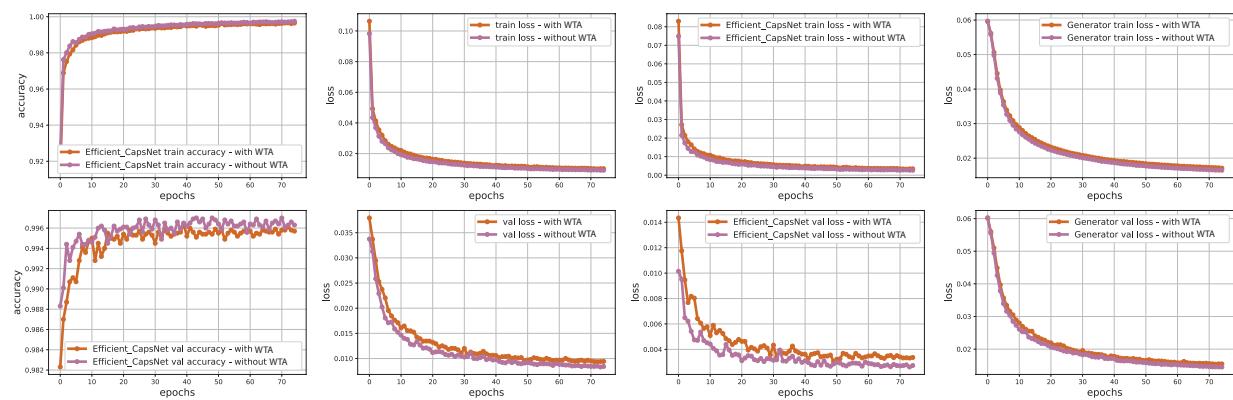


**Supplementary Figure 4:** Neural state machines implemented with coupled sWTA motifs. **A-B)** Input-dependent state transition conditions in a two-state (S1 and S2) finite state automaton (FSA) and its neuronal implementation via coupled sWTA through excitatory connection, gamma ( $\gamma$ ). The state transition is mediated through state pointer neurons. **C)** Population-level implementation of the two-state FSA as coupled sWTA in spiking-mode configuration for the two states (S1 and S2) and state-transition pointer (P12, P21) sWTA population as shown in (B) showing excitatory (red) and inhibitory (blue) neurons along with state pointer neurons (black) and state transition inputs. **D)** Impact of coupling strength ( $\gamma$ ) between sWTA and threshold mask noise (TM) is shown in different colors. **E)** Demonstration of persistent activity for state maintenance with optimal parameters in the presence of only input X or input Y. S2 output in blue and S1 output in green. **F)** Raster plot of the state pool neurons showing state transitions from S1 (green) to S2 (blue) along with inputs for optimal parameters. **G)** Heatmap of the synaptic weight matrix implementing population-level two-state FSA in TN neurons. **H)** Hardware load analysis for various core size activation runtime for coupled sWTA network and comparing it with Truenorth versus software emulation and programming platform implementation. The  $\log_{10}$  plot shows the differences in execution time as a function of computational load for hardware vs simulation. The graph shows the execution time for 1) NSCS (only TN Compass simulator) in green. 2) Native TN (only hardware) in red. 3) Total emulation+simulation run time in blue. For each case, various configuration is explored with coupled sWTA implemented using within core neurons (intra-core) vs external core neurons (inter-core).

