

# Why Pruning and Conditional Computation Work: A High-Dimensional Perspective

**Erdem Koyuncu**

*University of Illinois Chicago*

EKOYUNCU@UIC.EDU

## Abstract

We analyze the processes of pruning and conditional computation for the case of a single neuron in the asymptotic learning regime of large input dimension and training set size. For this purpose, we introduce conditional neurons, which implement an early exit strategy at the neuron level. Specifically, a conditional neuron considers the local field induced by a subset of its inputs. If this sub-local field is strong enough, then the rest of the inputs are ignored, saving computation. Conditional neurons provide an archetype of the well-known early exit or conditional computation architectures. As such, we formally analyze their generalization performance to understand why conditional computation is so effective in preserving performance despite significantly reduced average amount of computation. In the process, we introduce a concentration theorem for one-shot neuron-wise pruning, which is recently popularized in the context of large language models.

## 1. Introduction

In the context of neural networks, conditional computation refers to the idea of adapting the network computations based on the inputs or intermediate features produced at different layers [3, 6, 8, 11]. A recent example is mixture of experts in transformers [13, 17, 19]. Of particular interest to this work is one of the simplest cases of conditional computation, which is commonly referred to as early exit networks [10, 14, 16, 22]. The concept of early exit involves the utilization of intermediate classifiers, which are located in non-terminal layers of a neural network. When an intermediate classifier exhibits sufficient confidence in its decisions, it can perform an “early exit,” bypassing the the subsequent layers and thus conserving computational resources. Previous works have shown that early exit architectures can significantly improve upon ordinary neural networks in terms of the trade-off between computational complexity and accuracy. However, a theoretical justification as to why early exit networks perform so well has remained elusive, which is the goal of the present work. To the best of our knowledge, our work is the first to formally study a conditional computation scheme. Thus, our results shed light into the performance of other conditional architectures as well.

The rest of this paper is organized as follows: In Section 2, we introduce the conditional perceptron. We prove a concentration theorem for pruning in Section 3. We provide our main generalization bounds for the conditional perceptron in Section 4. The proofs are provided in the appendices.

## 2. Conditional Perceptrons

In this section, we begin by introducing conditional perceptrons, which implement the idea of early exit at the neuron level. Conditional perceptrons thus serve as one of the simplest special cases of a deep neural network with an early exit.

Consider inputs  $x_1, \dots, x_n$  to a neuron with corresponding weights  $w_1, \dots, w_n$ , respectively. Let  $\mathbf{x} = [x_1 \cdots x_n]^T$  and  $\mathbf{w} = [w_1 \cdots w_n]^T$  be the corresponding input and weight vectors. The Heaviside step function can be defined as  $\sigma(v) = 1$  for  $v \geq 0$  and  $\sigma(v) = -1$  for  $v < 0$ . We recall that an ordinary, “unconditional” perceptron provides the output  $y_{uc} \triangleq \sigma(v_0)$ ,  $v_0 \triangleq \mathbf{w}^T \mathbf{x}$ , where the subscript “uc” stands for unconditional. To construct the conditional perceptron, let us order the weights  $w_1, \dots, w_n$  from the smallest to the largest in magnitude as  $|w_{i_1}| \leq \cdots \leq |w_{i_n}|$ , where  $i_1, \dots, i_n$  is a permutation of  $1, \dots, n$ . The largest  $n - k$  weights in magnitude, where  $k \in \{1, \dots, n\}$  is a design parameter, are thus given by  $w_{i_{k+1}}, \dots, w_{i_n}$ . Ignoring the inputs from the remaining weights, we consider the local field induced by a normalized version of these weights. Namely, we define  $v_1 \triangleq \mathbf{w}_{ee}^T \mathbf{x}$ , where  $\mathbf{w}_{ee} \triangleq f_{ee}^k(\mathbf{w}) \triangleq \frac{\|\mathbf{w}\|}{\|\mathbf{w}_p\|} \mathbf{w}_p$ , and  $\mathbf{w}_p \triangleq f_p^k(\mathbf{w}) \triangleq [(w_p)_1 \cdots (w_p)_n]^T$  with  $(w_p)_{i_j} = 0$ ,  $1 \leq j \leq k$ , and  $(w_p)_{i_j} = w_{i_j}$ ,  $k + 1 \leq j \leq n$ . The subscript “ee” indicates that  $\mathbf{w}_{ee}$  corresponds to the early exit weight vector of the conditional perceptron, and “p” indicates a pruned (but unnormalized) version of the weight vector. The normalization factor  $\|\mathbf{w}\|/\|\mathbf{w}_p\|$  ensures that both  $v_0$  and  $v_1$  have the same variance of  $\|\mathbf{w}\|^2$  when the input  $\mathbf{x}$  is considered random with zero mean and identity covariance.

The input-output relationship of the conditional perceptron is then finally expressed as

$$y_c \triangleq \begin{cases} \sigma(v_1), & |v_1| \geq \tau, \\ \sigma(v_0), & \text{otherwise,} \end{cases} \quad (1)$$

where  $\tau > 0$  is a threshold hyperparameter, and the subscript “c” stands for conditional. Ideally, the conditional perceptron wishes to achieve the same class decision  $\sigma(v_0)$  as the ordinary perceptron. However, by definition, if the local field  $|v_1|$  is large enough, then an “early exit” is performed with the class decision  $\sigma(v_1)$ . We expect that the ignored weights, since they are relatively small in magnitude, will be unable to sway this decision to the disagreement  $\sigma(v_1) \neq \sigma(v_0)$ . Thus,  $\sigma(v_1) = \sigma(v_0)$  holds for most inputs, at least when  $k$  and  $\tau$  are both chosen to be large enough.

The motivation behind the early exit mechanism is to preserve computational resources. To demonstrate this fact, let us calculate the total floating point operations (FLOPs) required to calculate the conditional perceptron output  $v_1$ . Since only  $n - k$  of the entries of  $\mathbf{w}_{ee}$  are non-zero,  $v_1$  can be calculated using  $n - k$  multiplications and  $n - k - 1$  additions, for a total of  $2(n - k) - 1$  floating point operations (FLOPs). Checking the condition  $|v_1| \geq \tau$  in (1) is a mere extra FLOP. Hence, if  $|v_1| \geq \tau$ , the conditional perceptron consumes  $2(n - k) + 1$  FLOPs, where we have assumed that another FLOP is spent on the activation function. Otherwise,  $1 + 2k$  more FLOPs need to be performed to obtain  $v_0$  out of  $v_1$  via the relationship  $v_0 = (\|\mathbf{w}_p\|/\|\mathbf{w}\|)v_1 + \sum_{j=1}^{n-k} w_{i_j} x_{i_j}$ . Hence, letting  $\mu_{uc}$  and  $\mu_c$  denote the FLOPs required to implement an ordinary perceptron and a conditional perceptron, respectively, we have  $\mu_{uc} \triangleq 2n$ , and

$$\mu_c \triangleq \begin{cases} 2(n - k) + 1, & |v_1| \geq \tau, \\ 2n + 2, & \text{otherwise.} \end{cases} \quad (2)$$

On average, we thus expect the FLOPs with a conditional perceptron to be notably less than the FLOPs with an ordinary unconditional perceptron, without any significant penalty in terms of the classification performance.

### 3. A Concentration Theorem for Pruning

The operation of the conditional perceptron in (1) is intimately related to weight pruning in neural networks. Indeed,  $\mathbf{w}_{ee}$  is a pruned version of  $\mathbf{w}$ , followed by normalization. Pruning is achieved by retaining only the top  $k$  components of  $\mathbf{w}$  with the highest magnitudes, while setting all other components to zero. Hence, a conditional perceptron first evaluates the local field using a pruned version of its weights. If this local field is confident enough, an early exit is performed. The main difference between the conditional perceptron and pruning is the former's switch to the unpruned full set of weights whenever the local field provided by the pruned weights is not confident enough.

Regardless of the particular emphasis, whether it is the conditional perceptron or sole pruning, the relationship between the original feature extracting vector  $\mathbf{w}$  and its pruned counterpart  $\mathbf{w}_{ee}$  becomes a crucial point of interest. In particular, there has been a lot of recent interest on neuron-wise pruning for large language models [7, 21]. This line of work assumes that fine-tuning or retraining of models is not possible after pruning due to computational complexity, a condition we also adopt in our setting. To understand the statistical properties of the similarity  $\mathbf{w}_{ee}^T \mathbf{w}$ , we assume  $\mathbf{w}$  is uniform on the unit hypersphere, and derive the corresponding statistics of  $\mathbf{w}_{ee}^T \mathbf{w}$ . We will consider the regime where the dimension of  $\mathbf{w}$  grows to infinity. In this regime of high dimensions, an unexpected outcome emerges. Assuming that a positive fraction of the weights is retained during pruning, we demonstrate that the similarity denoted as  $\mathbf{w}_{ee}^T \mathbf{w}$  tends to converge in probability towards a non-zero constant. In other words, it essentially becomes deterministic. We will show that a similar result holds for conventional pruning in the sense that  $\mathbf{w}_p^T \mathbf{w}$  also concentrates to its mean.

