Brain network science modelling of sparse neural networks enables Transformers and LLMs to perform as fully connected

Editors: List of editors' names

Abstract

The Cannistraci-Hebb training (CHT) is a brain-inspired method that is used in Dynamic Sparse Training (DST) for growing synaptic connectivity in sparse neural networks. CHT leverages a gradient-free, topology-driven link regrowth mechanism, which has been shown to achieve ultra-sparse (1% connectivity or lower) advantage across various tasks compared to fully connected networks (FCs). Yet, CHT suffers two main drawbacks: (i) its time complexity is $\mathcal{O}(N \cdot d^3)$. N node network size, d node degree - hence it can be efficiently applied only to ultra-sparse networks. (ii) it rigidly selects top link prediction scores, which is inappropriate for the early training epochs, when the network topology presents many unreliable connections. Here, we design the first brain-inspired network model - termed bipartite receptive field (BRF) - to initialize the connectivity of sparse artificial neural networks. Then, we propose a matrix multiplication GPU-friendly approximation of the CH link predictor, which reduces the computational complexity to $\mathcal{O}(N^3)$, enabling a fast implementation of link prediction in large-scale models. Moreover, we introduce the Cannistraci-Hebb training soft rule (CHTs), which adopts a flexible strategy for sampling connections in both link removal and regrowth, balancing the exploration and exploitation of network topology. We also propose a sigmoid-based gradual density decay strategy, leading to an advanced framework referred to as CHTss. Empirical results show that using 1% of connections, CHTs outperform FCs in MLP architectures on visual classification tasks and compress some networks to less than 30% of the nodes. Using only 5% of the connections, CHTss outperforms FCs in two Transformer-based machine translation tasks. Finally, using 30% of the connections, CHTs and CHTss achieve superior performance compared to other dynamic sparse training methods in language modeling across different sparsity levels on LLaMA 60M, 130M, and 1B, and CHTs outperforms FC on the LLaMA1B model.

Keywords: Dynamic sparse training, Network science, Efficient training

1. Introduction

Fully connected layers in large models pose computational challenges during training and deployment. In contrast, the brain's neural networks exhibit sparse connectivity Drachman (2005); Walsh (2013), suggesting more scalable architectures. Dynamic sparse training (DST) Mocanu et al. (2018); Jayakumar et al. (2020); Evci et al. (2020); Yuan et al. (2021); Zhang et al. (2024b) reduces computational and memory costs while maintaining performance. Unlike pruning methods Han et al. (2016); Frantar and Alistarh (2023); Zhang et al. (2024a), DST starts with sparse networks and evolves their topology during training. Key innovations of DST focus on regrowth criteria, such as the gradient-free Cannistraci-Hebb training (CHT) Zhang et al. (2024b), inspired by brain-inspired network science Cannistraci et al. (2013); Daminelli et al. (2015); Durán et al. (2017); Cannistraci (2018); Narula (2017). CHT excels in ultra-sparse ANNs but faces challenges such as stacking in epitopological local minima and high time complexity of link prediction, making it impractical for large-scale models.

This article introduces the <u>Cannistraci-Hebb Training soft rule</u> (CHTs), which addresses CHT's limitations. CHTs 1) uses a multinomial distribution for both link removal and regrowth that balances the exploration and exploitation of network topology, 2) reduces the

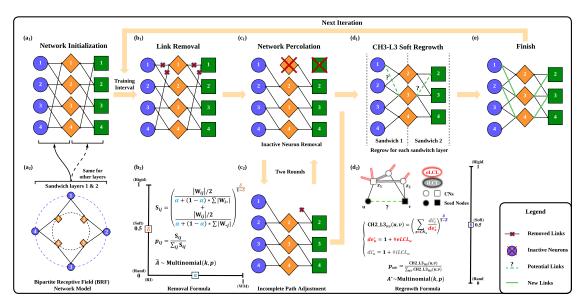


Figure 1: Illustration of the CHTs process. One training iteration follows the steps of $(a1) \rightarrow (b1) \rightarrow (c1) \rightarrow (c2) \rightarrow (d1) \rightarrow (e)$. (a1) Network initialization with each of the sandwich layers being a bipartite receptive field (BRF) network. (a2) BRF network representation with r=0. (b1) Link removal process. (b2) Formula for link removal. (c1) Inactive neurons removal. (c2) Adjust and remove incomplete links caused by inactive neuron removal. (d1) Regrowth of links according to the CH2-L3 node-based soft rule. (d2) Detailed illustration of the CH2-L3 node-based soft rule. (e) Finished state of the network after one iteration. The next iteration repeats the steps (b1) - (e) from this finished state. \tilde{A} indicates the removal set of the iteration and A^* is the regrown set.

time complexity of the path-based link predictor to $\mathcal{O}(N^3)$ with a node-based solution, and **3)** initializes the sparse topologies for bipartite networks with the brain-like receptive field, demonstrating superior performance compared to the traditional methods. Additionally, we propose a sigmoid gradual density decay strategy, forming an enhanced framework termed CHTss.

From the experimental results, CHTs and CHTss both outperform fully connected Transformers with only 5% of the links on Machine Translation tasks of Multi30k and IWSLT and achieves performance comparable to the fully connected LLaMA-60M, 130M, 1B in language modeling tasks on OpenWebText. These findings underscore the potential of CHTs and CHTss in enabling highly efficient and effective large-scale sparse neural network training.

