Spike-Based Winner-Take-All Computation: Fundamental Limits and Order-Optimal Circuits

Lili Su

Computer Science & Artificial Intelligence Laboratory Massachusetts Institute of Technology lilisu@mit.edu

> Chia-Jung Chang Brain and Cognitive Sciences Massachusetts Institute of Technology chiajung@mit.edu

Nancy Lynch Computer Science & Artificial Intelligence Laboratory Massachusetts Institute of Technology lynch@csail.mit.edu

March 1, 2024

Abstract

Winner-Take-All (WTA) refers to the neural operation that selects a (typically small) group of neurons from a large neuron pool. It is conjectured to underlie many of the brain's fundamental computational abilities. However, not much is known about the robustness of a spike-based WTA network to the inherent randomness of the input spike trains. In this work, we consider a spike-based k-WTA model wherein nrandomly generated input spike trains compete with each other based on their underlying statistics, and k winners are supposed to be selected. We slot the time evenly with each time slot of length 1 ms, and model the n input spike trains as n independent Bernoulli processes. The Bernoulli process is a good approximation of the popular Poisson process but is more biologically relevant as it takes the refractory periods into account. Due to the randomness in the input spike trains, no circuits can guarantee to successfully select the correct winners in finite time. We focus on analytically characterizing the minimal amount of time needed so that a target minimax decision accuracy (success probability) can be reached.

We first derive an information-theoretic lower bound on the decision time. We show that to have a (minimax) decision error $\leq \delta$ (where $\delta \in (0, 1)$), the computation time of any WTA circuit is at least

 $((1-\delta)\log(k(n-k)+1)-1)T_{\mathcal{R}},$

where $T_{\mathcal{R}}$ is a difficulty parameter of a WTA task that is independent of δ , n, and k. We then design a simple WTA circuit whose decision time is

$$O\left(\left(\log\left(\frac{1}{\delta}\right) + \log k(n-k)\right)T_{\mathcal{R}}\right)$$

It turns out that for any fixed $\delta \in (0, 1)$, this decision time is order-optimal in terms of its scaling in n, k, and $T_{\mathcal{R}}$.

1 Introduction

Humans and animals can form a stable perception and make robust judgments under ambiguous conditions. For example, we can easily recognize a dog in a picture regardless of its posture, hair color, and whether it stands in the shadow or is occluded by other objects. One fundamental feature of brain computation is its robustness to the randomness introduced at different stages, such as sensory representations [KK01, HW59], feature integration [KTA⁺03, MCM07], decision formation [PG99, SN01], and motor planning [HW98, LCG⁺15]. It has been shown that neurons encode information in a stochastic manner in the brain [BAB⁺97, KRR00, MA09, FDMM18]; even when the exact same sensory stimulus is presented or when the same kinematics are achieved, no deterministic patterns in the spike trains exist. Facing environmental ambiguity, humans and animals adaptively refine their behaviors by incorporating prior knowledge with their current sensory measurements [FSW08, KP04, SS06, EB02, KW04]. Nevertheless, it remains relatively unclear how neurons carry out robust computation facing ambiguity. Sparse coding is a common strategy in brain computation; to encode a task-relevant variable, often only a small group of neurons from a large neuron pool are activated [OF04, POMT⁺02, HDZ08, QKKF08, KF08, RPG99]. Understanding the underlying neuron selection mechanism is highly challenging.

Winner-Take-All (WTA) is a hypothesized mechanism to select proper neurons from a competitive network of neurons, and is conjectured to be a fundamental primitive of cognitive functions such as attention and object recognition [RP99, IKN98, YG98, Maa00]. Among these studies, it is commonly assumed that neurons transmit information with a continuous variable such as the firing rate. This assumption, however, ignores how temporal coding may additionally contribute to cortical computations. For example, some neurons in the auditory cortex will respond to auditory events with bursts at a fixed latency [GKvHW96, Nel04]. This phaselocking property is also observed in the hippocampus as well as the prefrontal cortex [SLW05, HSM06, BC95]. Another feature that has been neglected in a rate-based model is the inherent noise in the inputs. Although some studies used additive Gaussian noise [KCF17, LLW13, LIKB99, RV06] to account for input randomness, such WTA circuits are very sensitive to noise and could not successfully select even a single winner unless extra robustness strategy such as an additional nonlinearity is introduced into the dynamics [KCF17]. Last but not least, neurons have a refractory period, which prevents spikes from back propagating in axons [BIM98], and such a feature is usually neglected in the rate-based models. In contrast, a spike-based model may capture these neglected features. Nevertheless, how WTA computation can be implemented and its algorithmic characterization remains relatively under-explored.

In this paper, we study a spike-based k-WTA model wherein n randomly generated input spike trains are competing with each other with their underlying statistics, and the true winners are the k input spike trains whose underlying statistics are higher than others. More precisely, we slot the time evenly with each time slot of length 1 ms. We assume that these n input spike trains are generated by n independent Bernoulli processes with different rates. An abstract example is depicted in Figure 1. We use Bernoulli processes to capture the randomness in the input spike trains rather than using the popular Poisson processes because a Bernoulli process can be viewed as the time-slotted version of a refractory-period-modified Poisson process; it is well-known that due to the existence of refractory periods, a neuron cannot spike twice within 1 ms.

We focus on analytically characterizing the minimal amount of time needed so that a target minimax decision accuracy (success probability) can be reached. We first derive a lower bound on the decision time for a given decision accuracy. We show that no WTA circuits can have a computation time strictly less than

$$((1-\delta)\log(k(n-k)+1)-1)T_{\mathcal{R}},$$
(1)

where $T_{\mathcal{R}}$ is a parameter defining the difficulty to distinguish between two spike trains with different statistics in a WTA task, n is the number of input spike trains, k is the number of winners, and δ is the given target decision accuracy. In many practical settings we care about the sparse coding region where $k \ll n$. Our lower-bound is obtained by an information-theoretic argument, and holds for all WTA circuits without restricting their circuit architectures and their adopted activation functions. Throughout this paper, we are interested in the decision time's scaling in n, k, and $T_{\mathcal{R}}$, while treating $\delta \in (0, 1)$ as a fixed small constant. Not surprisingly, the above lower bound grows with the network size n when other parameters are fixed. This is because the larger n, the noisier the WTA competition. Similarly, when n and k are fixed, the easier to distinguish two spike trains with different statistics (i.e., the smaller $T_{\mathcal{R}}$), the shorter the necessary decision time is.

We construct a simple circuit whose decision time is

$$O\left(\left(\log\left(\frac{1}{\delta}\right) + \log k(n-k)\right)T_{\mathcal{R}}\right).$$



Figure 1: In this figure, n randomly generated input spike trains are fed to a WTA circuit as the circuit input. Clearly, no deterministic patterns can be read off from these spike trains. Here, we do not specify the output of a WTA circuit because the detailed specifications of the circuits' outputs might vary with the corresponding applications.

It turns out that for any fixed $\delta \in (0,1)$, this decision time is order-optimal in terms of its scaling in n, k, and $T_{\mathcal{R}}$, i.e., its decision time matches the lower-bound in (1) up to a constant multiplicative factor. In our circuit, each output neuron is a thresholded accumulator unit whose threshold is determined by δ , k, n, and $T_{\mathcal{R}}$, and the circuit's output is the first group of k output neurons that spike in the same time. The typical dynamics under our circuit are: the number of output neurons that spike simultaneously (i.e., spike at the same time) is monotonically increasing until exactly k output neurons spike simultaneously. The simultaneous spikes of these k output neurons cause strong inhibition of other output neurons; in particular, no other output neuron can spike within a sufficiently long period $\Omega\left(\left(\log\left(\frac{1}{\delta}\right) + \log k(n-k)\right)T_{\mathcal{R}}\right)$.

In addition, our results also give a set of testable hypotheses on neural recordings and humans'/animals' behaviors in decision-making. For instance, given the number of input spike trains and the number of true winners, our results can provide an estimate of the minimum decision time needed, which can provide some insights on the efficiency of a WTA circuit in terms of decision time. As another example, when two animals are involved in the same experiment, if both animals reach the same accuracy in discriminating two objects, does the animal that decides faster have more heterogeneous distributions of input spiking activities, i.e., smaller $T_{\mathcal{R}}$? Our results provide partial answers to this question.

2 Computational Model: Spiking Neuron Networks

In this section, we provide a general description of the computation model used. There is much freedom in choosing the detailed specification of the model. In particular, in Section 5 we provide a circuit construction (for solving the k-WTA competition) under this computation model.

2.1 Network Structure

A spiking neuron network (SNN) $\mathcal{N} = (U, E)$ consists of a collection of neurons U that are connected through synapses E. We assume that a SNN can be conceptually partitioned into three non-overlapping layers: *input* layer N_{in} , hidden layer N_h , and output layer N_{out} ; the neurons in each of these layers are referred to as *input* neurons, hidden neurons, and output neurons, respectively. The synapses E are essentially directed edges, i.e., $E := \{(\nu, \nu') : \nu, \nu' \in U\}$. For each $\nu \in U$, define $\mathsf{PRE}_{\nu} := \{\nu' : (\nu', \nu) \in E\}$ and $\mathsf{POST}_{\nu} := \{\nu' : (\nu, \nu') \in E\}$. Intuitively, PRE_{ν} is the collection of neurons that can directly influence neuron ν ; similarly, POST_{ν} is the collection of neurons that can be directly influenced by neuron ν .¹ We assume that the input neurons cannot be influenced by other neurons in the network, i.e., $\mathsf{PRE}_{\nu} = \emptyset$ for all $\nu \in N_{in}$. Each edge (ν, ν') in E has a weight, denoted by $w(\nu, \nu')$. The strength of the interaction between neuron ν and neuron ν' is captured

 $^{^{1}}$ In the languages of computational neuroscience, the incoming neighbors and outgoing neighbors are often referred to as pre-synaptic units and post-synaptic units.



Figure 2: A SNN consists of three layers: the input layer, the output layer, and the hidden layer. The hidden neurons might connect to both the input neurons and the output neurons to assist the computation of the neuron network. Neurons are connected through synapses. WTA circuits is a family of SNNs in which the number of output neurons equals the number of the input neurons.

as $|w(\nu,\nu')|$. The sign of $w(\nu,\nu')$ indicates whether neuron ν excites or inhibits neuron ν' : In particular, if neuron ν excites neuron ν' , then $w(\nu,\nu') > 0$; if neuron ν inhibits neuron ν' , then $w(\nu,\nu') < 0$. The set E might contain self-loops with $w(\nu,\nu)$ capturing the self-excitatory/self-inhibitory effects. An example of SNNs can be found in Figure 2.

Generic network structure for WTA circuits The family of WTA circuits under consideration is rather generic. We only assume that $|N_{in}| = |N_{out}| = n$ the numbers of the input neurons and of the output neurons are equal. For ease of exposition, denote

$$N_{in} = \{u_1, \cdots, u_n\}, \text{ and } N_{out} = \{v_1, \cdots, v_n\}.$$

The hidden neuron subset N_h can be arbitrary. The output neurons and the hidden neurons may be connected to each other in an arbitrary manner.