Our results will follow from a concentration theorem related to order statistics of Gamma random variables. Hence, we consider a general Gamma random variable  $G$  with PDF

$$f_G(x) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-x/\theta}, \quad x \geq 0, \quad (3)$$

where  $s$  and  $\theta$  are known as the shape and the scale parameters, respectively, and  $\Gamma(\cdot)$  is the Gamma function. Let  $G_{\leq y}$  denote the random variable  $G$  conditioned on  $G \leq y$ , where  $y \in [0, \infty)$ . The notation  $G_{\geq y}$  is defined similarly.

**Theorem 1** *Let  $G_1, G_2, \dots, G_n$  be independent and identically distributed Gamma random variables with shape  $s$  and scale  $\theta$ . Denote the corresponding order statistics by  $G_{(1)} \leq \dots \leq G_{(n)}$ . For a given  $0 < q < 1$ , let*

$$\Xi_n \triangleq \frac{G_{(\lceil qn \rceil + 1)} + \dots + G_{(n)}}{G_1 + \dots + G_n}. \quad (4)$$

Let  $\omega_q$  denote the unique real number such that  $\mathbb{P}(G \leq \omega_q) = q$ . Define the corresponding threshold

$$\tau_q(s, \theta) \triangleq \left( 1 + \frac{q\mathbb{E}[G_{\leq \omega_q}]}{(1-q)\mathbb{E}[G_{\geq \omega_q}]} \right)^{-1} \in [0, 1]. \quad (5)$$

Then, we have,  $\Xi_n \rightarrow \tau_q(s, \theta)$  as  $d \rightarrow \infty$ .

**Corollary 2** *Suppose  $\mathbf{W}$  is uniform on  $\mathbb{S}^n$ . Given  $q \in (0, 1)$ , let  $\mathbf{W}_p = f_p^{\lceil nq \rceil}(\mathbf{W})$  and  $\mathbf{W}_{ee} = f_{ee}^{\lceil nq \rceil}(\mathbf{W})$ . As  $n \rightarrow \infty$ , we have  $\mathbf{W}^T \mathbf{W}_{ee} \rightarrow \sqrt{\tau_q}$ , and  $\mathbf{W}^T \mathbf{W}_p \rightarrow \tau_q$ , where  $\tau_q \triangleq \tau_q(s = \frac{1}{2}, \theta = 2)$ .*

We make two observations: First, both the pruned vector  $\mathbf{W}_p$  and the early exit vector  $\mathbf{W}_{ee}$  remain at a constant angle or similarity with respect to the original weight vector asymptotically for large feature dimension. Second, as Fig. 1 shows, the similarity has a non-linear relationship with respect to the pruning rate, and remains very high even for large pruning rates. For example, even when 60% of the components of a unit norm weight vector are pruned, which corresponds to  $q = 0.4$ , after normalizing the pruned vector to unit norm, the resulting vector will have a similarity of roughly 0.982, or an angle of only  $10.8^\circ$ .

The phenomenon of concentration towards large similarities, even at high pruning rates, appears to be fundamental. This helps explain why pruned networks often perform nearly as well as their unpruned versions, a fact frequently noted in existing research. To elaborate further, we shall proceed in an informal fashion: Pruning a deep neural network at a certain rate would be roughly equivalent to pruning all its neurons at the same rate. In the pruned network, all neurons would approximately operate in the same manner as if they were in the unpruned network, provided that the pruning rate is not very high. This is because the weights of the pruned and unpruned networks then have a high similarity as a result of Corollary 2. Consequently, we anticipate the pruned network’s performance to closely match that of the unpruned network. Next, we turn our attention back to the case of a single neuron, and in particular, the conditional perceptron, which is much easier to analyze in a formal fashion.

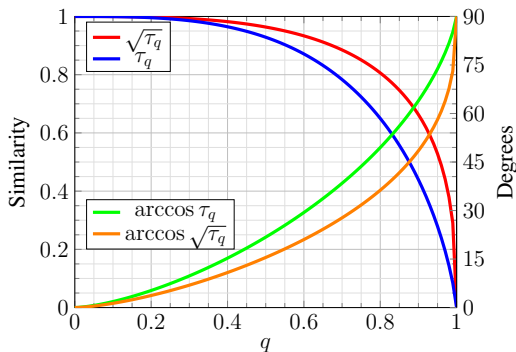


Figure 1: Concentrates for different sparsities.

#### 4. Generalization Performance of Conditional Perceptrons

In this section, we analyze the generalization error of the conditional perceptron in the classical student-teacher framework of learning theory, building on the concentration results in Section 3. To the best of our knowledge, this represents the first instance of a generalization analysis for a neural conditional computation or early exit system in the literature.

##### 4.1. Learning on the Unconditional Perceptron

We first provide an overview of learning on an ordinary, unconditional perceptron: Consider a dataset of  $N_t$  training vectors  $\{\mathbf{y}_1, \dots, \mathbf{y}_{N_t}\} \subset \mathbb{S}^n$  and a teacher  $\mathbf{t} \in \mathbb{S}^n$ . Given  $i \in \{1, \dots, N_t\}$ , let  $z_i \in \{-1, +1\}$  denote the desired output for training vector  $\mathbf{y}_i$ . The teacher determines the desired outputs in the sense that we set  $z_i = 1$  whenever  $\mathbf{t}^\dagger \mathbf{y}_i \geq 0$  and  $z_i = -1$  if  $\mathbf{t}^\dagger \mathbf{y}_i < 0$ .

We consider now a student  $\mathbf{w} \in \mathbb{R}^n$  acquired through some learning algorithm, as a function of only the input-output pairs  $(\mathbf{y}_1, z_1), \dots, (\mathbf{y}_{N_t}, z_{N_t})$ . The student weights coincide precisely with the “main” weights of the conditional perceptron as defined in Section 1, and hence we used the same notation. The specific learning algorithm to obtain  $\mathbf{w}$  out of the training vectors is not vital for our purposes. For example, the student can be chosen to be the vector that classifies the training data with the maximal margin, i.e.  $\mathbf{w} = \arg \max_{\mathbf{w}_0 \in \mathbb{S}^n} \min_{i \in \{1, \dots, N_t\}} z_i \mathbf{w}_0^T \mathbf{y}_i$ . The generalization error provided by the student is the probability  $P(\sigma(\mathbf{X}^T \mathbf{w}) \neq \sigma(\mathbf{X}^T \mathbf{t}))$  of mismatch between the student

and teacher decisions when the input  $\mathbf{X}$  to the perceptron is assumed to uniform on  $\mathbb{S}^n$ . For a fixed student and teacher vectors, the generalization error can simply be evaluated to be  $\frac{1}{\pi} \arccos \mathbf{w}^t \mathbf{t}$ .

We are often interested in the generalization error when averaged out over random datasets and teachers. For this purpose, suppose  $\mathbf{Y}_1, \dots, \mathbf{Y}_{N_t}, \mathbf{T} \sim N(\mathbf{0}_d, \mathbf{I}_d)$  are mutually independent. The student is then the random vector  $\mathbf{W} = \arg \max_{\mathbf{w}_0 \in \mathbb{S}^n} \min_{i \in \{1, \dots, N_t\}} \sigma(\mathbf{T}^T \mathbf{Y}_i) \mathbf{w}_0^T \mathbf{Y}_i$ , and the generalization error is given by  $\epsilon_{\text{uc}} \triangleq \mathbb{P}(\sigma(\mathbf{X}^T \mathbf{W}) \neq \sigma(\mathbf{X}^T \mathbf{T})) = \frac{1}{\pi} \mathbb{E}[\arccos \mathbf{W}^T \mathbf{T}]$ . It is difficult to calculate the generalization error exactly except for a few special cases. A special case is when both the number  $k$  of training vectors as well as the ambient dimension  $n$  grows to infinity. In other words, we have  $n, N_t \rightarrow \infty$ . If  $\alpha \triangleq \lim_{n \rightarrow \infty} \frac{N_t}{n}$  exists, then it is known [18, 20] that there is a constant  $C > 0$  such that  $\epsilon_{\text{uc}} \sim \bar{\epsilon}_{\text{uc}} = \frac{C}{\alpha}$  as  $\alpha \rightarrow \infty$ . Here, we used the notation  $\bar{\epsilon}_{\text{uc}}$  denote the asymptotic  $n, N_t \rightarrow \infty$  generalization error for the unconditional perceptron.