2. Cannistraci-Hebb training soft rule

Definition 1. Epitopological local minima. In the context of dynamic sparse training methods, we define an epitopological local minima (ELM) as a state where the sets of removed links and regrown links exhibit a significant overlap. See Appendix F for detailed descriptions.

Cannistraci-Hebb soft removal and regrowth. In this article, we adopt a probabilistic approach where the process of both regrowth and removal can be viewed as sampling from a $\{0,1\}$ multinomial distribution, with the score assigned by either removal metrics or link prediction scores, introducing a "soft sampling" mechanism. In this setup, each mask value is not rigidly determined by the scores but allows for selecting (with lower probability) low-score links as the target links to remove or regrow, facilitating the escape from the epitopological local minima (ELM).

Link removal alternating from weight magnitude and relative importance. We illustrate the link removal part of CHTs in Figure 1b1) and b2). We employ two methods, Weight Magnitude (WM) and Relative Importance (RI) Zhang et al. (2024a), to remove the connections during dynamic sparse training. Detailed information can be found in Appendix G.

Node-based link regrowth. In the original CHT framework, the time complexity of the path-based CH3-L3 (**CH3-L3p**, see Appendix C) metric is $O(N \cdot d^3)$, where N is the number of nodes and d is the network's average degree. This complexity is prohibitive for large models with numerous nodes and higher-density layers. To address this issue, we introduce a more efficient, node-based paradigm that eliminates the reliance on length-three paths between seed nodes, which also incorporates internal local community links (iLCL) to enhance the expressiveness of the formula. This variant, referred to as **CH2-L3n**, is formulated as:

$$\mathbf{CH2\text{-}L3n}(u,v) = \sum_{z \in L3} \frac{di_z^*}{de_z^*} \tag{1}$$

Here, u and v denote the seed nodes, while z_1 and z_2 are common neighbors on the L3 path Muscoloni et al. (2022). The term de_z^* and di_z^* represents the number of external local community links (eLCL) and iLCL of node z, with a default increment of 1 to prevent division by zero. We evaluate the runtime performance of CH3-L3p and CH2-L3n across different network sizes and sparsity levels, as illustrated in Figure 12. The results reveal that the node-based version achieves significantly faster runtime performance compared to the path-based methods.

Bipartite receptive field network modeling. We propose the Bipartite Receptive Field (BRF) network, a novel sparse topological initialization method for generating brain-like receptive field connectivity. During Bipartite Small World (BSW) initialization Zhang et al. (2024b), each output node is connected to its spatially nearest input nodes. This spatially local connectivity pattern aligns with the concept of receptive fields observed in biological neural systems, where neurons respond selectively to localized regions of input space. However, the rewiring process of BSW does not follow brain mechanisms: it simply deletes a set of links from the closer input nodes to rewire them uniformly at random anywhere on the input layer. In contrast, BRF directly generates connectivity with a tunable spatial-dependent randomness parameter $r \in [0, 1]$, controlling the clustering of links around the adjacency matrix diagonal $(r = 0 \rightarrow \text{fully local}, r = 1 \rightarrow \text{ER})$. Furthermore, the BRF model has the important property to conserve the degree distribution of the output neurons for each layer. In this study we consider fixed or uniformly at random, as shown in Appendix D.

The comparison of BSW, BSF, BRF-fixed, and BRF-uniform initializations are shown in Appendix D.

3. Sigmoid Gradual Decrease Density

As demonstrated in GraNet Liu et al. (2021) and MEST $_{EM\&S}$ Yuan et al. (2021), incorporating a density decrease strategy can significantly improve the performance of dynamic sparse training. In this article, we propose a sigmoid-based gradual density decrease strategy, defined as:

$$s_t = s_i + (s_i - s_f) \left(\frac{1}{1 + e^{-k\left(t - \frac{t_f + t_0}{2}\right)}} \right),$$
 (2)

where $t \in \{t_0, t_0 + \Delta t, \dots, t_0 + n\Delta t\}$, s_i is the initial sparsity, s_f is the target sparsity, t_0 is the starting epoch of gradual pruning, t_f is the end epoch of gradual pruning, and Δt is the pruning frequency. k controls the curvature of the decrease. We set k=6 for all the experiments in this article. This strategy ensures a smoother initial pruning phase, allowing the model to warm up and stabilize before undergoing significant pruning, thereby enhancing training stability and performance. A detailed discussion of the decay strategy can be found in Appendix I.

4. Experiments

Experimental details are provided in Appendix H. The baseline methods are detailed in Appendix K. We also demonstrate superiority with CHTs and CHTss using MLP on image classification datasets (See Appendix N).

4.1. Transformer on Machine Translation

We assess CHTs and CHTss using Transformer on classic machine translation tasks across three datasets. We report the BLEU in Table 5, which demonstrates that 1) CHTs surpasses other fixed density DST methods on all the sparsity scenrios. 2) Incorporating with the sigmoid density decrease, CHTss outperforms even the fully connected ones with only 5% density.

4.2. Natural Language Generation

Language modeling. We utilize LLaMA-60M, 130M, 1B (Touvron et al., 2023a) architecture as the baseline for our language generation experiments. We show the validation perplexity results of each algorithm across the different sparsities in Table 6. As shown, CHTs stably outperforms SET and RigL while CHTss is constantly better than GraNet and GMP. At 70% sparsity, CHTs and CHTss achieve superior performance compared to other dynamic sparse training methods in language modeling across different sparsity levels on LLaMA 60M, 130M, and 1B, and CHTs outperforms FC on the LLaMA1B model.