2.2 Network State

In a SNN, the communication among neurons is abstracted as spikes. We assume each neuron ν has two local variables: *spiking state* variable $S(\nu)$ and *memory state* variable $M(\nu)$. Nevertheless, for input neurons, we only consider their spiking states, assuming that their memory states are not influenced by the dynamics of the spiking neuron network under consideration. We slot the time evenly with each time slot of length 1 ms. Let $t = 1, 2, \cdots$ be the indices of the time slots. Henceforth, by saying time t, we mean the time interval [t-1,t) ms. For $t \ge 1$, let $S_t(\nu) \in \{0,1\}$ be the spiking state of neuron ν at time t indicating whether neuron ν spikes at time t or not. By convention, $S_0(\nu) := 0$. For a non-input neuron ν and for $t \ge 1$, let $M_t(\nu)$ be the memory state of neuron ν at time t summarizing the cumulative influence caused by the spikes of the neurons in PRE_i during the most recent m times, i.e., times $t - 1, t - 2, \cdots, t - m$. Concretely, let $V_t(\nu)$ be the charge of (non-input) neuron ν at time t (for $t \ge 1$) defined as

$$V_t(\nu) := \sum_{\nu' \in \mathsf{PRE}_\nu} w(\nu',\nu) S_t(\nu')$$

Clearly, $V_0(\nu) = 0$. Let V_t^{ν} be the sequence of length m such that

$$\boldsymbol{V}_t^{\nu} := [V_t(\nu), \cdots, V_{t-m+1}(\nu)]$$

and let $S_t(\nu)$ be the sequence of length m such that

$$\boldsymbol{S}_t^{\nu} := [S_t(\nu), \cdots, S_{t-m+1}(\nu)]$$

By convention, when $0 \le t \le m$, let

$$V_t^{\nu} := [V_t(\nu), \cdots, V_0(\nu), 0, \cdots, 0]$$

and

$$S_t^{\nu} := [S_t(\nu), \cdots, S_0(\nu), 0, \cdots, 0]$$

For $t \geq 1$, define the memory variable $M_t(\nu)$ as a pair of vectors S_{t-1}^{ν} and V_{t-1}^{ν} , i.e.,

$$M_t(\nu) := (S_{t-1}^{\nu}, V_{t-1}^{\nu})$$

By convention, let $M_0(\nu) := (0, 0)$, where **0** is the length *m* zero vector.

At time t + 1, the memory variable $M_{t+1}(\nu)$ is updated by shifting the two sequences forwards by one time unit – fetching in $S_t(\nu)$ and $V_t(\nu)$, respectively, and removing $S_{t-m}(\nu)$ and $V_{t-m}(\nu)$, respectively. The memory state $M_t(\nu)$ is known to neuron ν only, and it can influence the probability of generating a spike at time t through an activation function ϕ_{ν} , i.e.,

$$S_t(\nu) = \phi_{\nu} \left(M_t(\nu) \right), \forall t \ge 0.$$
⁽²⁾

Notably, ϕ_{ν} might be a random function.

In most neurons, the synaptic plasticity time window is about 80 -120 msec, but could also vary across brain regions, and vary across different time scales under different behavioral contexts. In a sense, the synaptic plasticity time window is closely related to m. As can be seen in Section 5, our order-optimal WTA circuit construction requires m to be sufficiently high. Nevertheless, this does not exclude the application of our WTA circuit to the contexts where m is small. This is because the memory variable can be implemented by a chain of hidden neurons near neuron ν . The detailed implementation of the local memory does not affect the order optimality of our WTA circuit.

3 Minimax Decision Accuracy/Success Probability

3.1 Random Input Spike Trains

We study the k-WTA model, wherein n randomly generated input spike trains are competing with each other, and, as a result of this competition, k out of them are selected to be the winners. In contrast, most existing works [VRP⁺18, Maa97, LMP16] assume deterministic input spike trains.

Recall that time is slotted into intervals of length 1 ms. We assume that the *n* input spike trains are generated from *n* independent Bernoulli processes with unknown parameters p_1, \dots, p_n , respectively. We refer to $\mathbf{p} = [p_1, \dots, p_n]$ as a *rate assignment* of the WTA competition. For example, suppose there are 2 input spike trains with rates 0.6 and 0.8, respectively, i.e., n = 2 and $\mathbf{p} = [0.6, 0.8]$. In each time, with probability 0.6 the first input spike train has a spike independently from whether the second input spike train has a spike or not; similarly for the second input spike train.

Note that in the most general scenario the spikes of the input neurons might be correlated; see Section 6 for detailed comments. We would like to explore the more general input spikes in our future work.

3.2 Minimax Performance Metric

We adopt the minimax framework [Wu17] (in which the circuit designer and nature play games against each other) to evaluate the performance (decision accuracy versus decision time) of a WTA circuit.

Let $\mathcal{R} \subseteq [c, C]$ be an arbitrary but finite set of rates where c and C are two absolute constants such that 0 < c < C < 1. A rate assignment p is chosen by nature from \mathcal{R}^n for which there exists a subset of $[n] := \{1, \dots, n\}$, denoted by $\mathcal{W}(p)$, such that

$$|\mathcal{W}(\boldsymbol{p})| = k, \text{ and } p_i > p_j \quad \forall i \in \mathcal{W}(\boldsymbol{p}), j \notin \mathcal{W}(\boldsymbol{p})$$
(3)

- recall that $|\cdot|$ is the cardinality of a set. We refer to set $\mathcal{W}(\mathbf{p})$ as the true winners with respect to the rate assignment \mathbf{p} . For example, suppose n = 5, k = 2, and

$$p = [p_1 = 0.2, p_2 = 0.1, p_3 = 0.2, p_4 = 0.8, p_5 = 0.85].$$

Here the true winners are 4 and 5, i.e., $\mathcal{W}(\mathbf{p}) = \{4, 5\}$. In this paper, we consider the following collection of rate assignments, denoted by \mathcal{AR} :

$$\mathcal{AR} := \{ \boldsymbol{p} : \boldsymbol{p} \in \mathcal{R}^n, \& \exists \mathcal{W}(\boldsymbol{p}) \subseteq [n] \text{ s.t. } |\mathcal{W}(\boldsymbol{p})| = k, \text{ and } p_i > p_j \forall i \in \mathcal{W}(\boldsymbol{p}), j \notin \mathcal{W}(\boldsymbol{p}) \}.$$
(4)

For each of reference, we refer to an element in \mathcal{AR} as an admissible rate assignment. Recall that the input of a WTA circuit is a collection of n independent spike trains. For a given rate assignment p, let $\{S_t(u_i)\}_{t=1}^T$ denote the spike train of length T at input neuron u_i . The circuit designer wants to design a WTA circuit that outputs a good guess/estimate \widehat{win} of $\mathcal{W}(p)$ for any choice of rate assignment p in \mathcal{AR} . Note that conditioning on

$$\boldsymbol{S} := \left[\{ S_t(u_1) \}_{t=1}^T, \cdots, \{ S_t(u_n) \}_{t=1}^T \right],$$

the estimate \widehat{win} is independent of p. Here S is used with a little abuse of notation as this notation hides its connection with T and the rate parameter p.² Later, we use the same notation to denote the n spike trains with random rate assignment, i.e., where p is randomly generated. Nevertheless, this abuse of notation significantly simplifies the exposition without sacrificing clarity. In particular, we will specify the underlying rate assignment when it is not clear from the context.

Under minimax framework, we are interested in the minimax error probability ³

$$\min_{\widehat{\boldsymbol{win}}} \max_{\boldsymbol{p} \in \mathcal{AR}} \mathbb{P}\left\{\widehat{\boldsymbol{win}}\left(\boldsymbol{S}\right) \neq \mathcal{W}(\boldsymbol{p})\right\}.$$
(5)

For a given deterministic WTA circuit \widehat{win} (i.e., the activation functions used are deterministic), the probability in $\mathbb{P}\left\{\widehat{win}(S) \neq \mathcal{W}(p)\right\}$ is taken w.r.t. the randomness in the stochastic spikes of each input neuron; for a randomized WTA circuit \widehat{win} (i.e., the activation functions are stochastic), in addition to the aforementioned source of randomness, the probability in $\mathbb{P}\left\{\widehat{win}(S) \neq \mathcal{W}(p)\right\}$ is also taken w.r.t. the randomness in the activation functions. In (5), the performance metric of a WTA circuit is the worst-case error probability

$$\max_{oldsymbol{p}\in\mathcal{AR}}\mathbb{P}\left\{\widehat{oldsymbol{win}}\left(oldsymbol{S}
ight)
eq\mathcal{W}(oldsymbol{p})
ight\}.$$

Essentially, the statistical inference problem can be viewed as a game between the circuit designer and nature.

4 Information-Theoretic Lower Bound on Decision Time

In this section, we provide a lower bound on the decision time for a given decision accuracy. The lower bounds derived in this section hold universally for all possible network structures (including the hidden layer), synapse weights, and the activation functions.

One observation is that the decision time is naturally lower bounded by the sample complexity, which is closely related to the Kullback-Leibler (KL) divergence⁴ between two Bernoulli distributions. The KL divergence between Bernoulli random variables with parameters r and r', respectively, is defined as

$$d(r \parallel r') := r \log\left(\frac{r}{r'}\right) + (1-r) \log\left(\frac{1-r}{1-r'}\right),$$
(6)

²A more rigorous notation should be $\mathbf{S}(T, \mathbf{p}) := \left[\{S_t(u_1)\}_{t=1}^T, \cdots, \{S_t(u_n)\}_{t=1}^T \right]$. We use \mathbf{S} for $\mathbf{S}(T, \mathbf{p})$ for ease of exposition. ³In the following expression, the min should really be an inf, but we abuse notation here for ease of exposition. In addition, the max really is a max as the set \mathcal{R} under consideration is of finite size.

⁴The Kullback-Leibler (KL) divergence gauges the **dissimilarity** between two distributions.

where, by convention, $\log \frac{0}{0} := 0$. Notably, $d(\cdot \| \cdot)$ is not symmetric in r and r'. In addition, when $r \neq 0$ and r' = 0 or 1, $d(r \| r') = \infty$. Recall that set \mathcal{R} is an arbitrary but finite set that are contained in the interval [c, C], where $c, C \in (0, 1)$. It holds that $d(r \| r') < \infty$ for all $r, r' \in \mathcal{R}$. For the more general distributions over a common discrete alphabet \mathcal{A} , say distributions P and Q, the Kullback-Leibler (KL) divergence between P and Q is defined as follows.

Definition 1 (KL-divergence). Let \mathcal{A} be a discrete alphabet (finite or countably infinite), and P and Q be two distributions over \mathcal{A} . Then define

$$D(P \parallel Q) := \sum_{a \in \mathcal{A}} P(a) \log \left(\frac{P(a)}{Q(a)}\right),$$

where $0 \cdot \log\left(\frac{0}{0}\right) = 0$ by convention.

Note that $D(P \parallel Q) \ge 0$ and $D(P \parallel Q) = 0$ if and only if P = Q almost surely. Similar to $d(\cdot \parallel \cdot)$, $D(P \parallel Q)$ is not symmetric in P and Q. In this paper, we choose the base to be 2. ⁵ Recall that the set of admissible rate assignments \mathcal{AR} is defined in (4).

Lemma 2. Fix a finite set \mathcal{R} . Let $\mathbf{p} = [p_1, \dots, p_n]$ and $\mathbf{q} = [q_1, \dots, q_n]$ be two rate assignments in \mathcal{AR} . Let $P_{\mathbf{S}}$ and $Q_{\mathbf{S}}$ be the distributions of the *n* spike sequences of the input neurons under rate assignments \mathbf{p} and \mathbf{q} , respectively. Then,

$$D(P_{\mathbf{S}} \parallel Q_{\mathbf{S}}) = T \sum_{i=1}^{n} d(p_i \parallel q_i)$$

Lemma 2 is proved in Appendix B.