## 4.2. Analyzing the Conditional Perceptron

Let us now analyze the generalization performance of the conditional perceptron. Consider the same learning formulation as in Section 4.1. Define the early exit vector  $\mathbf{W}_{\text{ee}} \triangleq f_{\text{ee}}^{[nq]}(\mathbf{W})$ , where  $0 < q < 1$ . We extend the definition in (1) to random variables via  $Y_c = \sigma(\mathbf{W}_{\text{ee}}^T \mathbf{X})$  if  $|\mathbf{W}_{\text{ee}}^T \mathbf{X}| \geq \tau$ , and  $Y_c = \sigma(\mathbf{W}^T \mathbf{X})$  if  $|\mathbf{W}_{\text{ee}}^T \mathbf{X}| < \tau$ . The generalization error of the conditional perceptron is  $\epsilon_c \triangleq \mathbb{P}(Y_c \neq \sigma(\mathbf{X}^T \mathbf{T}))$ . We expect  $\epsilon_c \geq \epsilon_{\text{uc}}$  as the conditional perceptron often operates with a pruned, normalized version  $\mathbf{W}_{\text{ee}}$  of the student vector  $\mathbf{W}$ . We anticipate that any decrease in accuracy will be offset by a reduction in computational demand. To determine the corresponding tradeoff between the average computation and accuracy, we calculate the average FLOPs for the conditional perceptron, by averaging out (2) over the random inputs  $\mathbf{X}$ . This yields

$$\bar{\mu}_c \triangleq (2(n-k) + 1) \mathbb{P}(|\mathbf{W}_{\text{ee}}^T \mathbf{X}| \geq \tau) + (2n+2)(1 - \mathbb{P}(|\mathbf{W}_{\text{ee}}^T \mathbf{X}| \geq \tau)). \quad (6)$$

It is convenient to normalize the FLOPs with respect to the FLOPs  $2n$  of the unconditional network in the  $n \rightarrow \infty$  asymptotic regime. For this purpose, we define  $\bar{\mu}'_c \triangleq \lim_{n \rightarrow \infty} \frac{\bar{\mu}_c}{2n}$ . Then, the maximum FLOPs is 1, spent by the unconditional perceptron. The following theorem provides a set of achievable pairs of FLOPs and generalization errors provided by the conditional perceptron.

**Theorem 3** *Let  $\epsilon > 0$ . For a given computation constraint  $\bar{\mu}'_c \leq 1 - \epsilon$ , a generalization error of  $\epsilon_{\text{uc}} + (\epsilon/q)^{\frac{1}{2(1-\rho^2)}}$  is achievable, where  $\rho = \cos(\arccos \sqrt{\tau q} + \pi \bar{\epsilon}_{\text{uc}})$ .*

Now, fix some “reasonable” sparsity rate  $q$  so that  $\rho$  is not too far from zero, and imagine  $\epsilon$  as the only variable. Then, the theorem shows that near the full FLOPs of 1, the system performance approaches the unconditional performance exponentially fast with rate  $O(\epsilon^A)$  for some constant  $A > 1$ . This helps explain why early exit networks do suffer significant performance loss despite one cuts of a decent chunk of the computation budget, i.e.  $\epsilon$  is not close to zero. In fact, the exponent grows to infinity if  $\rho \rightarrow 0$ . Our analysis suggests this can happen if the sparsity rate nears zero and the training rate is high, so that  $\bar{\epsilon}_{\text{uc}} \rightarrow 0$ . A similar situation arises in the context of pruning and the behavior of similarity scores in Fig. 1: They remain relatively unchanged up to pruning rates of 0.4, and only then begin to decrease. General conditional computation networks, such as mixtures of experts, exhibit a similar behavior: the network performance remains robust unless the computation budget is overly restricted. The results presented in this paper also shed light on this phenomenon.

## Acknowledgments

This work was supported in part by the Army Research Lab (ARL) under Grant W911NF-21-2-0272, in part by the Army Research Office (ARO) under Grant W911NF-24-1-0049, and in part by the National Science Foundation (NSF) under Grant CNS-2148182.

## References

- [1] Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55. US Government printing office, 1948.
- [2] Mohammad Ahsanullah, Valery B Nevzorov, Mohammad Shakil, Mohammad Ahsanullah, Valery Nevzorov, and Mohammad Shakil. Conditional distributions of order statistics. *An Introduction to Order Statistics*, pages 51–60, 2013.
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [4] Arnaud Buhot, Juan-Manuel Torres Moreno, and Mirta B Gordon. Finite size scaling of the bayesian perceptron. *Physical Review E*, 55(6):7434, 1997.
- [5] Seok-Ho Chang, Pamela C Cosman, and Laurence B Milstein. Chernoff-type bounds for the gaussian error function. *IEEE Transactions on Communications*, 59(11):2939–2944, 2011.
- [6] Andrew Davis and Itamar Arel. Low-rank approximations for conditional feedforward computation in deep neural networks. *arXiv preprint arXiv:1312.4461*, 2013.
- [7] Lucio Dery, Steven Kolawole, Jean-Francois Kagey, Virginia Smith, Graham Neubig, and Ameet Talwalkar. Everybody prune now: Structured pruning of llms with only forward passes. *arXiv preprint arXiv:2402.05406*, 2024.
- [8] David Eigen, Marc Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013.
- [9] Carl-Gustav Esseen. On the liapunoff limit of error in the theory of probability. *Arkiv för Matematik, Astronomi och Fysik*, A28:1–19, 1942. ISSN 0365-4133.
- [10] Alperen Gormez, Venkat Dasari, and Erdem Koyuncu. Class means as an early exit decision mechanism. In *IJCNN*, 2022.
- [11] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7436–7456, 2021.
- [12] James J Heckman and Bo E Honore. The empirical content of the roy model. *Econometrica: Journal of the Econometric Society*, pages 1121–1149, 1990.
- [13] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

- [14] Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. Shallow-deep networks: Understanding and mitigating network overthinking. In *International Conference on Machine Learning*, June 2019.
- [15] Erdem Koyuncu and Hamid Jafarkhani. Distributed beamforming in wireless multiuser relay-interference networks with quantized feedback. *IEEE Transactions on Information Theory*, 58(7):4538–4576, 2012. doi: 10.1109/TIT.2012.2191708.
- [16] Yuhang Li, Tamar Geller, Youngeun Kim, and Priyadarshini Panda. Seenn: Towards temporal spiking early exit neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- [17] Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- [18] Manfred Opper. Statistical mechanics of learning: Generalization. *The handbook of brain theory and neural networks*, pages 922–925, 1995.
- [19] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- [20] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari Morcos. Beyond neural scaling laws: beating power law scaling via data pruning. *Advances in Neural Information Processing Systems*, 35:19523–19536, 2022.
- [21] Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. *arXiv preprint arXiv:2306.11695*, 2023.
- [22] Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *2016 23rd international conference on pattern recognition (ICPR)*, pages 2464–2469. IEEE, 2016.
- [23] Shoufang Xu and Yu Miao. Limit behaviors of the deviation between the sample quantiles and the quantile. *Filomat*, 25(2):197–206, 2011.

## Appendix A. Proof of Theorem 1

### A.1. Some auxiliary results on Gamma random variables

We begin by presenting some auxiliary results concerning Gamma random variables that we will need to prove the theorem. As defined in Section 3, a general Gamma random variable  $G$  has PDF  $f_G(x) = x^{s-1}e^{-x/\theta}/\Gamma(s)/\theta^s$ ,  $x \geq 0$ , where  $s$  and  $\theta$  are known as the shape and the scale parameters, respectively, and  $\Gamma(\cdot)$  is the Gamma function. The lower and upper incomplete gamma functions are defined as  $\gamma(s, x) \triangleq \int_0^x t^{s-1}e^{-t}dt$  and  $\Gamma(s, x) \triangleq \int_x^\infty t^{s-1}e^{-t}dt$ , respectively.