References

- Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *science*, 286(5439):509–512, 1999.
- Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. arXiv preprint arXiv:1711.05136, 2017.
- Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, et al. Findings of the 2017 conference on machine translation (wmt17). Association for Computational Linguistics, 2017.
- Carlo Vittorio Cannistraci. Modelling self-organization in complex networks via a brain-inspired network automata theory improves link reliability in protein interactomes. *Sci Rep*, 8(1):2045–2322, 10 2018. doi: 10.1038/s41598-018-33576-8.
- Carlo Vittorio Cannistraci, Gregorio Alanis-Lobato, and Timothy Ravasi. From link-prediction in brain connectomes and protein interactomes to the local-community-paradigm in complex networks. *Scientific reports*, 3(1):1613, 2013.
- Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa Bentivogli, and Marcello Federico. Report on the 11th IWSLT evaluation campaign. In Marcello Federico, Sebastian Stüker, and François Yvon, editors, *Proceedings of the 11th International Workshop on Spoken Language Translation: Evaluation Campaign*, pages 2–17, Lake Tahoe, California, December 4-5 2014. URL https://aclanthology.org/2014.iwslt-evaluation.1.
- Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. Emnist: Extending mnist to handwritten letters. In 2017 international joint conference on neural networks (IJCNN), pages 2921–2926. IEEE, 2017.
- Simone Daminelli, Josephine Maria Thomas, Claudio Durán, and Carlo Vittorio Cannistraci. Common neighbours and the local-community-paradigm for topological link prediction in bipartite networks. New Journal of Physics, 17(11):113037, nov 2015. doi: 10.1088/1367-2630/17/11/113037. URL https://doi.org/10.1088/1367-2630/17/11/113037.
- Tri Dao, Beidi Chen, Nimit S Sohoni, Arjun Desai, Michael Poli, Jessica Grogan, Alexander Liu, Aniruddh Rao, Atri Rudra, and Christopher Ré. Monarch: Expressive structured matrices for efficient and accurate training. In *International Conference on Machine Learning*, pages 4690–4721. PMLR, 2022.
- David A Drachman. Do we have brain to spare?, 2005.
- Claudio Durán, Simone Daminelli, Josephine M Thomas, V Joachim Haupt, Michael Schroeder, and Carlo Vittorio Cannistraci. Pioneering topological methods for network-based drug-target prediction by exploiting a brain-network self-organization theory. *Briefings in Bioinformatics*, 19(6):1183–1202, 04 2017. doi: 10.1093/bib/bbx041. URL https://doi.org/10.1093/bib/bbx041.

- Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. Multi30k: Multilingual english-german image descriptions. In *Proceedings of the 5th Workshop on Vision and Language*, pages 70–74. Association for Computational Linguistics, 2016. doi: 10.18653/v1/W16-3210. URL http://www.aclweb.org/anthology/W16-3210.
- P ERDdS and A R&wi. On random graphs i. Publ. math. debrecen, 6(290-297):18, 1959.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 2943–2952. PMLR, 2020. URL http://proceedings.mlr.press/v119/evci20a.html.
- Elias Frantar and Dan Alistarh. Massive language models can be accurately pruned in one-shot. arXiv preprint arXiv:2301.00774, 2023.
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/ OpenWebTextCorpus, 2019.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. Advances in neural information processing systems, 28, 2015.
- Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In Yoshua Bengio and Yann LeCun, editors, 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016. URL http://arxiv.org/abs/1510.00149.
- Donald Hebb. The organization of behavior. emphnew york, 1949.
- Siddhant Jayakumar, Razvan Pascanu, Jack Rae, Simon Osindero, and Erich Elsen. Top-kast: Top-k always sparse training. *Advances in Neural Information Processing Systems*, 33: 20744–20754, 2020.
- Mike Lasby, Anna Golubeva, Utku Evci, Mihai Nica, and Yani Ioannou. Dynamic sparse training with structured sparsity. arXiv preprint arXiv:2305.02299, 2023.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proc. IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791. URL https://doi.org/10.1109/5.726791.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: single-shot network pruning based on connection sensitivity. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=B1VZqjAcYX.
- Ming Li, Run-Ran Liu, Linyuan Lü, Mao-Bin Hu, Shuqi Xu, and Yi-Cheng Zhang. Percolation on complex networks: Theory and application. *Physics Reports*, 907:1–68, 2021.