For the given \mathcal{R} , define task complexity $T_{\mathcal{R}}$ as

$$T_{\mathcal{R}} := \max_{r_1, r_2 \in \mathcal{R} \text{ s.t. } r_1 \neq r_2} \frac{1}{d(r_2 \parallel r_1) + d(r_1 \parallel r_2)}.$$
(7)

It is closely related to the smallest KL divergence between two distinct statistics in \mathcal{R} . The task complexity $T_{\mathcal{R}}$ kicks in due to the adoption of minimax decision framework (5).

The following lemma is used in the proof of our information-theoretic lower bound. This is a technical supporting lemma, and the choice of the specific rate assignments is due to some technical convenience in proving Theorem 4.

Lemma 3. For any finite set \mathcal{R} , let $r_1, r_2 \in \mathcal{R}$ such that $r_1 \neq r_2$. Let $p^0 = [p_1^0, \cdots, p_n^0]$ be

$$p_{\ell}^{0} = \begin{cases} r_{1}, & \text{if } \ell = 1, \cdots, k; \\ r_{2}, & \text{otherwise.} \end{cases}$$

$$\tag{8}$$

For $i = 1, \cdots, k$ and $j = k + 1, \cdots, n$, define rate assignment p^{ij} as

$$p_{\ell}^{ij} = \begin{cases} p_{\ell}^{0}, & \text{ if } \ell \neq i, \neq j; \\ p_{j}^{0}, & \text{ if } \ell = i; \\ p_{i}^{0}, & \text{ if } \ell = j. \end{cases}$$

Let $X_{\mathbf{p}}$ be a random rate assignment. If $X_{\mathbf{p}}$ is uniformly distributed over

$$\{p^0\} \cup \{p^{ij}: i = 1, \cdots, k, \& j = k+1, \cdots, n\},\$$

then the mutual information $I(X_p; S)$ satisfies the following:

$$I(X_{p}; S) \leq T (d(r_{2} || r_{1}) + d(r_{1} || r_{2})).$$

⁵Note that any base would work, see [PW14, Chapter 1.1].

See Appendix A for definition of $I(\cdot; \cdot)$. The proof of Lemma 3 can be found in Appendix C.

It turns out that if the input spike train length T is not sufficiently large (specified in Theorem 4), no matter how elegant the design of a WTA circuit is (no matter which activation function we choose, how many hidden neurons we use, and how we connect the hidden neurons and output neurons), its actual decision accuracy is always lower than the target decision accuracy $1 - \delta$.

Theorem 4. For any $1 \le k \le n-1$ and any set \mathcal{R} and any $\delta \in (0,1)$, if

$$T \le \left((1-\delta) \log(k(n-k)+1) - 1 \right) T_{\mathcal{R}},$$

then

$$\min_{\widehat{win}} \max_{\boldsymbol{p} \in \mathcal{AR}} \mathbb{P}\left\{ \widehat{win} \left(\boldsymbol{S} \right) \neq \mathcal{W}(\boldsymbol{p}) \right\} \geq \delta,$$

where the min is taken over all possible WTA circuits with different choices of activation functions and circuit architectures.

Theorem 4 says that if $T < ((1-\delta)\log(k(n-k)+1)-1)T_{\mathcal{R}}$, the worst case probability error of any WTA circuit is greater than δ , i.e., $\max_{\boldsymbol{p}\in\mathcal{AR}} \mathbb{P}\left\{\widehat{\boldsymbol{win}}(\boldsymbol{S})\neq\mathcal{W}(\boldsymbol{p})\right\} > \delta$. Theorem 4 is proved in Appendix E.

Remark 5 (Tightness of the lower bound in Theorem 4). Following our line of argument, by considering a richer family of critical rate assignments in Lemma 3, we might be able to obtain a tighter lower bound. Nevertheless, the constructed WTA circuit in Section 5 turn out to be order-optimal – its decision time matches the lower bound in Theorem 4 up to a multiplicative constant factor. This immediately implies that the lower bound obtained in Theorem 4 is tight up to a multiplicative constant factor.

5 Order-Optimal WTA Circuits

In Section 2.1, we provided a general description of the computation model we are interested in. In this section, we construct a specific WTA circuit under this general computation model. This WTA circuit turns out to be order-optimal in terms of decision time – the decision time of our WTA circuit matches the lower bound in Section 4 up to a multiplicative constant factor. To do that, we need to specify (1) the network structure, including the number of hidden neurons, the collection of synapses (directed communication links) between neurons, and the weights of these synapses; (2) the memorization capability of each neuron, i.e., the magnitude of m; and (3) ϕ_{ν} – the activation function used by neuron ν .

5.1 Circuit Design

In our designed circuit, there are four parameters \mathcal{R} , m, b, and δ , where $\mathcal{R} \subseteq [c, C]^{-6}$ is a finite set from which the p_i 's of the input spike trains are chosen, m is the memory range and b is the bias at the non-input neurons, and $(1 - \delta)$ is the target decision accuracy (i.e., success probability). Here, we assume that every non-input neuron has the same bias, i.e., $b_{\nu} = b$ for all non-input neurons ν . The four parameters \mathcal{R} , m, b, and δ can be viewed as some prior knowledge of the WTA circuit; they might be learned through some unknown network development procedure which is outside the scope of this work. In Sections 5.1.1, 5.1.3, and 5.1.4, we present the network structure and the activation functions adopted, and the requirement on m. For completeness, we specify the local memory update (in particular the vector \mathbf{V}) separately in Section 5.1.2. The dynamics of our WTA circuit is summarized in Section 5.1.5.

5.1.1 Network structure:

We propose a WTA circuit with the following network structure:

⁶Recall that $c, C \in (0, 1)$ are two absolute constants, i.e., they do not change with other parameters of the WTA circuit such as n, k, and δ .

- All output neurons are connected to each other by a complete graph. That is, $(v_i, v_j) \in E$ for all $v_i, v_j \in N_{out}$ such that $v_i \neq v_j$;
- Each edge from an input neuron to an output neuron has weight 1, i.e., $w(u_i, v_i) = 1$ for all $u_i \in N_h, v_i \in N_{out}$.
- All edges among the output neurons have weights $-\frac{1}{k}$. That is, $w(v_i, v_j) = -\frac{1}{k}$ for all $v_i, v_j \in N_{out}$ such that $v_i \neq v_j$.
- There are no hidden neurons, i.e., $N_h = \emptyset$;

5.1.2 Update local charge vector:

With the above choice of network structure, the charge $V_{t-1}(v_i)$ at the output neuron v_i at time t-1 is

$$V_{t-1}(v_i) = S_{t-1}(u_i) - \frac{1}{k} \sum_{j:1 \le j \le n, \& \ j \ne i} S_{t-1}(v_j).$$

Notably, $V_{t-1}(v_i) \in [-1, 1]$ for all $t \ge 1$ and output neuron v_i . When k = 1, the above update becomes

$$V_{t-1}(v_i) = S_{t-1}(u_i) - \sum_{j:1 \le j \le n, \& \ j \ne i} S_{t-1}(v_j).$$

which can be viewed as a spike model counterpart of the potential update under the traditional continuous rate model [KCF17, MM07] with lateral inhibition.

It is easy to see the following claims hold. For brevity, their proofs are omitted.

Claim 6. For $t \ge 1$ and for $i = 1, \dots, n$, $V_{t-1}(v_i) > 0$ if and only if $S_{t-1}(u_i) = 1$ and $\sum_{j:1 \le j \le n, \& j \ne i} S_{t-1}(v_j) \le k-1$, i.e., at time t-1, input neuron u_i spikes, and fewer than k-1 other output neurons spike.

Claim 7. For $t \ge 1$ and for $i = 1, \dots, n$, $V_{t-1}(v_i) \le -1$ only if $\sum_{j:1 \le j \le n, \& j \ne i} S_{t-1}(v_j) \ge k$, i.e., at time t-1, more than k other output neurons spike.

Note that $\sum_{j:1 \leq j \leq n, \& j \neq i} S_{t-1}(v_j) \geq k$ is not a sufficient condition to have $V_{t-1}(v_i) \leq -1$. To see this, suppose $\sum_{j:1 \leq j \leq n, \& j \neq i} S_{t-1}(v_j) = k$ and $S_{t-1}(u_i) = 1$. In this case it holds that $V_{t-1}(v_i) = 0$.

Claim 8. For $t \ge 1$ and for $i = 1, \dots, n$, if $V_{t-1}(v_i) = 0$, one of the following holds: (1) $S_{t-1}(u_i) = 1$ and $\sum_{j:1 \le j \le n, \& j \ne i} S_{t-1}(v_j) = k$, i.e., at time t-1, input neuron u_i spikes, and exactly k other output neurons spike;

(2) $S_{t-1}(u_i) = 0$ and $\sum_{j:1 \le j \le n, \& j \ne i} S_{t-1}(v_j) = 0$, i.e., at time t-1, input neuron u_i does not spike, and no other output neurons spike.

5.1.3 Activation functions:

There are many different choices of activation functions; see [wik] for a detailed list. In our construction, we use a simple threshold activation function, i.e.,

$$S_t(v_i) = \begin{cases} 1, & \text{if } (b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^m \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m\sum_{r=1}^m \mathbf{1}_{\{V_{t-r}(v_i)\leq-1\}}\right]_+ \ge b; \\ 0, & \text{otherwise,} \end{cases}$$

 $[\cdot]_{+} = \max[\cdot, 0]$, and b > 0 is the bias at neuron v_i for $i = 1, \dots, n$. It is easy to see that this activation function falls under the general form given by (2).

Remark 9. If the output neuron v_i does not spike at time t - 1, i.e., $S_{t-1}(v_i) = 0$, then in order for v_i to spike at time t, the following needs to hold:

$$\left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq -1\}}\right]_{+} \ge b.$$

In contrast, if the output neuron v_i does spike at time t-1, i.e., $S_{t-1}(v_i) = 1$, then

$$\left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq -1\}}\right]_{+} \ge 1$$

is enough for v_i to spike at time t. That is, under our activation rule, $S_{t-1}(v_i) = 1$ makes the activation of v_i much easier in the next round. However, if there exists $r \in \{1, 2, \dots, m\}$ such that

$$\mathbf{1}_{\{V_{t-r}(v_i) \le -1\}} = 1$$

then

$$\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq -1\}} \leq \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m$$
$$< m - m = 0.$$

Thus,

$$(b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq-1\}}\right]_{+}$$

= $(b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + 0$
< $b-1 < b$,

i.e., the output neuron v_i does not spike at time t. In other words, as long as there exists $r \in \{1, 2, \dots, m\}$ such that $\mathbf{1}_{\{V_{t-r}(v_i) \leq -1\}} = 1$, the activation of v_i is inhibited at time t.

5.1.4 Local memorization capability:

In our proposed circuit, we require that m satisfies the following:

$$m \ge \frac{8C^2(1-c)}{c^2(1-C)} \left(\log\left(\frac{3}{\delta}\right) + \log k(n-k) \right) T_{\mathcal{R}} := m^*$$
(9)

for target decision accuracy $1 - \delta \in (0, 1)$. In addition, we set $b = cm^*$. Recall that $c, C \in (0, 1)$ are two absolute constants that are lower bound and upper bound of any \mathcal{R} , respectively.

Intuitively, when other parameters are fixed, the higher the desired accuracy (i.e., the smaller δ), the larger m^* , i.e., the more memory is needed for selecting the winners in our WTA circuit. Similarly, the easier to distinguish two spike trains with different statistics (i.e., the lower $T_{\mathcal{R}}$), the smaller m^* . Interesting, with other parameters fixed, m^* depends on k as follows: m^* is increasing in k when $k \in \{1, \dots, \lfloor \frac{n}{2} \rfloor\}$, and m^* is decreasing in k when $k \in \{\lceil \frac{n}{2} \rceil, \dots, n-1\}$. In many practical settings we care about the region where $k \ll n$. Besides, with the choice of bias $b = cm^*$, the larger m^* also implies longer time is needed for our WTA circuit to declare k winners; details can be found (1) in Theorem 10.