We note the asymptotic expressions [1, Eqs. 6.5.29 and 6.5.32]

$$\gamma(s, x) \sim x^s/s, \quad x \rightarrow 0, \quad (7)$$

$$\gamma(s, x) \rightarrow \Gamma(s), \quad x \rightarrow \infty, \quad (8)$$

$$\Gamma(s, x) \rightarrow \Gamma(s), \quad x \rightarrow 0, \quad (9)$$

$$\Gamma(s, x) \sim x^{s-1}e^{-x}, \quad x \rightarrow \infty. \quad (10)$$

According to (7), we have

$$P(G \leq y) = \frac{1}{\Gamma(s)}\gamma(s, y/\theta) \sim \frac{(y/\theta)^s}{\Gamma(1+s)}, \quad y \rightarrow 0, \quad (11)$$

and by (10), we obtain,

$$P(G \geq y) = \frac{1}{\Gamma(s)}\Gamma(s, y/\theta) \sim \frac{(y/\theta)^{s-1}e^{-y/\theta}}{\Gamma(s)}, \quad y \rightarrow \infty, \quad (12)$$

Let  $G_{\leq y}^a$  denote a truncated Gamma random variable, obtained by conditioning  $G$  on the event  $G \leq y$ , where  $y \in [0, \infty)$ . The notation  $G_{\geq y}^a$  is defined similarly. A straightforward calculation reveals that for  $a \geq 0$ , we have

$$EG_{\leq y}^a = \frac{\int_0^y x^a f_G(x) dx}{\int_0^y f_G(x) dx} = \frac{\int_0^y x^{a+s-1} e^{-x/\theta} dx}{\int_0^y x^{s-1} e^{-x/\theta} dx} = \frac{\int_{-\infty}^{y/\theta} u^{a+s-1} \theta^{a+s} e^{-u} du}{\int_{-\infty}^{y/\theta} u^{s-1} \theta^s e^{-u} du} = \frac{\theta^a \gamma(a+s, y/\theta)}{\gamma(s, y/\theta)}, \quad (13)$$

and similarly,

$$EG_{\geq y}^a = \frac{\theta^a \Gamma(a+s, y/\theta)}{\Gamma(s, y/\theta)}. \quad (14)$$

In (13), the first equality is by definition. To obtain the second equality, we substituted the PDF of  $G$ . The third equality is by a change of variables  $u = x/\theta$ . The final equality is by the definition of the lower incomplete Gamma function. For  $a = 1$ , using (7), we can then obtain

$$EG_{\leq y}^a \sim \frac{\theta^a (y/\theta)^{a+s}/(a+s)}{(y/\theta)^s/s} = \frac{sy^a}{a+s}, \quad y \rightarrow 0, \quad (15)$$

and using (8) yields

$$EG_{\leq y}^a \rightarrow \frac{\theta^a \Gamma(a+s)}{\Gamma(s)}, \quad y \rightarrow \infty. \quad (16)$$



In a similar vein, using (9), we have

$$\mathbb{E}G_{\geq y}^a \rightarrow \frac{\theta^a \Gamma(a+s)}{\Gamma(s)}, y \rightarrow 0. \quad (17)$$

and applying (10), we obtain

$$\mathbb{E}G_{\geq y}^a \sim \frac{\theta^a (y/\theta)^{a+s-1} e^{-y/\theta}}{(y/\theta)^{s-1} e^{-y/\theta}} = y^a, y \rightarrow \infty. \quad (18)$$

We begin with a useful lemma on the expected values of truncated Gamma random variables.

**Lemma 4** *The derivative of the function*

$$y \mapsto \mathbb{E}[G_{\leq y}]. \quad (19)$$

*is bounded.*

**Proof** We have

$$\mathbb{E}[G_{\leq y}] = \frac{\int_0^y x f(x) dx}{\int_0^y f(x) dx} \quad (20)$$

Hence, by the fundamental theorem of calculus,

$$\frac{d\mathbb{E}[G_{\leq y}]}{dy} = \frac{y f(y) P(G \leq y) - \mathbb{E}[G_{\leq y}] P(G \leq y) f(y)}{P^2(G \leq y)} \quad (21)$$

$$= \frac{f(y)}{P(G \leq y)} (y - \mathbb{E}[G_{\leq y}]) \quad (22)$$

As  $y \rightarrow 0$ , according to (3), (11), and (15), we obtain

$$\lim_{y \rightarrow 0} \frac{d\mathbb{E}[G_{\leq y}]}{dy} = \frac{s}{1+s} \quad (23)$$

On the other hand, by (3) and (16), we have

$$\lim_{y \rightarrow \infty} \frac{d\mathbb{E}[G_{\leq y}]}{dy} = 0 \quad (24)$$

The statement of the lemma then follows from the continuity of  $y \mapsto \mathbb{E}[G_{\leq y}]$ . ■

The following useful lemma is a standard bound on powers of linear functions.

**Lemma 5** ([15, Lemma 7]) *For any real numbers  $x_1, \dots, x_n \geq 0$ , we have*

$$\left(\sum_{i=1}^n x_i\right)^\beta \leq n^{\beta-1} \sum_{i=1}^n x_i^\beta. \quad (25)$$

Our main technical result, Theorem 1, shows that a certain ratio related to order statistics of Gamma random variables converges in probability to a threshold given by (5). The following lemma shows that the derivatives of a more general form of the threshold function is bounded from above.

**Lemma 6** Let  $q \in (0, 1)$ ,  $n \geq 1$  and  $k = \lceil nq \rceil$ . Define the function

$$h(x) \triangleq \frac{(n-k)\mathbb{E}[G_{\geq x}]}{k\mathbb{E}[G_{\leq x}] + (n-k)\mathbb{E}[G_{\geq x}]}.$$
 (26)

There is a constant  $C_1 > 0$  that is independent of  $n$  such that  $|\frac{dh}{dx}| < C_1$ ,  $\forall x \in \mathbb{R}$  and for all large enough  $n$ .

**Proof** Let  $\nu(x) = \mathbb{E}G_{\leq x}/\mathbb{E}G_{\geq x}$ . We can rewrite (80) as

$$h(x) = \left(1 + \frac{k}{n-k}\nu(x)\right)^{-1},$$
 (27)

The derivative of (27) is calculated to be

$$\frac{dh}{dx} = -\frac{k}{n-k} \left(1 + \frac{k}{n-k}\nu(x)\right)^{-2} \frac{d\nu}{dx}$$
 (28)

It follows that

$$\left|\frac{dh}{dx}\right| \leq \frac{k}{n-k} \left|\frac{d\nu}{dx}\right| \leq \frac{2q}{1-q} \left|\frac{d\nu}{dx}\right|.$$
 (29)

The inequality follows since  $\lim_{n \rightarrow \infty} \frac{k}{n-k} = \frac{q}{1-q}$ . We used twice the limit as an upper bound, which will be valid for every large enough  $n$ .

What is now left to show is that the derivative  $\frac{d\nu}{dx}$  is bounded. It is sufficient to prove that the limits  $\lim_{x \rightarrow 0} d\nu/dx$  and  $\lim_{x \rightarrow \infty} d\nu/dx$  exist and they are finite, and that  $d\nu/dx$  is continuous on  $(0, \infty)$ . First, we calculate the derivative. Let

$$N(x) \triangleq \int_0^x yf(y)dy \int_x^\infty f(y)dy = [\mathbb{E}[G_{\leq x}]P(G \leq x)]P(G \geq x),$$
 (30)

$$D(x) \triangleq \int_x^\infty yf(y)dy \int_0^x f(y)dy = [\mathbb{E}[G_{\geq x}]P(G \geq x)]P(G \leq x).$$
 (31)

We note the alternate representation  $\nu(x) = \frac{N(x)}{D(x)}$ . By the fundamental theorem of calculus, we obtain

$$\frac{dN}{dx} = xf(x)P(G \geq x) - \mathbb{E}[G_{\leq x}]P(G \leq x)f(x),$$
 (32)

$$\frac{dD}{dy} = -xf(x)P(G \leq x) + \mathbb{E}[G_{\geq x}]P(G \geq x)f(x).$$
 (33)

Using the division rule for derivatives, after some cumbersome but straightforward calculations, we can obtain

$$\frac{d\nu}{dx} = f(x) \frac{xP(G \geq x)\mathbb{E}[G_{\geq x}] + xP(G \leq x)\mathbb{E}[G_{\leq x}] - \mathbb{E}[G_{\geq x}]\mathbb{E}[G_{\leq x}][P(G \leq x) + P(G \geq x)]}{\mathbb{E}^2[G_{\geq x}]P(G \leq x)P(G \geq x)}$$
 (34)

$$= \frac{f(x)[xE[G] - \mathbb{E}[G_{\leq x}]\mathbb{E}[G_{\geq x}]]}{\mathbb{E}^2[G_{\geq x}]P(G \leq x)P(G \geq x)}.$$
 (35)

Since all functions involved in (35) are continuous,  $d\nu/dx$  is continuous except possibly at 0 and  $\infty$ . Using (3), (11), (15), (17), and noting that  $E[G] = k\theta$ , we obtain

$$\lim_{x \rightarrow 0} \frac{d\nu}{dx} = \frac{1}{\theta(1+s)} \quad (36)$$

Also, substituting (3), (12), (16), and (18), to (35), we can show that the derivative at  $\infty$  is zero. Together with (36), this shows that the derivative is bounded. This concludes the proof.  $\blacksquare$

Let us now derive upper and lower bounds on the variances of truncated Gamma random variables.