- Shiwei Liu, Tianlong Chen, Xiaohan Chen, Zahra Atashgahi, Lu Yin, Huanyu Kou, Li Shen, Mykola Pechenizkiy, Zhangyang Wang, and Decebal Constantin Mocanu. Sparse training via boosting pruning plasticity with neuroregeneration. *Advances in Neural Information Processing Systems*, 34:9908–9922, 2021.
- Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):1–12, 2018.
- Alessandro Muscoloni, Umberto Michieli, Yingtao Zhang, and Carlo Vittorio Cannistraci. Adaptive network automata modelling of complex networks. *Preprints*, May 2022. doi: 10.20944/preprints202012.0808.v3. URL https://doi.org/10.20944/preprints202012.0808.v3.
- Vaibhav et al Narula. Can local-community-paradigm and epitopological learning enhance our understanding of how local brain connectivity is able to process, learn and memorize chronic pain? *Applied network science*, 2(1), 2017. doi: 10.1007/s41109-017-0048-x.
- Aleksandra I. Nowak, Bram Grooten, Decebal Constantin Mocanu, and Jacek Tabor. Fantastic weights and how to find them: Where to prune in dynamic sparse training, 2023. URL https://arxiv.org/abs/2306.12230.
- Ameya Prabhu, Girish Varma, and Anoop Namboodiri. Deep expander networks: Efficient deep networks from graph theory. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 20–35, 2018.
- James Stewart, Umberto Michieli, and Mete Ozay. Data-free model pruning at initialization via expanders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4518–4523, 2023.
- 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.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023b.
- Christopher A Walsh. Peter huttenlocher (1931–2013). Nature, 502(7470):172–172, 2013.
- Duncan J Watts and Steven H Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393(6684):440–442, 1998.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:1708.07747, 2017.

- Geng Yuan, Xiaolong Ma, Wei Niu, Zhengang Li, Zhenglun Kong, Ning Liu, Yifan Gong, Zheng Zhan, Chaoyang He, Qing Jin, et al. Mest: Accurate and fast memory-economic sparse training framework on the edge. *Advances in Neural Information Processing Systems*, 34:20838–20850, 2021.
- Yingtao Zhang, Haoli Bai, Haokun Lin, Jialin Zhao, Lu Hou, and Carlo Vittorio Cannistraci. Plug-and-play: An efficient post-training pruning method for large language models. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=Tr0lPx9woF.
- Yingtao Zhang, Jialin Zhao, Wenjing Wu, Alessandro Muscoloni, and Carlo Vittorio Cannistraci. Epitopological learning and cannistraci-hebb network shape intelligence brain-inspired theory for ultra-sparse advantage in deep learning. In *The Twelfth International Conference on Learning Representations*, 2024b. URL https://openreview.net/forum?id=iayEcORsGd.
- Michael Zhu and Suyog Gupta. To prune, or not to prune: exploring the efficacy of pruning for model compression, 2017. URL https://arxiv.org/abs/1710.01878.

Table 1: Hyperparameters of MLP on Image Classification Tasks.

Hyper-parameter	MLP
Hidden Dimension	1568 (3072 for CIFAR10)
# Hidden layers	3
Batch Size	32
Training Epochs	100
LR Decay Method	Linear
Start Learning Rate	0.025
End Learning Rate	$2.5e^{-4}$
ζ (fraction of removal)	0.3
Update Interval (for DST)	1 epoch
Momentum	0.9
Weight decay	$5e^{-4}$

Table 2: Hyperparameters of Transformer on Machine Translation Tasks. inoam refers to a learning rate scheduler that incorporates iterative warm-up phases, specifically designed for dynamic sparse training (DST) methods. The purpose is to allow newly regrown connections to accumulate momentum, preventing potential harm to the training process. For the fully connected (FC) baseline, only the standard noam scheduler is used.

Hyper-parameter	Multi30k	IWSLT14	WMT17
Embedding Dimension	512	512	512
Feed-forward Dimension	1024	2048	2048
Batch Size	1024 tokens	10240 tokens	12000 tokens
Training Steps	5000	20000	80000
Dropout	0.1	0.1	0.1
Attention Dropout	0.1	0.1	0.1
Max Gradient Norm	0	0	0
Warmup Steps	1000	6000	8000
Learning rate Decay Method	inoam	inoam	inoam
Iterative warmup steps	20	20	20
Label Smoothing	0.1	0.1	0.1
Layer Number	6	6	6
Head Number	8	8	8
Learning Rate	0.25	2	2
ζ (fraction of removal)	0.3	0.3	0.3
Update Interval (for DST)	100 steps	100 steps	100 steps

Appendix A. Conclusion

In this article, we propose the <u>Cannistraci-Hebb Training soft rule</u> with <u>sigmoid gradual</u> density decay (CHTss). First, we introduce a matrix multiplication mathematical formula

Table 3: Hyperparameters of LLaMA-60M, LLaMA-130M, and LLaMA-1B on OpenWebText. inoam refers to a learning rate scheduler that incorporates iterative warm-up phases, specifically designed for dynamic sparse training (DST) methods. The purpose is to allow newly regrown connections to accumulate momentum, preventing potential harm to the training process. For the fully connected (FC) baseline, only the standard noam scheduler is used.

Hyper-parameter	LLaMA-60M	LLaMA-130M	LLaMA-1B
Embedding Dimension	512	768	2048
Feed-forward Dimension	1376	2048	5461
Global Batch Size	512	512	512
Sequence Length	256	256	256
Training Steps	10000	30000	100000
Learning Rate	3e-3 (1e-3 for FC)	3e-3 (1e-3 for FC)	3e-3 (4e-4 for FC)
Warmup Steps	1000	10000	10000
Learning rate Decay Method	inoam	inoam	inoam
Iterative warmup steps	20	20	20
Optimizer	Adam	Adam	Adam
Layer Number	8	12	24
Head Number	8	12	32
ζ (fraction of removal)	0.1	0.1	0.1
Update Interval (for DST)	100 steps	100 steps	100 steps