On the other hand, in most neurons the synaptic plasticity time window is about 80-120 ms, and it is unclear whether (9) can be immediately satisfied or not. Fortunately, even if (9) is not immediately satisfied by a neuron due to its local bio-plausibility, it is possible that its local memory might be realized using a chain of hidden neurons.

5.1.5 Algorithm 1

The dynamics of our WTA circut is summarized in Algorithm 1, which is fully determined by what has been described in Sections 5.1.1, 5.1.2, 5.1.3, and 5.1.4. For Algorithm 1, we declare the first k output neurons that spike simultaneously to be winners.

Algorithm 1: *k*–WTA

1 Input: $\mathcal{R}, m, b, \text{ and } \delta$. 2 for $t \ge 1$ do At output neuron v_i for $i = 1, \dots, n$: $V_{t-1}(v_i) \leftarrow S_{t-1}(u_i) - \frac{1}{k} \sum_{i:1 \le j \le n} \sum_{k: j \ne i} S_{t-1}(v_j);$ 3 $V_{t-1}(v_i) \leftarrow [V_{t-1}(v_i), V_{t-2}(v_i), \cdots, V_{t-m}(v_i)];$ 4 $S_{t-1}(v_i) \leftarrow [S_{t-1}(v_i), S_{t-2}(v_i), \cdots, S_{t-m}(v_i)];$ 5 $M_t(v_i) \leftarrow (V_{t-1}(v_i), S_{t-1}(v_i));$ 6 $\mathbf{if} \ (b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^m \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m\sum_{r=1}^m \mathbf{1}_{\{V_{t-r}(v_i)\leq -1\}}\right]_+ \ge b \ \mathbf{then}$ 7 $S_t(v_i) \leftarrow 1.$ 8 else 9 $S_t(v_i) \leftarrow 0.$ 10

5.2 Circuit Performance

Recall that $\mathcal{W}(p)$ and m^* are defined in (3) and (9), respectively.

Theorem 10. Fix $\delta \in (0,1]$, and $1 \le k \le n-1$. Choose $m \ge m^*$ and $b = \max\{cm^*, 2\}$. Then for any admissible rate assignment p, with probability at least $1 - \delta$, the following hold:

- (1) There exist k output neurons that spike simultaneously by time m^* .
- (2) The first set of such k output neurons are the true winners $\mathcal{W}(\mathbf{p})$.
- (3) From the first time in which these k output neurons spike simultaneously, these k output neurons spike consecutively for at least b times, and no other output neurons can spike within b times.

The proof of Theorem 10 can be found in Appendix F. The first bullet in Theorem 10 implies that our WTA circuit can provide an output (a selection of k output neurons) by time m^* ; the second bullet in Theorem 10 says that the circuit's output indeed corresponds to the k true winners; and the third bullet says that the k simultaneous spikes of the selected winners are stable – the k selected winners continue to spike consecutively for at least b times. The proof of Theorem 10 essentially says that with high probability, under Algorithm 1, the number of output neurons that spike simultaneously is monotonically increasing until it reaches k. Upon the simultaneous spike of k output neurons, by our threshold activation rule, we know that the other output neurons are likely to be inhibited. In particular, if these k output neurons are likely to be inhibited for at least b times.

Remark 11 (Controlling stability). As can be seen from the proof of Theorem 10, in the activation function of Algorithm 1

$$(b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq-1\}}\right]_{+} \ge b$$

the first term $(b-1)\mathbf{1}_{\{S_{t-r}(v_i)=1\}}$ is crucial in achieving (3) in Theorem 10. In fact, we can increase the stability period by introducing a stability parameter s such that $1 < s \leq m$ and modifying the activation rule. Details can be found in Algorithm 2. It is easy to see that the activation function falls under the general form in (2). In the new activation function in Algorithm 2, for output neuron v_i , once it spikes, it continues to spike for at least s times. Following our line of analysis in the proof of Theorem 10, it can be seen that the declared k winners, from the first time they spike simultaneously, continue to spike consecutively for at least s times.

Remark 12 (Order-optimality). The decision time performance stated in (1) of Theorem 10 matches the information-theoretical lower bound in Theorem 4 up to a multiplicative constant factor both (a) when δ is sufficiently small and does not depend on $n, k, T_{\mathcal{R}}, c$, and C, and (b) when δ decays to zero at a speed at most $\frac{1}{(k(n-k))^{c_0}}$ where $c_0 > 0$ is some fixed constant. The detailed order-optimality argument is given next.

Algorithm 2: *k*–WTA

1 Input: \mathcal{R} , m, b, δ , and s where $1 < s \le m$. 2 for $t \ge 1$ do At output neuron v_i for $i = 1, \dots, n$: 3 $V_{t-1}(v_i) \leftarrow S_{t-1}(u_i) - \frac{1}{k} \sum_{j:1 \le j \le n, \& j \ne i} S_{t-1}(v_j);$ 4 5 6 $M_t(v_i) \leftarrow (V_{t-1}(v_i), S_{t-1}(v_i)).$ 7 if $\left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq -1\}}\right]_{+} \ge b$ then 8 $S_t(v_i) \leftarrow 1.$ 9 else 10 if $S_{t-1}(v_i) = 1$ and $\exists r \in \{2, \dots, s\}$ such that $S_{t-r}(v_i) = 0$ then 11 $S_t(v_i) \leftarrow 1.$ 1213 else $S_t(v_i) \leftarrow 0.$ 14

Suppose that δ is sufficiently small and does not depend on n, k, $T_{\mathcal{R}}$, c, and C Here, for ease of exposition, we illustrate the order-optimality with a specific choice of δ . In fact, the order-optimality holds generally for constant $\delta \in (0, 1)$ as long as it does not depend on n, k, $T_{\mathcal{R}}$, c, and C.

Suppose the target decision accuracy is $1 - \delta = 0.9$, i.e., $\delta = 0.1$. Then as long as $n \ge 31$, for any $1 \le k \le n - 1$,

$$m^* = \frac{8C^2(1-c)}{c^2(1-C)} \left(\log \frac{3}{0.1} + \log k(n-k) \right) T_{\mathcal{R}} \le \frac{16C^2(1-c)}{c^2(1-C)} \log k(n-k) T_{\mathcal{R}}.$$

On the other hand, recall from Theorem 4 that to have $\delta = 0.1$, the decision time is no less than

$$((1-\delta)\log(k(n-k)+1)-1)T_{\mathcal{R}} \ge \frac{1}{2}\log(k(n-k)+1)T_{\mathcal{R}} \ge \frac{1}{2}\log k(n-k)T_{\mathcal{R}}$$

where the first inequality holds as long as $n \ge 8$. Thus, when $n \ge 31$, in order to achieve the decision accuracy $1 - \delta = 0.9$, the decision time of our WTA circuit is on the same order of the information-theoretic lower bound in Theorem 4.

Suppose δ decays to zero at a moderate speed The decision time of our WTA circuit is order-optimal even for diminishing decision error δ as long as $\delta = \Omega(\frac{3}{(k(n-k))^{c_0}})$ where $c_0 > 0$ – it does not decay to zero "too fast" in k(n-k). To see this, let $\delta = \frac{3}{(k(n-k))^{c_0}}$ for some constant $c_0 > 0$. We have

$$\frac{8C^2(1-c)}{c^2(1-C)} \left(\log\left(\frac{3}{\frac{3}{(k(n-k))^{c_0}}}\right) + \log k(n-k) \right) T_{\mathcal{R}} = \frac{8C^2(1-c)(c_0+1)}{c^2(1-C)} \log k(n-k) T_{\mathcal{R}}.$$
 (10)

Resetting circuit when the input spike trains become quiescent In Algorithm 1, if the input spike trains become quiescent, then the corresponding circuits also become quiescent despite some delay in this response.

Lemma 13. If all input neurons are quiescent at time t_0 , and remain to be quiescent for all $t \ge t_0$, then $V_t(v_i) = 0$ and $S_t(v_i) = 0$ for any $t > t_0 + m$.

Lemma 13 is proved in Appendix D.

6 Discussion

In this paper, we investigated how k-WTA computation is robustly achieved in the presence of inherent noise in the input spike trains. In a spike-based k-WTA model, n randomly generated input spike trains are competing with each other, and the top k neurons with highest underlying statistics are the true winners. Given the stochastic nature of the spike trains, it is not trivial to properly select winners among a group of neurons. We derived an information-theoretic lower bound on the decision time for a given decision accuracy. Notably, this lower bound holds universally for any WTA circuit that falls within our model framework, regardless of their circuit architectures or their adopted activation functions. Furthermore, we constructed a circuit whose decision time matches this lower bound up to a constant multiplicative factor, suggesting that our derived lower bound is order-optimal. Here the order-optimality is stated in terms of its scaling in n, k, and T_R .

6.1 Comparison to previous WTA models

Randomness is introduced at different stages of brain computation and the stochastic nature of the spike trains are well observed [BAB+97, KRR00, MA09, FDMM18]. In our work, we focused on how to robustly achieve k-WTA computation in face of the intrinsic randomness in the spike trains. A common WTA model assumes that neurons transmit information by a continuous variable such as firing rate [DA01], which ignores the intrinsic randomness in spiking trains. Although some studies used additive Gaussian noise [KCF17, LLW13, LIKB99, RV06] in their rate-based WTA circuits to account for input randomness, these circuits are usually very sensitive to noise and could not successfully select even a single winner unless additional non-linearity is added [KCF17]. In fact, a neuron with a second non-linearity is similar to an output neuron in our constructed WTA circuit in that they both integrate their local inputs. Unfortunately, only simulation results were provided in [KCF17]; a theoretical justification of why such second non-linearity makes their WTA circuit robust to input noise is lacking. Though we focused on spike-based model, we hope our results can provide some insights for the rate-based model as well. On top of that, a rate-based model would require a high communication bandwidth, yet communication bandwidth is limited in the brain. Our spiking neural network model captures this feature by having a low communication cost, since it broadcasts 1 bit only.

However, we did not try to model every biologically relevant feature. In several studies using spiking network models, individual units are often modeled with details like ion channels and specific synaptic connectivity. Though more biologically relevant than our spiking neuron network model, those details significantly complicate the analysis. In fact, it could be challenging and intricate to move beyond computer simulation to characterize the model dynamics (such as the spiking nature of each unit, the time it takes to stabilize, etc.) analytically.

6.2 Potential applications for physiological experiments

Our work further provided testable hypotheses on how network size, similarities between input spike trains, and synaptic memory capacity would affect this lower bound. For example, in behavioral experiments using electrolytic lesions or pharmacological inhibition [CMA⁺03, HDS06, YLS13, KYPH16], the changes in performance are often highly variable and nonlinear. One possible reason comes from the difficulty of precisely manipulating network size as well as a lack of theoretical description of the relationship between network size and performance. With our analytical characterization, one might be able to estimate changes in the effective network size given performance in a decision-making task.

Besides the effect of network size, the distribution of feature representations (i.e., different set $\mathcal{R}s$ of different individual animals) could be used to account for between-subject variability in decision making. Consider a random-dot coherent motion task where animals need to decide which of two directions the majority of dots are moving [SN01]. In this task, performance accuracy and reaction time vary across animals. If we perform neural recordings in their visual cortex (i.e., to record their $\mathcal{R}s$), we might be able to decide their reaction time or accuracy, given population representations of dot motion in these cortical neurons [SN96, JM06]. For example, an animal whose stimulus-evoked responses are more heterogeneous in the visual cortex might be able to react faster given the same accuracy, governed by our derived lower-bound.