**Lemma 7** *For every  $x \geq 0$ , the variances of truncated Gamma random variables follow the bounds*

$$\max\{s, 1\}\theta^2 \geq \text{var}(G_{\geq x}) \geq \min\{s, 1\}\theta^2, \quad (37)$$

$$C_2 \geq \text{var}(G_{\leq x}), \quad (38)$$

where  $C_2$  is a constant that is independent of  $x$ .

**Proof** Assume that the shape parameter  $s$  of the Gamma random variable  $G$  satisfies  $s \in (0, 1]$ . Then,  $G$  is log-convex, and according to [12, Proposition 2], the function  $x \mapsto \text{var}(G_{\geq x})$  is monotonically increasing. In particular,  $\text{var}(G_{\geq x}) \geq \text{var}(G_{\geq 0}) = \text{var}(G) = s\theta^2$ . On the other hand, when  $s \in [1, \infty)$ , the density  $G$  is log-concave. In this case, [12, Proposition 1] shows that  $x \mapsto \text{var}(G_{\geq x})$  is monotonically decreasing. Hence, we have  $\text{var}(G_{\geq x}) \geq \lim_{x \rightarrow \infty} \text{var}(G_{\geq x})$ . In what follows, we calculate the limit. We have

$$\text{var}(G_{\geq x}) = E[G_{\geq x}^2] - E^2[G_{\geq x}] \quad (39)$$

$$= \theta^2 \frac{\Gamma(2+s, \frac{x}{\theta})\Gamma(s, \frac{x}{\theta}) - \Gamma^2(1+s, \frac{x}{\theta})}{\Gamma^2(s, \frac{x}{\theta})} \quad (40)$$

We have a  $\frac{0}{0}$  indeterminacy as  $x \rightarrow \infty$ . We can thus apply L'Hôpital's rule to obtain

$$\lim_{x \rightarrow \infty} \text{var}(G_{\geq x}) = \lim_{x \rightarrow \infty} \frac{x^2\Gamma(s, \frac{x}{\theta}) - 2\theta x\Gamma(1+s, \frac{x}{\theta}) + \theta^2\Gamma(2+s, \frac{x}{\theta})}{2\Gamma(s, \frac{x}{\theta})} \quad (41)$$

$$= \lim_{x \rightarrow \infty} \theta^s \frac{t\Gamma(1+s, \frac{x}{\theta}) - x\Gamma(s, \frac{x}{\theta})}{e^{-\frac{x}{\theta}} x^{s-1}} \quad (42)$$

$$= \lim_{x \rightarrow \infty} \frac{\theta^{1+s}\Gamma(s, \frac{x}{\theta})}{e^{-\frac{x}{\theta}} x^{-2+s}(\theta - s\theta + x)} \quad (43)$$

$$= \theta^2 \quad (44)$$

The second and the third equalities also follow from L'Hôpital's rule. In order to obtain the final equality, we have applied (10). Note that the derivatives of the upper incomplete Gamma function can be evaluated via the formulae  $\frac{d\Gamma(s, x)}{dx} = \frac{d}{dx} \int_x^\infty t^{s-1} e^{-t} dt = -x^{s-1} e^{-x}$ , by the fundamental theorem of calculus. Therefore, for any  $s$ , we obtain  $\text{var}(G_{\geq x}) \geq \min\{s, 1\}\theta^2$ .

Since  $\text{var}(G_{\leq x}) = E[G_{\leq x}^2] - [EG_{\leq x}]^2 \leq E[G_{\leq x}^2]$ , according to (15) and (16), the lower conditional variance  $\text{var}(G_{\leq x})$  is bounded by a constant that is independent of  $x$  from above. For the

upper conditional variance  $\text{var}(G_{\geq x})$ , we consider the cases of  $s \in (0, 1)$  and  $s \in [1, \infty)$  separately. In the former scenario, the monotonically increasing nature of  $x \mapsto \text{var}(G_{\geq x})$ , as established in [12, Proposition 2], in conjunction with (44), demonstrates that  $\text{var}(G_{\geq x}) \leq \theta^2$  for every  $x$ . For  $s \geq 1$ , we obtain  $\text{var}(G_{\geq x}) \leq \text{var}(G_{\geq 0}) = s\theta^2$ , according to [12, Proposition 1]. Hence, for any  $s$ , we have  $\text{var}(G_{\geq x}) \leq \theta^2 \max\{1, s\}$ . This concludes the proof of the lemma.  $\blacksquare$

As a corollary, we obtain lower and upper bounds on a linear combination of variances truncated Gamma random variables.

**Corollary 8** *Let  $q \in (0, 1)$ ,  $n \geq 1$  and  $k = \lceil nq \rceil$ . Let*

$$\sigma^2 \triangleq y^2(k-1)\text{var}(G_{\leq x}) + (1-y)^2(n-k)\text{var}(G_{\geq x}). \quad (45)$$

*Then, for every  $\tau \in (0, 1)$  and  $y \in (0, 1)$  with  $|y - \tau| \leq \frac{1-\tau}{2}$ , we have*

$$C_3 n \leq \sigma^2 \leq C_4 n \quad (46)$$

*for all sufficiently large  $n$ , where  $C_3$  and  $C_4$  are constants.*

**Proof** According to Lemma 7, for any  $s$ , we obtain

$$\sigma^2 \geq (1-y)^2(n-k)\text{var}(G_{\geq x}) \quad (47)$$

$$\geq (1-y)^2(n-k) \min\{s, 1\}\theta^2 \quad (48)$$

The bound  $|y - \tau| \leq \frac{1-\tau}{2}$  implies

$$\sigma^2 \geq \left(\frac{1-\tau}{2}\right)^2 (n-k) \min\{s, 1\}\theta^2 \quad (49)$$

Also, substituting  $k = \lceil nq \rceil$ , we obtain

$$\sigma^2 \geq \left(\frac{1-\tau}{2}\right)^2 (n - \lceil nq \rceil) \min\{s, 1\}\theta^2 \quad (50)$$

$$\geq \frac{1-q}{2} \left(\frac{1-\tau}{2}\right)^2 \min\{s, 1\}\theta^2 n, \quad (51)$$

for sufficiently large  $n$ . This proves the lower estimate on the variance.

For the upper estimate, we can first obtain

$$\sigma^2 \leq \underbrace{y^2}_{\leq 1} \underbrace{(k-1)}_{\leq n} \text{var}(G_{\leq x}) + \underbrace{(1-y)^2}_{\leq 1} \underbrace{(n-k)}_{\leq n} \text{var}(G_{\geq x}) = n[\text{var}(G_{\leq x}) + \text{var}(G_{\geq x})], \quad (52)$$

and applying Lemma 7 proves the upper estimate on  $\sigma^2$ . This concludes the proof of the corollary.  $\blacksquare$

The following lemma is utilized to bound the error terms resulting from the Berry-Esseen estimates.

**Lemma 9** Let  $q, y \in (0, 1)$ ,  $n \geq 1$  and  $k = \lceil nq \rceil$ . Let  $\sigma^2$  be as defined in (45) of Corollary 8, and

$$\rho \triangleq y^3(k-1)\mathbb{E}|G_{\leq x} - \mathbb{E}G_{\leq x}|^3 + (1-y)^3(n-k)\mathbb{E}|G_{\geq x} - \mathbb{E}G_{\geq x}|^3. \quad (53)$$

For every large enough  $n$ , we have

$$\int_0^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx \leq C_6 n^{-\frac{1}{2}} \quad (54)$$

for some constant  $C_6 > 0$  that is independent of  $y, k, n$ .