A

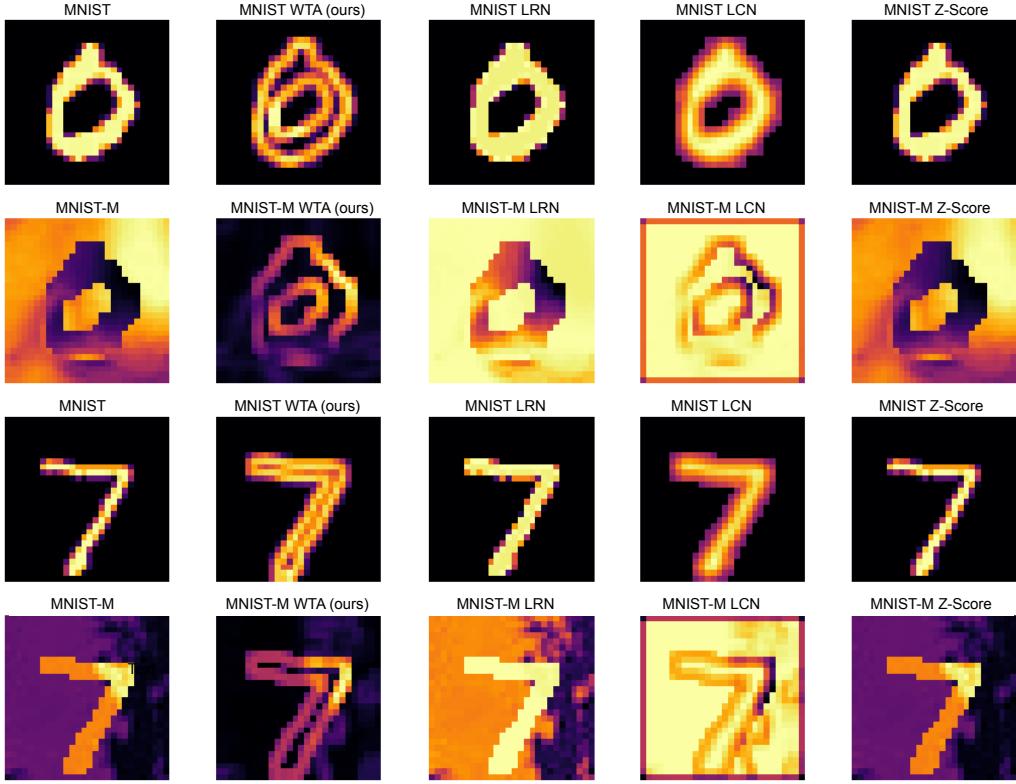


B

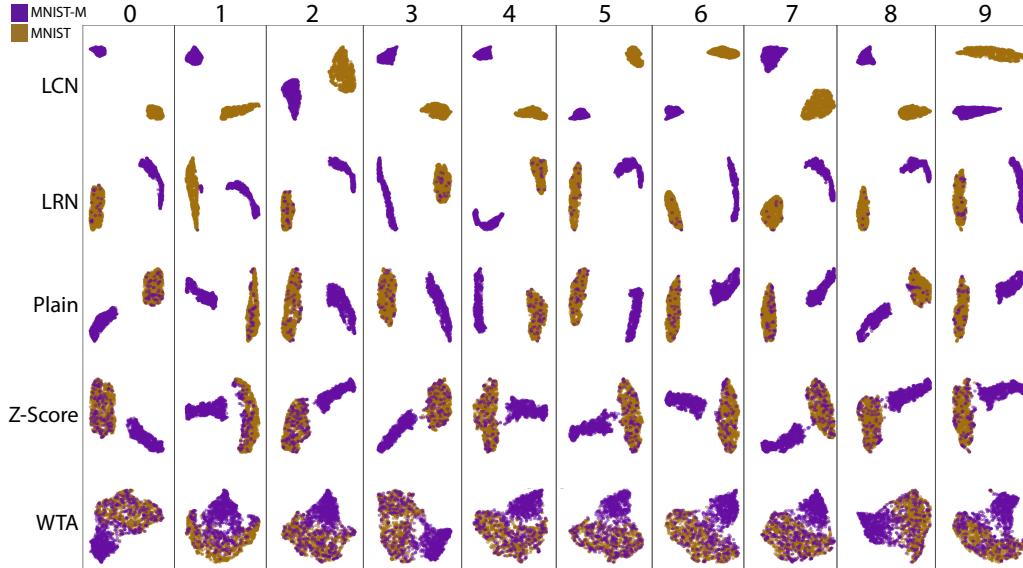


**Supplementary Figure 5:** A. Training and validation performance (top-5 accuracy, total accuracy, and loss) of ViT architecture with and without adding WTA-layer. B. Training and validation performance (accuracies and losses) of CapsuleNet architecture with and without adding WTA-layer.

**A. Processing of WTA (ours) and several image processing techniques for domain generalization on sample images of 0 and 7**



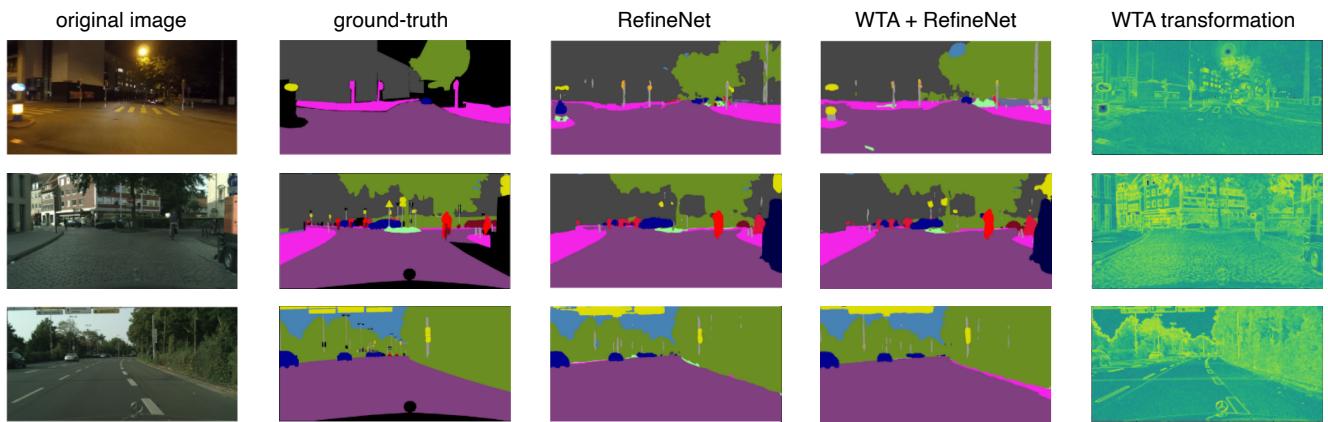
**B. UMAP visualizations of WTA (ours) and other techniques to observe domain shift robustness between clusters of datasets**