Last but not least, our work also offered predictions on how local memory capacity could affect performance in decision-making. For example, when there is more ambiguity in input representations, to obtain the same performance (both accuracy and decision time), a larger time window for memory storage in synapses [KPS10] is required. From previous experimental work [BMG⁺17], we know that synaptic plasticity has time scale ranging from milliseconds to seconds across different brain regions, and such plasticity could efficiently store entire behavioral sequences within synaptic weights. Combining with our analytical characterization, when performance accuracy changes over time, assuming other parameters such as input statistics, decision time and network size are fixed, one might be able to predict how synaptic plasticity changes. Overall, our work not only provided a theoretical framework, but also provided a set of testable hypotheses on neural recordings and behaviors in decision-making under ambiguity.

6.3 Limitations and extensions

When δ is a constant, our lower bound is order-optimal in terms of its scaling in n, k, and $T_{\mathcal{R}}$. Nevertheless, the scaling of the derived lower bound in terms of δ is not tight. It would be interesting to know the optimal scaling in δ when other parameters $(n, k, \text{ and } T_{\mathcal{R}})$ are fixed. We leave it as one future direction.

To simplify complexity, our model posed a few assumptions that ignored some features in the brain. One of these assumptions is that each input neuron is independent. However, various degrees of average noise correlations between cortical neurons have been reported. For example, average noise correlations in primary visual cortex could be close to 0.1 [SSB+15], 0.18 [SK08], or even much larger as 0.35 [GD08]. Similarly, noise correlations have been observed in other sensory brain regions [CK11]. In our work, we ignored correlations between these neurons, but it would be interesting as a future direction to extend in our spiking network model.

Second, our model used a threshold activation function by assuming the synaptic transmission is basically noise-free and that the only noise source comes from the input in this paper. However, synaptic transmission is highly unreliable in biological networks [AS94, FSW08, Bor10], and a deterministic activation function would fail to capture this feature compared to a stochastic activation function. Moreover, failure in synaptic transmission could serve a computational role [BS09, Maa97].

Another assumption in our circuit is that the output neurons can inhibit each other. In common scenarios, an output neuron is usually excitatory, and does not inhibit other neurons directly without recruiting inhibitory cells. We incorporate stability in these output neurons by assuming they can inhibit each other in our circuit implementation. For a model where an output neuron is limited to be excitatory only, we can add a chain of inhibitory neurons to achieve stability WTA computation.

Last but not least, in our k-WTA circuit, the number of output neurons that spike simultaneously increases monotonically until there are exactly k output neurons that spike simultaneously. We acknowledge that this might not be biologically plausible in most cases in the brain. From large-scale neural recordings, we know that the number of neurons that spike simultaneously is usually variable, so this could be a future direction to construct a circuit that better matches experimental observations.

Acknowledgement

We would like to thank Christopher Quinn at Purdue University and Zhi-Hong Mao at University of Pittsburgh for the helpful discussions and references.

References

- [AS94] Christina Allen and Charles F Stevens. An evaluation of causes for unreliability of synaptic transmission. *Proceedings of the National Academy of Sciences*, 91(22):10380–10383, 1994. 14
- [BAB⁺97] Roland Baddeley, Larry F Abbott, Michael CA Booth, Frank Sengpiel, Tobe Freeman, Edward A Wakeman, and Edmund T Rolls. Responses of neurons in primary and inferior temporal visual cortices to natural scenes. Proceedings of the Royal Society of London B: Biological Sciences, 264(1389):1775–1783, 1997. 2, 13

- [BC95] György Buzsáki and James J Chrobak. Temporal structure in spatially organized neuronal ensembles: a role for interneuronal networks. *Current opinion in neurobiology*, 5(4):504–510, 1995. 2
- [BIM98] Michael J Berry II and Markus Meister. Refractoriness and neural precision. In Advances in Neural Information Processing Systems, pages 110–116, 1998. 2
- [BMG⁺17] 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. 14
- [Bor10] J Gerard G Borst. The low synaptic release probability in vivo. Trends in neurosciences, 33(6):259–266, 2010. 14
- [BS09] Tiago Branco and Kevin Staras. The probability of neurotransmitter release: variability and feedback control at single synapses. *Nature Reviews Neuroscience*, 10(5):373, 2009. 14
- [CK11] Marlene R Cohen and Adam Kohn. Measuring and interpreting neuronal correlations. Nature neuroscience, 14(7):811, 2011. 14
- [CMA⁺03] Luke Clark, Facundo Manes, Nagui Antoun, Barbara J Sahakian, and Trevor W Robbins. The contributions of lesion laterality and lesion volume to decision-making impairment following frontal lobe damage. *Neuropsychologia*, 41(11):1474–1483, 2003. 13
- [DA01] Peter Dayan and Laurence F Abbott. Theoretical neuroscience: computational and mathematical modeling of neural systems. 2001. 13
- [EB02] Marc O Ernst and Martin S Banks. Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870):429, 2002. 2
- [FDMM18] Ulisse Ferrari, Stephane Deny, Olivier Marre, and Thierry Mora. A simple model for low variability in neural spike trains. arXiv preprint arXiv:1801.01362, 2018. 2, 13
- [FSW08] A Aldo Faisal, Luc PJ Selen, and Daniel M Wolpert. Noise in the nervous system. *Nature reviews neuroscience*, 9(4):292, 2008. 2, 14
- [GD08] Diego A Gutnisky and Valentin Dragoi. Adaptive coding of visual information in neural populations. *Nature*, 452(7184):220, 2008. 14
- [GKvHW96] Wulfram Gerstner, Richard Kempter, J Leo van Hemmen, and Hermann Wagner. A neuronal learning rule for sub-millisecond temporal coding. *Nature*, 383(6595):76, 1996.
- [HDS06] Timothy D Hanks, Jochen Ditterich, and Michael N Shadlen. Microstimulation of macaque area lip affects decision-making in a motion discrimination task. *Nature neuroscience*, 9(5):682, 2006. 13
- [HDZ08] Tomáš Hromádka, Michael R DeWeese, and Anthony M Zador. Sparse representation of sounds in the unanesthetized auditory cortex. *PLoS biology*, 6(1):e16, 2008. 2
- [HSM06] Thomas TG Hahn, Bert Sakmann, and Mayank R Mehta. Phase-locking of hippocampal interneurons' membrane potential to neocortical up-down states. *Nature neuroscience*, 9(11):1359, 2006. 2
- [HW59] David H Hubel and Torsten N Wiesel. Receptive fields of single neurones in the cat's striate cortex. *The Journal of physiology*, 148(3):574–591, 1959. 2
- [HW98] Christopher M Harris and Daniel M Wolpert. Signal-dependent noise determines motor planning. Nature, 394(6695):780, 1998. 2

- [IKN98] Laurent Itti, Christof Koch, and Ernst Niebur. A model of saliency-based visual attention for rapid scene analysis. *IEEE Transactions on pattern analysis and machine intelligence*, 20(11):1254–1259, 1998. 2
- [JB67] I Jacobs and E Berlekamp. A lower bound to the distribution of computation for sequential decoding. *IEEE Transactions on Information Theory*, 13(2):167–174, 1967. 21
- [JM06] Mehrdad Jazayeri and J Anthony Movshon. Optimal representation of sensory information by neural populations. *Nature neuroscience*, 9(5):690, 2006. 13
- [KCF17] Birgit Kriener, Rishidev Chaudhuri, and Ila Fiete. How fast is neural winner-take-all when deciding between many options? *bioRxiv*, page 231753, 2017. 2, 9, 13
- [KF08] Mattias P Karlsson and Loren M Frank. Network dynamics underlying the formation of sparse, informative representations in the hippocampus. *Journal of Neuroscience*, 28(52):14271–14281, 2008. 2
- [KK01] Masaharu Kinoshita and Hidehiko Komatsu. Neural representation of the luminance and brightness of a uniform surface in the macaque primary visual cortex. *Journal of neurophysiology*, 86(5):2559–2570, 2001. 2
- [KP04] David C Knill and Alexandre Pouget. The bayesian brain: the role of uncertainty in neural coding and computation. *TRENDS in Neurosciences*, 27(12):712–719, 2004. 2
- [KPS10] Andreas Knoblauch, Günther Palm, and Friedrich T Sommer. Memory capacities for synaptic and structural plasticity. *Neural Computation*, 22(2):289–341, 2010. 14
- [KRR00] Prakash Kara, Pamela Reinagel, and R Clay Reid. Low response variability in simultaneously recorded retinal, thalamic, and cortical neurons. *Neuron*, 27(3):635–646, 2000. 2, 13
- [KTA⁺03] Zoe Kourtzi, Andreas S Tolias, Christian F Altmann, Mark Augath, and Nikos K Logothetis. Integration of local features into global shapes: monkey and human fmri studies. Neuron, 37(2):333–346, 2003. 2
- [KW04] Konrad P Körding and Daniel M Wolpert. Bayesian integration in sensorimotor learning. Nature, 427(6971):244, 2004. 2
- [KYPH16] Leor N Katz, Jacob L Yates, Jonathan W Pillow, and Alexander C Huk. Dissociated functional significance of decision-related activity in the primate dorsal stream. Nature, 535(7611):285, 2016. 13
- [LCG⁺15] Nuo Li, Tsai-Wen Chen, Zengcai V Guo, Charles R Gerfen, and Karel Svoboda. A motor cortex circuit for motor planning and movement. *Nature*, 519(7541):51, 2015. 2
- [LIKB99] Dale K Lee, Laurent Itti, Christof Koch, and Jochen Braun. Attention activates winner-take-all competition among visual filters. *Nature neuroscience*, 2(4):375, 1999. 2, 13
- [LLW13] Shuai Li, Yangming Li, and Zheng Wang. A class of finite-time dual neural networks for solving quadratic programming problems and its k-winners-take-all application. *Neural Networks*, 39:27–39, 2013. 2, 13
- [LMP16] Nancy Lynch, Cameron Musco, and Merav Parter. Computational tradeoffs in biological neural networks: Self-stabilizing winner-take-all networks. arXiv preprint arXiv:1610.02084, 2016. 5
- [MA09] Gaby Maimon and John A Assad. Beyond poisson: increased spike-time regularity across primate parietal cortex. *Neuron*, 62(3):426–440, 2009. 2, 13
- [Maa97] Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. Neural networks, 10(9):1659–1671, 1997. 5, 14