**Proof** Using the inequalities  $y \leq 1$  and  $k \leq n$ , we obtain

$$\rho \leq n\mathbb{E}|G_{\leq x} - \mathbb{E}G_{\leq x}|^3 + n\mathbb{E}|G_{\geq x} - \mathbb{E}G_{\geq x}|^3. \quad (55)$$

Applying Lemma 5 yields

$$\rho \leq 4n(\mathbb{E}G_{\leq x}^3 + \mathbb{E}^3G_{\leq x} + \mathbb{E}G_{\geq x}^3 + \mathbb{E}^3G_{\geq x}) \quad (56)$$

$$\leq 8n(\mathbb{E}G_{\geq x}^3 + \mathbb{E}^3G_{\geq x}) \quad (57)$$

$$\leq C_7 n(1+x^3), \quad (58)$$

where  $C_7$  is a constant that is independent of  $n$  and  $x$ . The last inequality is a consequence of (17) and (18).

It is a standard result in probability theory that the exact PDF for the  $k$ th order statistic  $G^{(k)}$  is given by

$$f_{G^{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} f_G(x) [F_G(x)]^{k-1} [1-F_G(x)]^{n-k} \quad (59)$$

$$\leq n2^{n-1} f_G(x) [1-F_G(x)]^{n-k} \quad (60)$$

Since  $F_G(x) \rightarrow 1$ , there exists  $x_0 > 0$  and  $c \in (0, 1)$  such that for every  $x \geq x_0$ , we have  $n2^{n-1}[1-F_G(x)]^{n-\lceil nq \rceil} \leq c^{-n}$ . As a result, we obtain  $f_{G^{(k)}}(x) \leq f_G(x)c^{-n}, \forall x \geq x_0$ . Combining this with (58) and the lower bound on  $\sigma$  in Corollary 8, we obtain

$$\int_0^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx = \int_0^{x_0} \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx + \int_{x_0}^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx \quad (61)$$

$$\leq \int_0^{x_0} \frac{C_7 n(1+x_0^3)}{(C_3 n)^{1.5}} f_{G^{(k)}}(x) dx + \int_{x_0}^\infty \frac{C_7 n(1+x^3)}{(C_3 n)^{1.5}} f_G(x) c^{-n} dx \quad (62)$$

$$\leq \frac{C_7(1+x_0^3)}{C_3^{1.5}} n^{-\frac{1}{2}} + \frac{C_7(1+x_0^3)}{C_3^{1.5}} (\mathbb{E}G + \mathbb{E}G^3) o(n^{-\frac{1}{2}}). \quad (63)$$

The proof is now complete since the moments of a Gamma random variable are also finite.  $\blacksquare$

### A.2. Other auxiliary results

As the first result, we recall the central limit theorem for quantiles.

**Proposition 10 ([23])** *Let  $X_1, X_2, \dots, X_n$  be IID copies of a random variable  $X$ . Given  $p \in (0, 1)$ , let  $\xi_p = \inf\{x : F_X(x) \geq 1 - p\}$ . Suppose  $F_X$  has a continuous first derivative  $f_X$  in the neighborhood of  $\xi_p$  and  $f(\xi_p) > 0$ . Then,*

$$\frac{\sqrt{n}f_X(\xi_p)}{\sqrt{p(1-p)}} \left( X_{(\lceil n(1-p) \rceil)} - \xi_p \right) \sim N(0, 1) \text{ as } n \rightarrow \infty. \quad (64)$$

We then present the general form of Berry-Esseen theorem for non-identically distributed random variables.

**Proposition 11 (Berry-Esseen Theorem [9])** *Let  $X_1, X_2, \dots$ , be independent random variables with  $E[X_i] = 0$  and  $E[|X_i|^3] < \infty$  for every  $i \in \mathbb{Z}_{>0}$ . Let  $\Phi$  denote the CDF of the standard normal distribution. For all  $n$ , we have*

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P} \left( \frac{X_1 + X_2 + \dots + X_n}{\sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2}} \leq x \right) - \Phi(x) \right| \leq 8 \left( \sum_{i=1}^n \text{var} X_i \right)^{-3/2} \sum_{i=1}^n E|X_i|^3. \quad (65)$$

### A.3. Proof of the theorem

We are now ready to prove Theorem 1. We organize the proof in four steps.

**Step-1:** In the first step, we approximate the cumulative distribution function of  $\Xi_n$  via a normal random variable using the Berry-Esseen theorem. Let  $y \in [0, 1]$  and  $k = \lceil nq \rceil$ . We have

$$\mathbb{P}(\Xi_n \leq y) = \mathbb{P}(G_{(k+1)} + \dots + G_{(n)} \leq y(G_{(1)} + \dots + G_{(n)})) \quad (66)$$

$$= \int_0^\infty \mathbb{P} \left( (1-y) \sum_{j=k+1}^n G_{(j)} - y \sum_{j=1}^{k-1} G_{(j)} \leq yx \mid G_{(k)} = x \right) f_{G_{(k)}}(x) dx. \quad (67)$$

According to [2], the two sets of order statistics  $(G_{(k+1)}, \dots, G_{(n)})$  and  $(G_{(1)}, \dots, G_{(k-1)})$  that appear in (67) are conditionally independent given  $G_{(k)} = x$ . Moreover, we have

$$\left[ \sum_{j=k+1}^n G_{(j)} \mid G_{(k)} = x \right] \sim \sum_{j=1}^{n-k} G_{\geq x, j}, \quad (68)$$

and

$$\left[ \sum_{j=1}^{k-1} G_{(j)} \mid G_{(k)} = x \right] \sim \sum_{j=1}^{k-1} G_{\leq x, j}, \quad (69)$$

where  $G_{\leq x, 1}, G_{\geq x, 1}, G_{\leq x, 2}, G_{\geq x, 2}, \dots$  is a sequence of IID random variables with

$$G_{\leq x, j} \sim G_{\leq x}, \quad (70)$$

and

$$G_{\geq x,j} \sim G_{\geq x} \quad (71)$$

for every  $j \in \mathbb{Z}_{>0}$ .

Let us now normalize the mean of (70) and (71) as

$$G'_{\leq x,j} \triangleq G_{\leq x,j} - \mathbb{E}[G_{\leq x,j}] = G_{\leq x,j} - \mathbb{E}[G_{\leq x}], \quad (72)$$

$$G'_{\geq x,j} \triangleq G_{\geq x,j} - \mathbb{E}[G_{\geq x,j}] = G_{\geq x,j} - \mathbb{E}[G_{\geq x}], \quad (73)$$

where  $j \in \mathbb{Z}_{>0}$ . We have

$$\mathbb{P}(\Xi_n \leq y) = \int_0^\infty \mathbb{P}\left((1-y) \sum_{j=1}^{n-k} G'_{\geq x,j} - y \sum_{j=1}^{k-1} G'_{\leq x,j} \leq y'\right) f_{G^{(k)}}(x) dx, \quad (74)$$

where

$$y' \triangleq yx + y(k-1)\mathbb{E}[G_{\leq x}] - (1-y)(n-k)\mathbb{E}[G_{\geq x}]. \quad (75)$$

By definition, the random variables  $G'_{\leq x,j}$  and  $G'_{\geq x,j}$  have zero mean. According to Lemma 7, they also have finite normalized moments. Hence, Proposition 11 is applicable, and we have

$$\mathbb{P}(\Xi_n \leq y) \leq \int_0^\infty \Phi(y'/\sigma) f_{G^{(k)}}(x) dx + 8 \int_0^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx, \quad (76)$$

where  $\Phi$  is the CDF of the standard normal random variable with zero mean and unit variance,

$$\sigma^2 \triangleq y^2(k-1)\text{var}(G_{\leq x}) + (1-y)^2(n-k)\text{var}(G_{\geq x}), \quad (77)$$

and

$$\rho \triangleq y^3(k-1)\mathbb{E}\left[\left|G_{\leq x} - \mathbb{E}[G_{\leq x}]\right|^3\right] + (1-y)^3(n-k)\mathbb{E}\left[\left|G_{\geq x} - \mathbb{E}[G_{\geq x}]\right|^3\right]. \quad (78)$$

Let us now rewrite the mean-normalized threshold  $y'$  defined in (75) as

$$y' = yx - y\mathbb{E}[G_{\leq x}] - \left[k\mathbb{E}[G_{\leq x}] + (n-k)\mathbb{E}[G_{\geq x}]\right] \left(h(x) - y\right), \quad (79)$$

where

$$h(x) \triangleq \frac{(n-k)\mathbb{E}[G_{\geq x}]}{k\mathbb{E}[G_{\leq x}] + (n-k)\mathbb{E}[G_{\geq x}]}. \quad (80)$$

**Step-2:** Let  $\omega_q = \inf\{x : F_G(x) \geq q\}$  be as defined in the theorem statement. Let us also recall the threshold

$$\tau \triangleq \tau_q(s, \theta) = \left(1 + \frac{q\mathbb{E}[G_{\leq \omega_q}]}{(1-q)\mathbb{E}[G_{\geq \omega_q}]}\right)^{-1} = \frac{(1-q)\mathbb{E}[G_{\geq \omega_q}]}{q\mathbb{E}[G_{\leq \omega_q}] + (1-q)\mathbb{E}[G_{\geq \omega_q}]}. \quad (81)$$

from the theorem statement. In this step, we find an upper bound on  $y'$  for  $y = \tau - \epsilon$  and  $|x - \omega_q| \in O(\epsilon)$ .