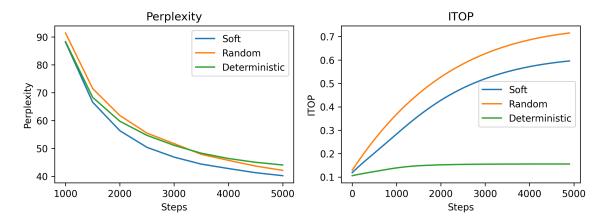


Figure 2: Comparison of link regrowth strategies in CHTs using a LLaMA-60M model trained on OpenWebText for 5000 steps. The left plot shows validation perplexity (lower is better), while the right plot reports the in-time over-parameterization (ITOP) rate, which measures the cumulative proportion of links activated during training. Results are presented for three strategies: Soft, Random, and Deterministic regrowth.

for GPU-friendly approximation of the CH link predictor. This significantly reduces the computational complexity of CHT and speeds up the running time, enabling the implemen-

Table 4: Performance comparison of different dynamic sparse training methods on MNIST, Fashion MNIST (FMNIST), and EMNIST datasets trained on MLP at 99% sparsity. ACC represents accuracy, and ANP denotes the active neuron percolation rate, indicating the final size of the network. Accuracies present a standard error taken over three seeds. The best dynamic sparse training method for each dataset is highlighted in bold, and the performances that surpass the fully connected model are marked with "*".

Method	MNIST		FMNIST		EMNIST	
Modfod	ACC (%)	ANP	ACC (%)	ANP	ACC (%)	ANP
FC	98.78 ± 0.02	-	90.88 ± 0.02	_	87.13 ± 0.04	-
CHTs ^p CHTs ⁿ CHT RigL SET	$98.81 \pm 0.04*$ 98.76 ± 0.05 98.48 ± 0.04 98.61 ± 0.01 98.14 ± 0.02	20% 27% 29% 29% 100%	$90.93 \pm 0.03*$ 90.67 ± 0.05 88.70 ± 0.07 89.91 ± 0.07 89.00 ± 0.09	89% 73% 30% 55% 100%	$87.61 \pm 0.07^*$ $87.82 \pm 0.04^*$ 86.35 ± 0.08 86.94 ± 0.08 86.31 ± 0.08	24% 28% 21% 28% 100%
CHTss ⁿ GraNet GMP	$98.83 \pm 0.02*$ $98.81 \pm 0.00*$ 98.62 ± 0.03	32% 35% 58 %	90.81 ± 0.11 89.98 ± 0.06 90.29 ± 0.19	40% 53% 69%	$87.52 \pm 0.04*$ 86.94 ± 0.03 86.93 ± 0.09	35 % 45% 75 %

^p Refers to the regrowth method CH3 L3p.

Table 5: Performance comparison on machine translation tasks of Multi30k, IWSLT, and WMT with varying final sparsity levels. The scores indicate BLEU scores, which is the higher the better. CHTs (GMP) represents CHTs with GMP's density decay strategy. Bold values denote the best performance among fixed sparsity DST methods or density decay DST methods. The performances that surpass the fully connected model are marked with "*". The density decay of GMP, GraNet, and CHTss starts with a sparsity of 50%. The scores are averaged over three seeds \pm their standard error.

Method	Multi30k		IWS	LT	WMT		
mounoa	95%	90%	95%	90%	95%	90%	
FC	31.38 ± 0.38		24.48	24.48 ± 0.30		25.49 ± 0.15	
SET	28.99 ± 0.28	29.73 ± 0.10	18.53 ± 0.05	20.13 ± 0.08	20.19 ± 0.12	21.52 ± 0.28	
RigL	29.94 ± 0.27	30.26 ± 0.34	20.53 ± 0.21	21.52 ± 0.15	20.71 ± 0.21	22.22 ± 0.10	
CHT	27.79	28.38	18.59	19.91	19.03	21.08	
CHTs	28.94 ± 0.57	29.81 ± 0.37	21.15 ± 0.10	21.92 ± 0.17	20.94 ± 0.63	22.40 ± 0.06	
MEST	28.89 ± 0.26	30.04 ± 0.52	19.56 ± 0.10	21.05 ± 0.21	20.70 ± 0.07	22.22 ± 0.10	
GMP	30.51 ± 0.82	30.49 ± 0.40	22.76 ± 0.82	22.82 ± 0.53	22.47 ± 0.10	23.37 ± 0.08	
GraNet	31.31 ± 0.31	$31.62 \pm 0.48*$	22.53 ± 0.12	22.43 ± 0.09	22.51 ± 0.21	23.46 ± 0.09	
CHTss :	$32.03 \pm 0.29*$	$32.86 \pm 0.16*$	$24.51 \pm 0.02*$	24.31 ± 0.04	23.73 ± 0.43	24.61 ± 0.14	

tation of CHTs in large-scale models. Second, we propose a Cannistraci-Hebb training soft

ⁿ Refers to the regrowth method CH2 L3n.

Table 6: Validation perplexity of different dynamic sparse training (DST) methods on OpenWebText using LLaMA-60M, LLaMA-130M, and LLaMA-1B across varying sparsity levels. Bold values denote the best performance among fixed sparsity DST methods or density decay DST methods. Lower perplexity corresponds to better model performance. GMP, GraNet, and CHTss are run with an initial sparsity of $s_i = 0.5$. The test of CHT over LLaMA-1B is missing due to its excessive runtime. The performances that surpass the fully connected model are marked with "*".