- [Maa00] Wolfgang Maass. On the computational power of winner-take-all. Neural computation, 12(11):2519–2535, 2000. 2
- [MCM07] Najib J Majaj, Matteo Carandini, and J Anthony Movshon. Motion integration by neurons in macaque mt is local, not global. *Journal of Neuroscience*, 27(2):366–370, 2007. 2
- [MM07] Zhi-Hong Mao and Steve G Massaquoi. Dynamics of winner-take-all competition in recurrent neural networks with lateral inhibition. *IEEE transactions on neural networks*, 18(1):55–69, 2007. 9
- [Nel04] Israel Nelken. Processing of complex stimuli and natural scenes in the auditory cortex. *Current opinion in neurobiology*, 14(4):474–480, 2004. 2
- [OF04] Bruno A Olshausen and David J Field. Sparse coding of sensory inputs. *Current opinion in neurobiology*, 14(4):481–487, 2004. 2
- [PG99] Michael L Platt and Paul W Glimcher. Neural correlates of decision variables in parietal cortex. Nature, 400(6741):233, 1999. 2
- [POMT⁺02] Javier Perez-Orive, Ofer Mazor, Glenn C Turner, Stijn Cassenaer, Rachel I Wilson, and Gilles Laurent. Oscillations and sparsening of odor representations in the mushroom body. *Science*, 297(5580):359–365, 2002. 2
- [PW14] Yury Polyanskiy and Yihong Wu. Lecture notes on information theory. Lecture Notes for ECE563 (UIUC) and, 6:2012–2016, 2014. 7, 18, 20
- [QKKF08] R Quian Quiroga, Gabriel Kreiman, Christof Koch, and Itzhak Fried. Sparse but not ?grandmother-cell?coding in the medial temporal lobe. Trends in cognitive sciences, 12(3):87– 91, 2008. 2
- [RP99] Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex. Nature neuroscience, 2(11):1019, 1999. 2
- [RPG99] Peter Redgrave, Tony J Prescott, and Kevin Gurney. The basal ganglia: a vertebrate solution to the selection problem? *Neuroscience*, 89(4):1009–1023, 1999. 2
- [RV06] Nicolas P Rougier and Julien Vitay. Emergence of attention within a neural population. *Neural Networks*, 19(5):573–581, 2006. 2, 13
- [SK08] Matthew A Smith and Adam Kohn. Spatial and temporal scales of neuronal correlation in primary visual cortex. *Journal of Neuroscience*, 28(48):12591–12603, 2008. 14
- [SLW05] Athanassios G Siapas, Evgueniy V Lubenov, and Matthew A Wilson. Prefrontal phase locking to hippocampal theta oscillations. *Neuron*, 46(1):141–151, 2005. 2
- [SN96] Michael N Shadlen and William T Newsome. Motion perception: seeing and deciding. Proceedings of the national academy of sciences, 93(2):628–633, 1996. 13
- [SN01] Michael N Shadlen and William T Newsome. Neural basis of a perceptual decision in the parietal cortex (area lip) of the rhesus monkey. *Journal of neurophysiology*, 86(4):1916–1936, 2001. 2, 13
- [SS06] Alan A Stocker and Eero P Simoncelli. Noise characteristics and prior expectations in human visual speed perception. *Nature neuroscience*, 9(4):578, 2006. 2
- [SSB+15] Marieke L Schölvinck, Aman B Saleem, Andrea Benucci, Kenneth D Harris, and Matteo Carandini. Cortical state determines global variability and correlations in visual cortex. Journal of Neuroscience, 35(1):170–178, 2015. 14

- [VRP⁺18] Stephen J Verzi, Fredrick Rothganger, Ojas D Parekh, Tu-Thach Quach, Nadine E Miner, Craig M Vineyard, Conrad D James, and James B Aimone. Computing with spikes: The advantage of fine-grained timing. *Neural computation*, 30(10):2660–2690, 2018. 5
- [wik] Activation function. https://en.wikipedia.org/wiki/Activation_function. Accessed: 2018-08-08. 9
- [Wu17] Yihong Wu. Lecture notes on information-theoretic methods for high-dimensional statistics. Lecture Notes for ECE598YW (UIUC), 2017. 5
- [YG98] AL Yuille and D Geiger. The handbook of brain theory and neural networks, 1998. 2
- [YLS13] Eric A Yttri, Yuqing Liu, and Lawrence H Snyder. Lesions of cortical area lip affect reach onset only when the reach is accompanied by a saccade, revealing an active eye-hand coordination circuit. Proceedings of the National Academy of Sciences, 110(6):2371–2376, 2013. 13

Appendices

A Preliminaries

In this section, we present some preliminaries on information measures and Fano's inequality. Interested readers are referred to [PW14] for comprehensive background.

A.1 Information Measures

Let X and Y be two random variables. The mutual information between X and Y, denoted by I(X;Y), measures the dependence between X and Y, or, the information about X (resp. T) provided by Y (resp. X).

Definition 14 (Mutual information). Let X and Y be two random variables.

$$I(X;Y) := D(P_{XY} || P_X P_Y), D(P || Q) := \sum_{a \in \mathcal{A}} P(a) \log \frac{P(a)}{Q(a)},$$

where P_{XY} denotes the joint distribution of X and Y, and $P_X P_Y$ denotes the product of the marginal distributions of X and Y.

In the following, we use the notation $X \to Y$ to denote that Y is a (possibly random) function of X. Thus, $W \to X \to Y \to \widehat{W}$ means that X is a (possibly random) function of W; Y is a (possibly random) function of X; and \widehat{W} is a (possibly random) function of Y. Fano's inequality:

Theorem 15. [*PW14*, Corollary 5.1] Let $T: \Theta \to [M]$, and let $\theta \to X \to Y \to \widehat{T}(\theta)$ be an arbitrary Markov chain. Suppose both θ and $T(\theta)$ are uniformly distributed over a set of size M. Then

$$P_e := \mathbb{P}\left\{T(\theta) \neq \widehat{T}(\theta)\right\} \ge 1 - \frac{I(X;Y) + 1}{\log M}.$$

Theorem 16 (Chernoff Bound). Let X_1, \dots, X_n be *i.i.d.* with $X_i \in \{0,1\}$ and $\mathbb{P}\{X_1=1\} = p$. Set $X = \sum_{i=1}^n X_i$. Then

- for any $t \in [0, 1-p]$, we have $\mathbb{P}\{X \ge (p+t)n\} \le \exp(-nd(p+t || p))$.
- for any $t \in [0, p]$, we have $\mathbb{P} \{ X \le (p t) n \} \le \exp(-nd(p t \| p)).$

B Proof of Lemma 2

Proof of Lemma 2. Lemma 2 follows easily from the independence between input spike trains and the assumption that the spikes in each input spike train are i.i.d.. For completeness, we present the proof as follows.

Recall that

$$\boldsymbol{S} := \left[\{ S_t(u_1) \}_{t=1}^T, \cdots, \{ S_t(u_n) \}_{t=1}^T \right]$$

Let $\mathbf{s} = [s_1, \dots, s_n]$ such that each component s_i is a binary sequence of length T, i.e.,

$$s_i = [b_1^i, \cdots, b_T^i] \in \{0, 1\}^T$$

For each $i = 1, \dots, n$, let $P_{\mathbf{S}}(\{S_t(u_i)\}_{t=1}^T)$ and $Q_{\mathbf{S}}(\{S_t(u_i)\}_{t=1}^T)$ be the marginal distributions of $\{S_t(u_i)\}_{t=1}^T$ under joint distributions $P_{\mathbf{S}}$ and $Q_{\mathbf{S}}$ respectively. Similarly, $P_{\mathbf{S}}(S_t(u_i))$ and $Q_{\mathbf{S}}(S_t(u_i))$ are the corresponding two marginal distributions of $S_t(u_i)$. Thus, we have

$$\begin{split} & D\left(P_{S}(\{S_{t}(u_{i})\}_{t=1}^{T}) \parallel Q_{S}(\{S_{t}(u_{i})\}_{t=1}^{T})\right) \\ & \stackrel{(a)}{=} \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} P_{S}(\{S_{t}(u_{i})\}_{t=1}^{T} = [b_{1}^{i}, \cdots, b_{T}^{i}]) \log \frac{P_{S}(\{S_{t}(u_{i})\}_{t=1}^{T} = [b_{1}^{i}, \cdots, b_{T}^{i}])}{Q_{S}(\{S_{t}(u_{i})\}_{t=1}^{T} = [b_{1}^{i}, \cdots, b_{T}^{i}])} \\ & \stackrel{(b)}{=} \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} \left(\prod_{t'=0}^{T-1} P_{S}(S_{t'}(u_{i}) = b_{t'}^{i})\right) \log \frac{\prod_{t=1}^{T} P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{\prod_{t=1}^{T} Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} \left(\prod_{t'=0}^{T-1} P_{S}(S_{t'}(u_{i}) = b_{t'}^{i})\right) \sum_{t=1}^{T} \log \frac{P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{t=1}^{T} \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} \left(\prod_{t'=0}^{T-1} P_{S}(S_{t'}(u_{i}) = b_{t'}^{i})\right) \log \frac{P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{t=1}^{T} \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} \left(\prod_{t'=0}^{T-1} P_{S}(S_{t'}(u_{i}) = b_{t'}^{i})\right) \log \frac{P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{t=1}^{T} \sum_{[b_{1}^{i}, \cdots, b_{T}^{i}]} \left(\prod_{t'=0}^{T-1} P_{S}(S_{t'}(u_{i}) = b_{t'}^{i})\right) P_{S}(S_{t}(u_{i}) = b_{t}) \log \frac{P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{t=1}^{T} \sum_{b_{1}^{i}} P_{S}(S_{t}(u_{i}) = b_{t}^{i}) \log \frac{P_{S}(S_{t}(u_{i}) = b_{t}^{i})}{Q_{S}(S_{t}(u_{i}) = b_{t}^{i})} \\ & = \sum_{t=1}^{T} \left(p_{i} \log \frac{p_{i}}{q_{i}} + (1 - p_{i}) \log \frac{1 - p_{i}}{1 - q_{i}}\right) \\ & = \sum_{t=1}^{T} d(p_{i} \parallel q_{i}) = T \cdot d(p_{i} \parallel q_{i}). \end{split}$$

where $\sum_{[b_1^i, \dots, b_T^i]}$ is the summation over all binary sequences of length T. In the last displayed equation, equality (a) follows from the definition of KL divergence; equality (b) is true because of independence of spikes; equality (c) follows from the fact that for any fixed b_t^i ,

$$\sum_{\begin{bmatrix} b_1^i, \cdots, b_t^i \end{bmatrix} \setminus \{t\}} \left(\prod_{t'=1 \& t' \neq t}^T P_{\boldsymbol{S}}(S_{t'}(u_i) = b_{t'}^i) \right) = 1,$$

where we use $\sum_{[b_1^i, \dots, b_T^i] \setminus \{t\}}$ to denote the summation over all binary sequences of length T with the t-th entry fixed.

Similarly, we get

$$D(P_{\mathbf{S}} \parallel Q_{\mathbf{S}}) = \sum_{\mathbf{s} = [s_1, \cdots, s_n]} P_{\mathbf{S}}(\mathbf{S} = \mathbf{s}) \log \frac{P_{\mathbf{S}}(\mathbf{S} = \mathbf{s})}{Q_{\mathbf{S}}(\mathbf{S} = \mathbf{s})}$$
$$= \sum_{i=1}^n D\left(P_{\mathbf{S}}(\{S_t(u_i)\}_{t=1}^T) \parallel Q_{\mathbf{S}}(\{S_t(u_i)\}_{t=1}^T)\right)$$
$$= \sum_{i=1}^n Td(p_i \parallel q_i) = T\sum_{i=1}^n d(p_i \parallel q_i),$$

proving the lemma.

C Proof of Lemma 3

Proof of Lemma 3. Since mutual information can be viewed as distance to product distributions, by [PW14, Theorem 3.4], we have

$$I(X_{\boldsymbol{p}}; \boldsymbol{S}) = \min_{Q_{X_{\boldsymbol{p}}}Q_{\boldsymbol{S}}} D\left(P_{X_{\boldsymbol{p}}, \boldsymbol{S}} \parallel Q_{X_{\boldsymbol{p}}}Q_{\boldsymbol{S}}\right)$$

where $P_{X_p,S}$ is the joint distribution of X_p and S, and Q_{X_p} and Q_S are any distributions of X_p and S, respectively.