According to Lemma 6, we have

$$\sup_{x \in \mathbb{R}} \left| \frac{d}{dx} h(x) \right| \leq C_1, \quad (82)$$

for some constant  $C_1 > 0$  that is independent of  $n$ . As a result, if  $|x - \omega_q| \leq \delta \triangleq \frac{\epsilon}{2C_1}$ , we have

$$\left| h(x) - h(\omega_q) \right| < C_1 \delta = \frac{\epsilon}{2}. \quad (83)$$

In comparison with

$$h(\omega_q) = \frac{(n - \lceil nq \rceil) \mathbb{E}[G_{\geq \omega_q}]}{\lceil nq \rceil \mathbb{E}[G_{\leq \omega_q}] + (n - \lceil nq \rceil) \mathbb{E}[G_{\geq \omega_q}]}, \quad (84)$$

it follows that

$$\lim_{n \rightarrow \infty} h(\xi_q) = \tau \quad (85)$$

Hence, for large enough  $n$ , we can argue that  $|h(\omega_q) - \tau| \leq \frac{\epsilon}{4}$ . Combining with (83), for  $y = \tau - \epsilon$ , we can thus obtain

$$h(x) - y \geq h(\omega_q) - \frac{\epsilon}{2} - y \quad (86)$$

$$= h(\omega_q) - \tau + \frac{\epsilon}{2} \quad (87)$$

$$\geq \frac{\epsilon}{4}. \quad (88)$$

On the other hand, we have the lower bounds

$$k\mathbb{E}[G_{\leq x}] + (n - k) \underbrace{\mathbb{E}[G_{\geq x}]}_{\geq \mathbb{E}[G_{\leq x}]} \geq n\mathbb{E}[G_{\leq x}] \quad (89)$$

$$\geq n\mathbb{E}[G_{\leq \omega_q - \delta}] \quad (90)$$

$$\geq n \left( \mathbb{E}[G_{\leq \omega_q}] - C_8 \delta \right) \quad (91)$$

$$\geq n \frac{\mathbb{E}[G_{\leq \omega_q}]}{2}, \quad (92)$$

where  $C_8 > 0$  is a constant. In the above derivation, (90) follows since  $x \mapsto \mathbb{E}[G_{\leq x}]$  is monotonically increasing. Inequality (91) is a consequence of Lemma 4, and the final inequality holds for all small enough  $\epsilon$ . Substituting (88) and (92) to (79), we obtain

$$y' \leq yx - \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon \quad (93)$$

$$\leq (\tau - \epsilon)(\omega_q + \delta) - \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon \quad (94)$$

$$\leq 2\tau\omega_q - \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon, \quad (95)$$



where the last inequality holds for all small enough  $\epsilon$ .

**Step-3:** Let  $y = \tau - \epsilon$  as in Step 2. We decompose the first integral in the upper bound in (76) as

$$\mathbb{P}(\Xi_n \leq y) \leq \int_{|x-\omega_q| \leq \delta} \Phi(y'/\sigma) f_{G^{(k)}}(x) dx + \int_{|x-\omega_q| \geq \delta} \underbrace{\Phi(y'/\sigma)}_{\leq 1} f_{G^{(k)}}(x) dx + 8 \int_0^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx. \quad (96)$$

For the first term, we implement the upper limit on  $y'$  specified in (95), which holds a negative value for sufficiently large  $n$ . In this regime, we can apply the upper bound on  $\sigma$  as provided by Corollary 8 to obtain a valid upper limit for the first term in (96). Furthermore, an upper bound for the final term of (96) is obtained in Lemma 9. These estimates yield

$$\mathbb{P}(\Xi_n \leq y) \leq \int_{|x-\omega_q| \leq \delta} \Phi \left[ \frac{1}{\sqrt{C_4 n}} \left( 2\tau\omega_q - \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon \right) \right] f_{G^{(k)}}(x) dx + \int_{|x-\omega_q| \geq \delta} f_{G^{(k)}}(x) dx + 8C_6 n^{-\frac{1}{2}}. \quad (97)$$

In (97), the  $\Phi[\cdot]$ -term decays to zero as  $n \rightarrow \infty$ . Moreover, the second integral also vanishes as  $n \rightarrow \infty$  as a result of Proposition 10. Therefore, we obtain

$$\mathbb{P}(\Xi_n \leq y) \leq o(1) \int_{|x-\omega_q| \leq \delta} f_{G^{(k)}}(x) dx + o(1) \quad (98)$$

$$\leq o(1). \quad (99)$$

This shows that for any  $\epsilon > 0$ , we have  $\mathbb{P}(\Xi_n \leq \tau - \epsilon) \rightarrow 0$ .

**Step-4:** Now, let  $y = \tau + \epsilon$ . We will show for any  $\epsilon > 0$ , we have  $\mathbb{P}(\Xi_n \leq \tau + \epsilon) \rightarrow 1$ .

As in Step-2, for large enough  $n$ , we can argue that  $|h(\xi_q) - \tau| \leq \frac{\epsilon}{4}$ . Combining with (83), for  $y = \tau + \epsilon$ , we obtain

$$h(x) - y \leq h(\xi_q) + \frac{\epsilon}{2} - y \quad (100)$$

$$= h(\xi_q) - \tau - \frac{\epsilon}{2} \quad (101)$$

$$\leq -\frac{\epsilon}{4}. \quad (102)$$

Combining with (92) and substituting to (79), we obtain

$$y' \geq yx - y \underbrace{\mathbb{E}[G_{\leq x}]}_{\leq x} + \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon \quad (103)$$

$$= \frac{\mathbb{E}[G_{\leq \omega_q}]}{8} n\epsilon \quad (104)$$

Analogous to the upper bound in (76), applying Proposition 11 yields the lower estimate

$$\mathbb{P}(\Xi_n \leq y) \geq \int_0^\infty \Phi(y'/\sigma) f_{G^{(k)}}(x) dx - 8 \int_0^\infty \frac{\rho}{\sigma^3} f_{G^{(k)}}(x) dx \quad (105)$$

$$\geq \int_{|x-\omega_q| \leq \delta} \Phi(y'/\sigma) f_{G^{(k)}}(x) dx - o(1) \quad (106)$$

$$\geq \int_{|x-\omega_q| \leq \delta} \Phi\left(\frac{1}{\sqrt{C_4 n}} \frac{\mathbb{E}[G_{\leq \omega_q}] n \epsilon}{8}\right) f_{G^{(k)}}(x) dx - o(1) \quad (107)$$

$$= (1 - o(1)) \int_{|x-\omega_q| \leq \delta} f_{G^{(k)}}(x) dx - o(1) \quad (108)$$

$$= (1 - o(1))(1 - o(1)) - o(1) \quad (109)$$

$$= 1 - o(1). \quad (110)$$

The third inequality follows from (104) and Corollary 8. This shows that for any  $\epsilon > 0$ , we have  $\mathbb{P}(\Xi_n \leq \tau + \epsilon)$ . Combining with the conclusion of Step-3, the proof of the theorem is now complete.

## Appendix B. Proof of Corollary 2

Let  $N_1, \dots, N_n$  be independent and identically distributed  $\mathcal{N}(0, 1)$  random variables. We can set  $\mathbf{W} = [N_1 \cdots N_n]^T / (\sum_{i=1}^n N_i^2)^{1/2}$ . Let  $|N_{i_1}| \leq \dots \leq |N_{i_n}|$ , where  $i_1, \dots, i_n$  is a permutation of  $1, \dots, n$ . By definition, the random vector  $\mathbf{W}_p$  is zero at all indices except at  $\{i_j : j = \lceil nq \rceil + 1, \dots, n\}$  where it equals  $\mathbf{W}$ . This implies  $\mathbf{W}^T \mathbf{W}_p = \sum_{j=\lceil nq \rceil + 1}^n N_{i_j}^2 / \sum_{i=1}^n N_i^2$ . The statement for  $\mathbf{W}^T \mathbf{W}_p$  now follows from Theorem 1 since  $\mathcal{N}(0, 1)^2$  is a Gamma random variable with shape  $\frac{1}{2}$  and scale 2. The convergence of  $\mathbf{W}^T \mathbf{W}_{ee}$  is proved similarly.