Method		LLaM	A-60M			LLaMA	A-130M		LLaMA-1B
111001104	70%	80%	90%	95%	70%	80%	90%	95%	70%
FC		26	.56			19	.27		14.62
SET RigL CHT	31.77 39.96 31.02	30.69 41.33 32.99	35.26 45.34 35.01	39.70 51.49 41.87	20.82 25.85 21.02	22.02 66.35 22.82	24.73 37.18 26.27	28.37 49.39 30.01	16.37 149.17 –
CHTs	28.12	29.84	33.03	36.47	20.10	21.33	23.71	26.45	14.53*
MEST GMP GraNet CHTss	28.26 29.22 30.55 27.62	29.94 30.59 31.51 29.00	33.60 33.68 33.76 31.42	37.87 39.00 39.98 35.10	21.32 20.49 22.84 19.85	22.21 22.28 29.03 20.70	24.98 23.61 26.81 22.51	27.96 27.16 61.31 25.07	60.36 31.76 79.44 15.41

Table 7: Perplexity (PPL) results across different sparsities (0.7, 0.8, 0.9, 0.95) for CHTs and CHTss under different regrowth strategies (Fixed and Uniform) and r settings on LLaMA60M.

			Fixed			Uniform			
	Sparsity	r = 0.0	r = 0.1	r = 0.2	r = 0.3	r = 0.0	r = 0.1	r = 0.2	r = 0.3
	70%	28.16	28.39	28.25	28.32	30.11	28.12	28.43	28.56
m CHTs	80%	30.22	29.84	30.04	30.03	30.19	29.86	30.23	30.06
CHIS	90%	33.32	33.37	33.03	33.77	33.45	33.36	33.88	33.72
	95%	37.29	37.51	37.24	37.46	37.23	36.47	37.33	37.67
	70%	27.62	30.05	27.82	28.43	27.62	27.74	27.74	27.68
CHTss	80%	29.00	29.00	29.66	32.91	29.49	29.69	29.09	29.24
CHISS	90%	31.51	31.67	31.65	31.59	31.66	32.61	31.68	31.42
	95%	38.66	35.31	36.24	37.50	42.20	37.40	35.36	35.10

rule (CHTs), which innovatively utilizes a soft sampling rule for both removal and regrowth links, striking a balance for epitopological exploration and exploitation. Third, we integrate CHTs with a sigmoid gradual density decay strategy. Empirically, CHTss surpasses the fully connected Transformer using only 5% density and achieves comparable language modeling performance. This represents a relevant result for dynamic sparse training.

Table 8: Perplexity (PPL) results across different sparsities (0.7, 0.8, 0.9, 0.95) for CHTs and CHTss under different regrowth strategies (Fixed and Uniform) and r settings on LLaMA-130M.

			Fixed			Uniform			
	Sparsity	r = 0.0	r = 0.1	r = 0.2	r = 0.3	r = 0.0	r = 0.1	r = 0.2	r = 0.3
	70%	20.24	20.16	20.10	20.25	20.62	20.18	20.15	20.20
m CHTs	80%	21.33	21.37	21.36	21.48	21.34	21.40	21.40	22.49
CHIS	90%	23.72	23.76	23.76	23.94	23.74	23.73	23.71	24.99
	95%	28.05	26.45	26.90	26.91	26.78	27.97	29.05	27.10
	70%	20.63	19.88	19.93	19.85	21.43	19.90	20.93	19.94
CHTss	80%	20.71	22.60	20.86	20.70	20.73	20.74	20.72	20.82
CHISS	90%	22.58	22.72	22.61	22.51	22.53	22.59	22.60	23.12
	95%	25.28	25.12	25.20	25.12	25.07	25.15	25.23	25.12

Appendix B. Related Work

B.1. Dynamic sparse training

Dynamic sparse training is a subset of sparse training methodologies. Unlike static sparse training methods (also known as pruning at initialization) Prabhu et al. (2018); Lee et al. (2019); Dao et al. (2022); Stewart et al. (2023), dynamic sparse training allows for the evolution of network topology during the training process. The pioneering method in this field is Sparse Evolutionary Training (SET) Mocanu et al. (2018), which removes links based on the magnitude of their weights and regrows new links randomly. Subsequent developments have sought to refine and expand upon this concept of dynamic topological evolution. One such advancement was proposed by DeepR Bellec et al. (2017), a method that adjusts network connections based on stochastic gradient updates combined with a Bayesian-inspired update rule. Another significant contribution is RigL Evci et al. (2020), which leverages the gradient information of non-existing links to guide the regrowth of new connections during training. MEST Yuan et al. (2021) utilizes both gradient and weight magnitude information to selectively remove and randomly regrow new links, analogously to SET. In addition, it introduces an EM&S strategy that allows the model to train at a higher density and gradually converge to the target sparsity. The Top-KAST Jayakumar et al. (2020) method maintains constant sparsity throughout training by selecting the top K parameters based on parameter magnitude at each training step and applying gradients to a broader subset B, where $B \supset A$. To avoid settling on a suboptimal sparse subset, Top-KAST also introduces an auxiliary exploration loss that encourages ongoing adaptation of the mask. Additionally, sRigL Lasby et al. (2023) adapts the principles of RigL to semi-structured sparsity, facilitating the training of vision models from scratch with actual speed-ups during training phases. Despite these advancements, the state-of-the-art method remains RigL-based, yet it is not fully sparse in backpropagation, necessitating the computation of gradients for non-existing links. Addressing this limitation, Zhang et al. Zhang et al. (2024b) propose CHT, a dynamic sparse training methodology that adopts a gradient-free regrowth strategy that relies solely on

topological information (network shape intelligence), achieving an ultra-sparse configuration that surpasses fully connected networks in some tasks.