For any fixed $Q_{\mathbf{S}}$, it holds that

$$\min_{Q_{X_{p}}} D\left(P_{X_{p},S} \parallel Q_{X_{p}}Q_{S}\right) = \min_{Q_{X_{p}}} D\left(P_{S|X_{p}}P_{X_{p}} \parallel Q_{X_{p}}Q_{S}\right)$$
$$\leq D\left(P_{S|X_{p}}P_{X_{p}} \parallel P_{X_{p}}Q_{S}\right),$$

where the equality follows from conditioning, and the inequality is true because the best choice over all Q_{X_p} cannot be worse than any specific choice of Q_{X_p} . Here $S \mid X_p$ denotes the *n* input spike trains conditioning on the choice of rate assignment.

For any fixed $Q_{\mathbf{S}}$, we have

$$\begin{split} D\left(P_{S|X_{p}}P_{X_{p}} \parallel P_{X_{p}}Q_{S}\right) &= P_{X_{p}}(X_{p} = p^{0})\sum_{s} P_{S|X_{p} = p^{0}}(S = s) \left[\log \frac{P_{S|X_{p} = p^{0}}(S = s)P_{X_{p}}(X_{p} = p^{0})}{Q_{S}(S = s)P_{X_{p}}(X_{p} = p^{0})}\right] \\ &+ \sum_{i=1}^{k}\sum_{j=k+1}^{n} P_{X_{p}}(X_{p} = p^{ij})\sum_{s} P_{S|X_{p} = p^{ij}}(S = s) \left[\log \frac{P_{S|X_{p} = p^{ij}}(S = s)P_{X_{p}}(X_{p} = p^{ij})}{Q_{S}(S = s)P_{X_{p}}(X_{p} = p^{ij})}\right] \\ &= \frac{1}{k(n-k)+1}\sum_{s} P_{S|X_{p} = p^{0}}(S = s) \left[\log \frac{P_{S|X_{p} = p^{0}}(S = s)}{Q_{S}(S = s)}\right] \\ &+ \frac{1}{k(n-k)+1}\sum_{i=1}^{k}\sum_{j=k+1}^{n}\sum_{s} P_{S|X_{p} = p^{ij}}(S = s) \left[\log \frac{P_{S|X_{p} = p^{ij}}(S = s)}{Q_{S}(S = s)}\right] \\ &= \frac{1}{k(n-k)+1}D\left(P_{S|X_{p} = p^{0}} \parallel Q_{S}\right) + \frac{1}{k(n-k)+1}\sum_{i=1}^{k}\sum_{j=k+1}^{n} D\left(P_{S|X_{p} = p^{ij}} \parallel Q_{S}\right), \end{split}$$

where \sum_{s} is summation over all possible *n* binary sequences of length *T*. Here $P_{S|X_p=p^0}$ is the distribution of *S* with the rate assignment p^0 , and $P_{S|X_p=p^{ij}}$ is the distribution of *S* with the rate assignment p^{ij} . Choosing Q_S to be the distribution of *S* with rate assignment p^0 defined in (8), then for any $i = 1, \dots, k$ and $j = k + 1, \dots, n$, we have

$$D\left(P_{S|X_{p^{ij}}} \parallel Q_{S}\right) = T(d(r_{2} \parallel r_{1}) + d(r_{1} \parallel r_{2}))$$

Therefore,

$$I(X_{p} \parallel S) \leq \frac{1}{k(n-k)+1} \sum_{i=1}^{k} \sum_{j=k+1}^{n} T(d(r_{2} \parallel r_{1}) + d(r_{1} \parallel r_{2}))$$

$$\leq T(d(r_{2} \parallel r_{1}) + d(r_{1} \parallel r_{2})).$$

D Proof of Lemma 13

By the activation rules in Algorithm 1, we know that

$$S_{t_0+m} = \begin{cases} 1, \text{ if } (b-1)\mathbf{1}_{\{S_{t_0+m-1}(v_i)=1\}} + \sum_{r=1}^m \left(\mathbf{1}_{\{V_{t_0+m-r}>0\}} - m\mathbf{1}_{\{V_{t_0+m-r}\leq-1\}}\right) > b;\\ 0, \text{ otherwise.} \end{cases}$$

As all input neurons are quiescent at time t_0 and remain to be quiescent for all $t \ge t_0$, it follows that

$$(b-1)\mathbf{1}_{\{S_{t_0+m-1}(v_i)=1\}} + \sum_{r=1}^{m} \left(\mathbf{1}_{\{V_{t_0+m-r}>0\}} - m\mathbf{1}_{\{V_{t_0+m-r}\leq-1\}}\right)$$
$$= (b-1)\mathbf{1}_{\{S_{t_0+m-1}(v_i)=1\}} - m\sum_{r=1}^{m}\mathbf{1}_{\{V_{t_0+m-r}\leq-1\}}$$
$$\leq b-1 < b.$$

Thus, $S_{t_0+m}(v_i) = 0$ for all $i = 1, \dots, n$. So we have $V_{t_0+m+1}(v_i) = 0$ for all $i = 1, \dots, n$, which again implies that $S_{t_0+m+1}(v_i) = 0$ for all $i = 1, \dots, n$. Therefore, we conclude that $S_t(v_i) = 0$ and $V_t(v_i) = 0$ for all $t > t_0 + m$.

E Proof of Theorem 4

Proof of Theorem 4. We prove this via a genie-aided argument [JB67] by assuming that there is a genie that can access the firing sequences of all the n input neurons. By assuming the existence of a genie, we are essentially considering the centralized setting. Clearly, if the error probability is high even in the centralized setting, then no SNNs (which are distributed algorithms) can achieve lower error probability.

Suppose that $T \leq ((1-\delta)\log(k(n-k)+1)-1)T_{\mathcal{R}}$. By (7) there exists r_1, r_2 such that $r_1 \neq r_2$ and

$$T \le ((1 - \delta) \log(k(n - k) + 1) - 1) \frac{1}{d(r_2 \parallel r_1) + d(r_1 \parallel r_2)}.$$

Without loss of generality, assume that $r_1 > r_2$.

Consider the k(n-k) + 1 possible rate assignments defined in Lemma 3. Let \mathcal{P} be the set of such rate assignments. By Yao's minimax principle, we know the minimax probability of error is always lower bounded by Bayes probability of error with any prior distribution:

$$\max_{\boldsymbol{p}\in\mathcal{AR}_{k}}\mathbb{P}\left\{\widehat{\boldsymbol{win}}\left(\boldsymbol{S}\right)\neq\mathcal{W}(\boldsymbol{p})\right\}\geq\mathbb{E}_{X_{\boldsymbol{p}}\sim Unif(\mathcal{P})}\left[\mathbb{P}\left\{\widehat{\boldsymbol{win}}\left(\boldsymbol{S}\right)\neq\mathcal{W}(X_{\boldsymbol{p}})\right\}\right],$$

where $X_{\mathbf{p}} \sim Unif(\mathcal{P})$ is uniformly distributed over set \mathcal{P} . In addition, by Fano's inequality, we have

$$\mathbb{E}_{X_{\boldsymbol{p}} \sim Unif(\mathcal{P})}\left[\mathbb{P}\left\{\widehat{\boldsymbol{win}}\left(\boldsymbol{S}\right) \neq \mathcal{W}(X_{\boldsymbol{p}})\right\}\right] \geq 1 - \frac{I(X_{\boldsymbol{p}};\boldsymbol{S}) + 1}{\log(k(n-k)+1)}.$$
(11)

Applying Lemma 3, we get

$$\max_{\boldsymbol{p}\in\mathcal{AR}_{k}} \mathbb{P}\left\{\widehat{\boldsymbol{win}}\left(\boldsymbol{S}\right)\neq\mathcal{W}(\boldsymbol{p}\right)\right\}\geq1-\frac{I(X_{\boldsymbol{p}};\boldsymbol{S})+1}{\log(k(n-k)+1)}$$
$$\geq1-\frac{T\left(d(r_{2}\parallel r_{1})+d(r_{1}\parallel r_{2})\right)+1}{\log(k(n-k)+1)}$$
$$\geq\delta.$$

The last inequality holds as $T \leq ((1 - \delta) \log(k(n - k) + 1) - 1) T_{\mathcal{R}}$.

F Proof of Theorem 10

The proof of Theorem 10 uses the following technical fact and lemma.

Fact 17. For any given $p \in (0,1)$ and b > 0, let $f_{p,b} : \mathbb{R} \to \mathbb{R}$, defined as: for all t > 0,

$$f_{p,b}(t) := \exp\left(-td\left(\frac{b}{t} \parallel p\right)\right).$$

Function $f_{p,b}(\cdot)$ is increasing when $t \in (0, \frac{b}{p})$ and decreasing when $t \ge \frac{b}{p}$.

This fact follows immediately from a simple algebra.

Lemma 18. Assume $u, v \in [c, C] \subseteq (0, 1)$. Then for any $\alpha \in (0, 1)$,

$$d((1 - \alpha)u + \alpha v \parallel u) \ge \frac{\alpha^2 c(1 - C)}{2C(1 - c)} (d(u \parallel v) + d(v \parallel u)).$$

Proof. Note that for any fixed $q \in [c, C]$, $d(x \parallel q)$ is a function of x, where $x \in [c, C]$. In addition, by simple algebra, we have

$$d'(x \parallel q) = \log \frac{(1-q)x}{q(1-x)}, \text{ and } d''(x \parallel q) = \frac{1}{x(1-x)}.$$
(12)

By Taylor expansion, we have

$$d((1 - \alpha)u + \alpha v || u) = d(u || u) + ((1 - \alpha)u + \alpha v - u) d'(u || u) + \frac{((1 - \alpha)u + \alpha v - u)^2}{2} d''(\xi || u),$$

where $\xi \in [\min\{u, (1-\alpha)u + \alpha v\}, \max\{u, (1-\alpha)u + \alpha v\}]$. By (12),

$$d\left((1-\alpha)u + \alpha v \parallel u\right) = 0 + 0 + \frac{1}{\xi(1-\xi)}\frac{\alpha^2(u-v)^2}{2} \ge \frac{\alpha^2(u-v)^2}{2C(1-c)}$$

On the other hand, since $d(u \parallel v) + d(u \parallel v)$ is symmetric in u and v, without loss of generality, assume that $u \ge v$. We have

$$\begin{aligned} d\left(u \parallel v\right) + d\left(u \parallel v\right) &= (u - v) \log \frac{u(1 - v)}{v(1 - u)} \\ &= (u - v) \log \left(1 + \frac{u - v}{v(1 - u)}\right) \\ &\leq (u - v) \frac{u - v}{v(1 - u)} = \frac{(u - v)^2}{v(1 - u)} \leq \frac{(u - v)^2}{c(1 - C)} \\ &\leq \frac{2C(1 - c)}{c(1 - C)\alpha^2} d\left((1 - \alpha)u + \alpha v \parallel u\right), \end{aligned}$$

proving the lemma.

Now we are ready to prove Theorem 10.