## Appendix C. Proof of Theorem 3

Let us first calculate and upper bound on the normalized FLOPs. It is easily seen that  $|\mathbf{W}_{ee}^T \mathbf{X}|^2$  is a Chi-squared random variable with 1 degree of freedom. We thus obtain  $\mathbb{P}(|\mathbf{W}_{ee}^T \mathbf{X}| \geq \tau) = \Gamma(\frac{1}{2}, \frac{\tau^2}{2})$ , where  $\Gamma(\cdot, \cdot)$  is the upper incomplete Gamma function. The normalized FLOPs is thus

$$\bar{\mu}'_c = (1 - q) \Gamma\left(\frac{1}{2}, \frac{\tau^2}{2}\right) + \gamma\left(\frac{1}{2}, \frac{\tau^2}{2}\right), \quad (111)$$

where  $\gamma(\cdot, \cdot)$  is the lower incomplete Gamma function. Further, we obtain

$$\bar{\mu}'_c = (1 - q) + q\gamma\left(\frac{1}{2}, \frac{\tau^2}{2}\right) \quad (112)$$

$$= 1 - q + q\operatorname{erf}\left(\sqrt{\frac{\tau^2}{2}}\right) \quad (113)$$

$$\leq 1 - q + q\left(1 - \underbrace{\sqrt{\frac{2e}{\pi}}}_{>1} \frac{\sqrt{\beta-1}}{\beta} e^{-\beta\tau^2/2}\right) \quad (114)$$

$$\leq 1 - q \frac{\sqrt{\beta-1}}{\beta} e^{-\beta\tau^2/2} \quad (115)$$

$$\leq 1 - qe^{-\tau^2} \quad (116)$$

The upper bound on the error function in (114) follows from [5], and is valid for any  $\beta > 1$ . In (116), we substituted  $\beta = 2$ .

We now analyze the generalization error  $\epsilon_c$ . We consider the following events:

- Let  $E_0$  be the event where the decisions of the conditional perceptron and the teacher do not match so that  $\epsilon_c = P(E_0)$ .
- We let  $E_1$  denote the event that the student  $\mathbf{W}$  and the teacher  $\mathbf{T}$  are at least  $\delta_1$ -close with respect to the angular distance. In other words, let  $E_1$  denote the event that  $\arccos \mathbf{W}^T \mathbf{T} \leq \delta_1$ , where  $\delta_1 \in [0, \frac{\pi}{2}]$ .
- Let  $E_2$  be the event that the student and the early exit vector (which are derived from the student weights) are  $\delta_2$ -close. In other words, let  $E_2$  represent the event  $\arccos \mathbf{W}^\dagger \mathbf{W}_{ee} \leq \delta_2$ , where  $\delta_2 \in [0, \frac{\pi}{2}]$ .
- Finally, let  $E_3$  be the event that  $|\mathbf{X}^T \mathbf{W}_{ee}| \geq \tau$ , encoding the criterion of early exit in the definition of the conditional perceptron in Section 4.2.

We have

$$\epsilon_c = P(E_0 E_1 E_2) + \underbrace{P(E_0 | E_1^c \text{ or } E_2^c)}_{\leq 1} \underbrace{P(E_1^c \text{ or } E_2^c)}_{\leq P(E_1^c) + P(E_2^c)} \quad (117)$$

$$= P(E_0 E_1 E_2 E_3) + \underbrace{P(E_0 E_1 E_2 E_3^c)}_{\leq P(E_0 E_3^c) \leq \epsilon_{uc}} + P(E_1^c) + P(E_2^c) \quad (118)$$

$$= \epsilon_{uc} + P(E_1^c) + P(E_2^c) + P(E_0 E_1 E_2 E_3) \quad (119)$$

Let  $\bar{\epsilon}_{uc}$  denote the asymptotic  $n, N_t \rightarrow \infty$  generalization error for the unconditional perceptron so that  $\epsilon_{uc} = \frac{1}{\pi} \mathbb{E}[\arccos \mathbf{W}^T \mathbf{T}] \rightarrow \bar{\epsilon}_{uc}$  as  $n \rightarrow \infty$ . Due to the self-averaging property learning [4], we have, in addition, the convergence in mean  $\frac{1}{\pi} \arccos \mathbf{W}^T \mathbf{T} \rightarrow \bar{\epsilon}_{uc}$ . Hence, if  $\delta_1 > \pi \bar{\epsilon}_{uc}$ , we have  $P(E_1^c) \rightarrow 0$ . With a similar argument, provided that  $\delta_2 > \arccos \sqrt{\tau} q$ , we have  $P(E_2^c) \rightarrow 0$  as a result of Corollary 2. What is left to analyze is the final term. By symmetry, we have

$$P(E_0 E_1 E_2 E_3) = 2P(\mathbf{X}^T \mathbf{T} < 0, \arccos \mathbf{W}^T \mathbf{T} \leq \delta_1, \arccos \mathbf{W}^T \mathbf{W}_{ee} \leq \delta_2, \mathbf{X}^T \mathbf{W}_{ee} \geq \tau). \quad (120)$$

Let us recall the triangle inequality for angular distances: For unit norm vectors  $a_1, a_2, a_3$ , we have  $\arccos a_1^\dagger a_2 \leq \arccos a_1^\dagger a_3 + \arccos a_3^\dagger a_2$ . Therefore,

$$P(E_0 E_1 E_2 E_3) \leq 2P(\mathbf{X}^T \mathbf{T} < 0, \arccos \mathbf{T}^T \mathbf{W}_{ee} \leq \delta_1 + \delta_2, \mathbf{X}^T \mathbf{W}_{ee} \geq \tau) \quad (121)$$

$$= 2P(\mathbf{X}^T \mathbf{T} < 0, \mathbf{T}^T \mathbf{W}_{ee} \geq \cos(\delta_1 + \delta_2), \mathbf{X}^T \mathbf{W}_{ee} \geq \tau) \quad (122)$$

$$\leq 2P(\mathbf{X}^T \mathbf{a} < 0, \mathbf{X}^T \mathbf{b} \geq \tau) \quad (123)$$

where  $\mathbf{a}$  and  $\mathbf{b}$  are arbitrary unit-norm deterministic vectors with  $\mathbf{a}^T \mathbf{b} = \cos(\delta_1 + \delta_2)$ . The random variables  $\mathbf{X}^T \mathbf{a}$  and  $\mathbf{X}^T \mathbf{b}$  are jointly Gaussian with zero mean, unit variance, and covariance  $\rho \triangleq \cos(\delta_1 + \delta_2)$ . We can thus evaluate the joint probability as

$$P(\mathbf{X}^T \mathbf{a} < 0, \mathbf{X}^T \mathbf{b} \geq \tau) = \int_{\tau}^{\infty} \int_{-\infty}^0 \frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)}\right) dx dy \quad (124)$$

$$= \frac{1}{2\sqrt{2\pi}} \int_{\tau}^{\infty} e^{-y^2/2} \operatorname{erfc}\left(\frac{\rho y}{\sqrt{2(1-\rho^2)}}\right) dy \quad (125)$$

Using the upper bound  $\operatorname{erfc}x \leq e^{-x^2}$ , we obtain

$$P(\mathbf{X}^T \mathbf{a} < 0, \mathbf{X}^T \mathbf{b} \geq \tau) \leq \frac{1}{2\sqrt{2\pi}} \int_{\tau}^{\infty} e^{-y^2/2} e^{-\frac{\rho^2 y^2}{2(1-\rho^2)}} dy \quad (126)$$

$$= \frac{1}{2\sqrt{2\pi}} \int_{\tau}^{\infty} e^{-\frac{y^2}{2(1-\rho^2)}} dy \quad (127)$$

$$= \frac{1}{4} \sqrt{1-\rho^2} \operatorname{erfc}\left(\frac{\tau}{\sqrt{2(1-\rho^2)}}\right) \quad (128)$$

$$\leq \frac{1}{2} e^{-\frac{\tau^2}{2(1-\rho^2)}}. \quad (129)$$

A joint consideration with the previous bounds concludes the proof of the theorem.