B.2. Cannistraci-Hebb Theory and Network Shape Intelligence

As the SOTA gradient-free link regrown method, CHT Zhang et al. (2024b) originates from a brain-inspired network science theory. Drawn from neurobiology, Hebbian learning was introduced in 1949 (Hebb, 1949) and can be summarized in the axiom: "neurons that fire together wire together." This could be interpreted in two ways: changing the synaptic weights (weight plasticity) and changing the shape of synaptic connectivity (Cannistraci et al., 2013; Daminelli et al., 2015; Durán et al., 2017; Cannistraci, 2018; Narula, 2017). The latter is also called *epitopological plasticity* (Cannistraci et al., 2013) because plasticity means "to change shape," and epitopological means "via a new topology." Epitopological Learning (EL) (Daminelli et al., 2015; Durán et al., 2017; Cannistraci, 2018) is derived from this second interpretation of Hebbian learning and studies how to implement learning on networks by changing the shape of their connectivity structure. One way to implement EL is via link prediction, which predicts the existence and likelihood of each nonobserved link in a network. CH3-L3 is one of the best-performing and most robust network automata, belonging to the Cannistraci-Hebb (CH) theory (Muscoloni et al., 2022), which can automatically evolve the network topology starting from a given structure. The rationale is that, in any complex network with local-community organization, the cohort of nodes tends to be co-activated (fire together) and to learn by forming new connections between them (wire together) because they are topologically isolated in the same local community (Muscoloni et al., 2022). This minimization of the external links induces a topological isolation of the local community, which is equivalent to forming a barrier around it. The external barrier is fundamental to maintaining and reinforcing the signaling in the local community, inducing the formation of new links that participate in epitopological learning and plasticity.

Appendix C. Cannistraci-Hebb epitopological rationale

The original CHT framework leverages the Cannistraci-Hebb link predictor on Length 3 paths (CH3-L3p) metric for link regrowth. Given two seed nodes u and v in a network, this metric assigns a score

$$\mathbf{CH3-L3p}(u,v) = \sum_{z_1, z_2 \in L3} \frac{1}{\sqrt{de_{z_1}^* \cdot de_{z_2}^*}}$$
(3)

Here, u and v denote the seed nodes, while z_1 and z_2 are common neighbors on the L3 path Muscoloni et al. (2022), a walk of three consecutive links that connects u to v via those two intermediate nodes. The term de_i^* represents the number of external local community links (eLCL) of node i, with a default increment of 1 to prevent division by zero. Path-based link prediction has demonstrated its effectiveness on both real-world networks Muscoloni et al. (2022) and artificial neural networks Zhang et al. (2024b). However, this method incurs a high computational cost due to the need to compute and store all length-three paths, resulting in a time complexity of $O(N \cdot d^3)$, where N is the number of nodes and d is the network's average degree. This complexity is prohibitive for large models with

numerous nodes and higher-density layers. To address this issue, we introduce a more efficient, node-based paradigm that eliminates the reliance on length-three paths between seed nodes. Instead, this approach focuses on the common neighbors of seed nodes. The node-based version of CH3-L3p, denoted as CH2-L3n, is defined as follows:

$$\mathbf{CH2\text{-}L3n}(u,v) = \sum_{z \in L3} \frac{di_z^*}{de_z^*}$$
 (4)

Here, u and v denote the seed nodes, while z is the common neighbor on the L3 path Muscoloni et al. (2022), a walk of three consecutive links that connects u to v via two of those intermediate nodes. The terms $di_z^* de_z^*$ represent the number of internal local community links (iLCLs) and external local community links (eLCLs) of node i, with a default increment of 1 to prevent division by zero. Internal local community links (iLCLs) are those that connect nodes belonging to the same local community. Contrarily, external local community links (eLCLs) connect nodes belonging to different communities. Figure 3 gives a visual representation of L2 and L3 paths between two seed nodes u and v, defining their local community.

Appendix D. Sparse topological initialization

Correlated sparse topological initialization. Correlated Sparse Topological Initialization (CSTI) is a physics-informed topological initialization. CSTI generates the adjacency matrix by computing the Pearson correlation between each input feature across the calibration dataset and then selects the predetermined number of links, calculated based on the desired sparsity level, as the existing connections. CSTI performs remarkably better when the layer can directly receive input information. However, for layers that cannot receive inputs directly, it cannot capture the correlations from the start since the model is initialized randomly, as in the case of the Transformer. Therefore, in this article, we aim to address this issue by investigating different network models to initialize the topology, to improve the performance for cases where CSTI cannot be directly applied.

Bipartite scale-free model. In artificial neural networks (ANNs), fully connected networks are inherently bipartite. This article explores initializing bipartite networks using models from network science. The Bipartite Scale-Free (BSF) Zhang et al. (2024b) network model extends the concept of scale-freeness to bipartite structures, making them suitable for dynamic sparse training. Initially, the BSF model generates a monopartite Barabási-Albert (BA) model Barabási and Albert (1999), a well-established method for creating scale-free networks in which the degree distribution follows a power law (γ =2.76 in Figure 4). Following the creation of the BA model, the BSF approach removes any connections between nodes of the same type (neuron in the same layer) and rewires these connections to nodes of the opposite type (neuron in the opposite layer). This rewiring is done while maintaining the degree of each node constant to preserve the power-law exponent γ .