Proof of Theorem 10. Without loss of generality, assume that

$$p_1 \ge \cdots \ge p_k > p_{k+1} \ge \cdots \ge p_n$$

For a given rate assignment $p \in \mathcal{AR}$, define $\tau_1, \tau_2, \cdots, \tau_n$ as

$$\tau_i := \inf_t \left\{ t : \sum_{r=0}^{\min\{t, m^*\}} S_r(u_i) \ge b \right\}, \quad \forall \ i = 1, \cdots, n.$$

To show Theorem 10, it is enough to show that with probability $1 - \delta$,

$$\tau_i < \tau_j \quad \forall \ i = 1, \cdots, k, \text{ and } j = k+1, \cdots, n;$$

$$(13)$$

and
$$\tau_i \le m^* \quad \forall \ i = 1, \cdots, k..$$
 (14)

Before diving into proving (13) and (14) hold with probability at least $1 - \delta$, let's check the sufficiency of (13) and (14). Let $t_0 := \max_{1 \le i \le k} \tau_i$. Let \mathcal{E} be the event on which (13) and (14) hold. Clearly, conditioning on event \mathcal{E} , we have

$$\left(\max_{1 \le i \le k} \tau_i \mid \mathcal{E}\right) = t_0 \mid \mathcal{E} \le m^* - 1 \le m - 1,$$

and

$$\left(\max_{1 \le i \le k} \tau_i \mid \mathcal{E}\right) = t_0 \mid \mathcal{E} < \tau_j \mid \mathcal{E} \quad \forall j = k+1, \cdots, n$$

Notably, for any $t \leq t_0 \leq m-1$ and for $i = 1, \dots, n$,

$$\left[\sum_{r=1}^{t} \mathbf{1}_{\{V_r(v_i)>0\}} - m \sum_{r=1}^{t} \mathbf{1}_{\{V_r(v_i)\leq -1\}}\right]_{+} \le \sum_{r=1}^{t} \mathbf{1}_{\{V_r(v_i)>0\}} \le \sum_{r=1}^{t} S_r(u_i)$$

Thus, conditioning on \mathcal{E} , at most k-1 output neurons ever spike by time t_0 . So we have (1) $\mathbf{1}_{\{V_t(v_i) \leq -1\}} = 0$, and (2) $\mathbf{1}_{\{V_t(v_i) > 0\}} = S_t(u_i)$, for all $i = 1, \dots, n$ and for all $t \leq t_0$. In addition, we have for all $t \leq t_0$,

$$(b-1)\mathbf{1}_{\{S_t(v_i)=1\}} + \left[\sum_{r=1}^t \mathbf{1}_{\{V_r(v_i)>0\}} - m \sum_{r=1}^t \mathbf{1}_{\{V_r(v_i)\leq-1\}}\right]_+$$

= $(b-1)\mathbf{1}_{\{S_t(v_i)=1\}} + \sum_{r=1}^t \mathbf{1}_{\{V_r(v_i)>0\}}$
= $(b-1)\mathbf{1}_{\{S_t(v_i)=1\}} + \sum_{r=1}^t S_r(u_i).$

By the activation rules in Algorithm 1, we know, conditioning on \mathcal{E} , at time $t_0 + 1 \leq m^*$, output neurons v_1, \dots, v_k spike simultaneously, and output neurons v_{k+1}, \dots, v_n do not spike, proving (1) in Theorem 10. By the choice of t_0 , we know that, on \mathcal{E} , $t_0 + 1$ is the first time that k output neurons spike simultaneously, and no other k output neurons ever spike simultaneously, proving (2) in Theorem 10.

By a simple induction argument, it can be shown that conditioning on \mathcal{E} , in each of the time slot t such that $t_0 + 1 \leq t \leq m + 1$, output neurons v_1, \dots, v_k spike, and no other output neurons (i.e., output neurons v_{k+1}, \dots, v_n do not spike). Let's consider the case when t = (m+1) + 1. As among output neurons, only v_1, \dots, v_k spike, and no other output neurons spike for any $t' \leq m + 1$, it follows that

$$m\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i) \le -1\}} = 0, \quad \forall \ v_1, \cdots, v_k.$$

Thus, for these k output neurons,

$$\begin{aligned} (b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq-1\}}\right]_{+} \\ &= (b-1) + \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} \\ &= (b-1) + \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} + \mathbf{1}_{\{V_{t-1}(v_i)>0\}} - \mathbf{1}_{\{V_{t-1-m}(v_i)>0\}} \\ &\geq b-2 + \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} \\ &= b-2 + \sum_{r=1}^{m} S_r(u_i) \geq 2b-2 \geq b, \end{aligned}$$

where the last inequality holds as long as $b \ge 2$. For output neurons v_{k+1}, \cdots, v_n , we have

$$(b-1)\mathbf{1}_{\{S_{t-1}(v_i)=1\}} + \left[\sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}} - m \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)\leq-1\}}\right]$$

$$\leq \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-r}(v_i)>0\}}$$

$$\stackrel{(a)}{=} \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} + \mathbf{1}_{\{V_{t-1}(v_i)>0\}} - \mathbf{1}_{\{V_{t-1-m}(v_i)>0\}}$$

$$= \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} - \mathbf{1}_{\{V_{t-1-m}(v_i)>0\}}$$

$$\leq \sum_{r=1}^{m} \mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} = \sum_{r=1}^{m} \mathbf{1}_{\{V_{r}(v_i)>0\}} < b.$$

+

Equality (a) follows because at time t - 1, output neurons v_1, \dots, v_k spike, resulting in $\mathbf{1}_{\{V_{t-1-r}(v_i)>0\}} = 0$ for $i \neq 1, \dots, k$. Thus, we know conditioning on event \mathcal{E} , at time (m+1)+1, the output neurons v_1, \dots, v_k spike, and no other output neuron spike. It can be shown by a simple induction that at each time t such that $t_0 + 1 \leq t \leq m+b$, the output neurons v_1, \dots, v_k spike, and no other output neurons spike. This proves (3) in Theorem 10.

Next we prove (13) and (14). By definition of τ_j , we know that $\tau_j \leq m^*$ for all $j = 1, \dots, n$. Thus, we only need to show that with probability $1 - \delta$,

$$\tau_i < \tau_j \quad \forall \ i = 1, \cdots, k, \text{ and } j = k+1, \cdots, n_j$$

which is the focus of the remainder of our proof.

Note that

$$\mathbb{P}\left\{\tau_{i} < \tau_{j}, \quad \forall i \in \{1, \cdots, k\}, \forall j \in \{k+1, \cdots, n\}\right\} \\
= \mathbb{P}\left\{\tau_{i} < \tau_{j}, \& \ \tau_{i} < m^{*}, \quad \forall i \in \{1, \cdots, k\}, \forall j \in \{k+1, \cdots, n\}\right\} \\
\geq 1 - \sum_{i=1}^{k} \sum_{j=k+1}^{n} \mathbb{P}\left\{\tau_{i} \ge \tau_{j}, \text{ or } \tau_{i} = m^{*}\right\}.$$
(15)

For each term in the summation of (15), we have

$$\mathbb{P}\left\{\tau_i \ge \tau_j, \text{ or } \tau_i = m^*\right\} = \mathbb{P}\left\{\tau_i = m^*\right\} + \mathbb{P}\left\{\tau_i \ge \tau_j, \& \tau_i < m^*\right\},\tag{16}$$

which follows from the fact that $\mathbb{P}\{A \cup B\} = \mathbb{P}\{A\} + \mathbb{P}\{B - A\}$ for any sets A and B. Note that $m^*p_i \ge b$. By Chernoff bound, the first term in (16) is bounded as

$$\mathbb{P}\left\{\tau_i = m^*\right\} = \mathbb{P}\left\{\sum_{r=0}^{m^*} S_r(u_i) \le b\right\} \le \exp\left(-m^* \cdot d\left(\frac{b}{m^*} \parallel p_i\right)\right).$$
(17)

For the second term in (16), we have

$$\mathbb{P}\left\{\tau_i \ge \tau_j \text{ and } \tau_i < m^*\right\} = \mathbb{P}\left\{\sum_{r=0}^{\tau_i} S_r(u_j) \ge b, \text{ and } \tau_i < m^*\right\}$$
$$\le \exp\left(-t^* \cdot d\left(\frac{b}{t^*} \parallel p_{k+1}\right)\right) + \exp\left(-t^* \cdot d\left(\frac{b}{t^*} \parallel p_k\right)\right),$$

where $t^* \in \left(\frac{b}{p_{k+1}}, \frac{b}{p_k}\right)$. Thus, (16) is upper bounded as

$$\mathbb{P}\left\{\tau_{i} \geq \tau_{j}, \text{ or } \tau_{i} = m^{*}\right\} \leq \exp\left(-m^{*} \cdot d\left(\frac{b}{m^{*}} \parallel p_{k+1}\right)\right) + \exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k}\right)\right) + \exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k}\right)\right) \\ \leq \exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k+1}\right)\right) + 2\exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k}\right)\right).$$

Eq (15) is bounded as

$$\mathbb{P}\left\{\tau_{i} < \tau_{j}, \quad \forall i \in \{1, \cdots, k\}, \forall j \in \{k+1, \cdots, n\}\right\}$$

$$\geq 1 - \sum_{i=1}^{k} \sum_{j=k+1}^{n} \mathbb{P}\left\{\tau_{i} \geq \tau_{j}, \text{ or } \tau_{i} = m^{*}\right\}$$

$$\geq 1 - \sum_{i=1}^{k} \sum_{j=k+1}^{n} \left(\exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k+1}\right)\right) + 2\exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k}\right)\right)\right)$$

$$= 1 - k(n-k) \left(\exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k+1}\right)\right) + 2\exp\left(-t^{*} \cdot d\left(\frac{b}{t^{*}} \parallel p_{k}\right)\right)\right).$$

Let $t^* = \frac{b}{(p_k + p_{k+1})/2}$, it holds that

$$\exp\left(-t^* \cdot d\left(\frac{b}{t^*} \parallel p_{k+1}\right)\right) = \exp\left(-\frac{b}{(p_k + p_{k+1})/2} \cdot d\left(\frac{p_k + p_{k+1}}{2} \parallel p_{k+1}\right)\right),$$
$$2\exp\left(-t^* \cdot d\left(\frac{b}{t^*} \parallel p_k\right)\right) = 2\exp\left(-\frac{b}{(p_k + p_{k+1})/2} \cdot d\left(\frac{p_k + p_{k+1}}{2} \parallel p_k\right)\right).$$

By Lemma 18, we know

$$d\left(\frac{p_k + p_{k+1}}{2} \parallel p_{k+1}\right) \ge \frac{c(1-C)}{8C(1-c)} \left(d(p_{k+1} \parallel p_k) + d(p_k \parallel p_{k+1})\right),$$

and,

$$d\left(\frac{p_k + p_{k+1}}{2} \parallel p_k\right) \ge \frac{c(1-C)}{8C(1-c)} \left(d(p_{k+1} \parallel p_k) + d(p_k \parallel p_{k+1})\right).$$

Thus, we get

$$\mathbb{P}\left\{\tau_{i} < \tau_{j}, \quad \forall i \in \{1, \cdots, k\}, \forall j \in \{k+1, \cdots, n\}\right\}$$

$$\geq 1 - 3k(n-k) \exp\left(-\frac{2b}{p_{k} + p_{k+1}} \frac{c(1-C)}{8C(1-c)} \left(d(p_{k} \parallel p_{k+1}) + d(p_{k+1} \parallel p_{k})\right)\right)$$

Since $b = \frac{8C^2(1-c)}{c(1-C)} \left(\log \frac{3}{\delta} + \log k(n-k) \right) T_{\mathcal{R}}$, we have

$$3k(n-k)\exp\left(-\frac{2b}{p_k+p_{k+1}}\frac{c(1-C)}{8C(1-c)}\left(d(p_k \parallel p_{k+1}) + d(p_{k+1} \parallel p_k)\right)\right) \le \delta.$$

Thus, $\mathbb{P}\left\{\tau_i < \tau_j, \quad \forall i \in \{1, \cdots, k\}, \forall j \in \{k+1, \cdots, n\}\right\} \leq 1 - \delta$. In addition,

 $t^* = \frac{2b}{p_k + p_{k+1}} \le \frac{1}{c}b = m^* \le m,$

completing the proof of Theorem 10.