**C. Quantitative performance of WTA (ours) and other techniques on domain generalization image classification tasks**

Technique	MNIST → MNIST	MNIST → MNIST-M
LRN	0.9869	0.2103
LCN	0.988	0.1822
Z-score	0.985	0.4342
WTA (ours)	0.9849	0.70

**Supplementary Figure 6:** (A). Comparison of WTA and other normalization (LCN, LRN, Z-Score) techniques on 0 and 7 sample images from MNIST and MNIST-M. (B). UMAP embeddings of MNIST (brown) and MNIST-MM (purple) datasets, each row shows the embeddings of digits (0: leftmost - 9: rightmost) drawn from MNIST and MNIST-MM datasets for different techniques. (C). ViT is trained on MNIST and tested on both MNIST and MNIST-M datasets for comparisons between WTA and other techniques.



**Supplementary Figure 7:** Qualitative outputs of Night-time driving data set against the RefineNet trained with WTA and without WTA layer. The WTA representation is demonstrated in the last column which is fed to RefineNet for training.

**Supplementary Table 1:** Performance comparison of various DNN architectures after and before adding WTA-layer. **Bold green** shows the benchmarking results, whereas normal green shows the performance improvement after adding the WTA layer but no benchmarking. The gray shows the cases where model performance is not improved by adding a WTA layer in the architecture (best seen in color).

Models	MNIST → USPS (%)	USPS → MNIST (%)	SVHN → MNIST (%)	MNIST → MNIST-M (%)	MNIST-M → MNIST (%)
Vision Transformer	<b>84.26 / 75.2</b> $k = [3,3]$	<b>78.0 / 72.0</b> $k = [3,3]$	<b>73.6 / 66.6</b> $k = [3,3]$	70.0 / 42.0 $k = [3,3]$	98.1 / 98.3 $k = [3,3]$
EfficientNet	<b>83.5 / 77.9</b> $k = [4,4]$	74.2 / 50.7 $k = [4,4]$	69.1 / 61.0 $k = [3,3]$	48.8 / 18.8 $k = [4,4]$	96.3 / 95.0 $k = [3,3]$
CapsuleNet	94.1 / 96.4 $k = [3,3]$	<b>87.8 / 87.2</b> $k = [3,3]$	<b>75.9 / 58.1</b> $k = [4,4]$	57.2 / 22.0 $k = [4,4]$	<b>98.6 / 98.48</b> $k = [4,4]$
MobileNet	82.5 / 84.4 $k = [4,4]$	70.9 / 60.0 $k = [4,4]$	<b>73.5 / 72.2</b> $k = [4,4]$	<b>73.4 / 33.9</b> $k = [4,4]$	<b>97.8 / 97.3</b> $k = [3,3]$
ResNet	<b>82.8 / 82.5</b> $k = [3,3]$	66.0 / 58.5 $k = [3,3]$	<b>71.2 / 63.4</b> $k = [3,3]$	70.2 / 38.2 $k = [3,3]$	<b>97.8 / 97.4</b> $k = [3,3]$
	previous best: 82.2	previous best: 77.5	previous best: 70.58	previous best: 70.28	previous best: 97.8

**Supplementary Table 2:** Training settings for each model for object classification and segmentation tasks.

Architecture	Optimizer	Learning Rate	Loss Function	Epochs (saving best model by validation loss)
Transformer	Adam	0.001	crossentropy	200
CapsuleNet	Adam	0.0005	margin loss + 'MSE'	75
EfficientNet	Adam	0.001	crossentropy	100
MobileNet	Adam	0.01	crossentropy	100
ResNet	Adam	0.01	crossentropy	100
U-Net	Adam	0.0001	DICE-based loss	100
Refine Net	SGD	Decaying	crossentropy	120