Bipartite small-world model. The Bipartite Small-World (BSW) network model Zhang et al. (2024b) is designed to incorporate small-world properties and a high clustering coefficient into bipartite networks. Initially, the model constructs a regular ring lattice and assigns two distinct types of nodes to it. Each node is connected by an equal number of links to

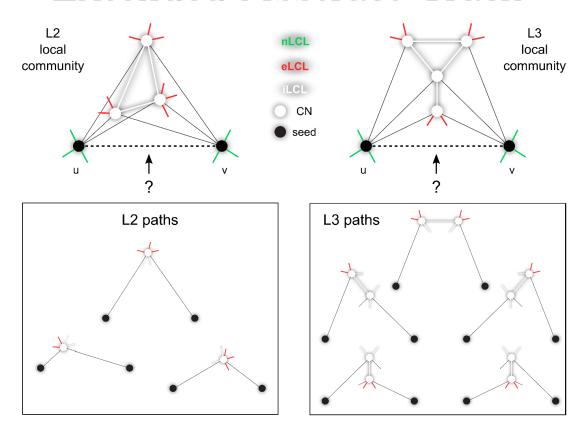


Figure 3: Cannistraci-Hebb epitopological rationale. Muscoloni et al. (2022) The figure illustrates an explanatory example of topological link prediction using the Cannistraci-Hebb epitopological rationale based on either L2 or L3 paths. The two black nodes represent the seed nodes whose unobserved interaction is to be assigned a likelihood score. White nodes denote the common neighbours (CNs) of the seed nodes at either L2 or L3 distance. Together, the set of CNs and the internal local community links (iLCL) constitute the local community. Different link types are color-coded: green for nLCLs, red for external local community links (eLCLs), and white for iLCLs. The L2 (path length 2) and L3 (path length 3) paths associated with the illustrated communities are highlighted. Notably, in artificial neural networks (ANNs), linear layers correspond to bipartite networks, which inherently support only L3 path predictions, as shown in Figure 1.

the nearest nodes of the opposite type, fostering high clustering but lacking the small-world property. Similar to the Watts-Strogatz model (WS) Watts and Strogatz (1998), the BSW model introduces a rewiring parameter, β , which represents the percentage of links randomly removed and then rewired within the network. At $\beta=1$, the model transitions into an Erdős-Rényi model ERDdS and R&wi (1959), exhibiting small-world properties but without a high clustering coefficient, which is popular as the topological initialization of the other dynamic sparse training methods Mocanu et al. (2018); Evci et al. (2020); Yuan et al. (2021).

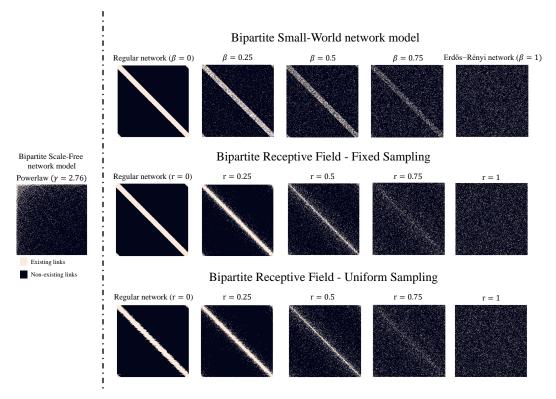
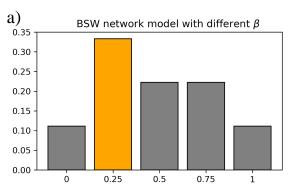


Figure 4: **The adjacency matrix** of the Bipartite Scale-Free (BSF) network model compared to those of the Bipartite Small-World (BSW) network, the Bipartite Receptive Field with fixed sampling (BRF_f), and the Bipartite Receptive field with uniform sampling (BRF_u) as parameters β and r vary between 0 and 1. a) The BSF model inherently forms a scale-free network characterized by a power-law distribution with $\gamma = 2.76$. b) As β changes from 0 to 1, the network exhibits reduced clustering. It is important to note that when $\beta = 0$, the BSW model does not qualify as a small-world network. c) As r increases towards 1, the adjacency matrix becomes more random, while sampling the output neurons' degrees from a fixed or uniform distribution.

Bipartite receptive field model. The Bipartite Receptive Field (BRF) model is a random network generation technique designed to mimic the receptive field phenomenon in the brain networks. The process involves adding links to the adjacency matrix of the bipartite network, with the connectivity structured around the main diagonal according to a parameter $r \in [0, 1]$. A low value of r results in links that are primarily clustered around the diagonal, while a higher value of r leads to a more random connectivity pattern. Specifically, a bipartite adjacency matrix with links near the diagonal indicates that adjacent nodes from the two layers are linked, whereas links far from the diagonal correspond to more distant node pairs. Mathematically, consider an $N \times M$ bipartite adjacency matrix $M_{i,j_{i=1,\ldots,N},j_{i=1,\ldots,N}}$, where M represents the input size and N represents the output size. Each entry of the matrix $m_{i,j}$ is set to 1 if input node i is connected to output node j, and 0 otherwise. A scoring function $S_{i,j}$ is assigned to each connection in the adjacency matrix based on its distance to the main



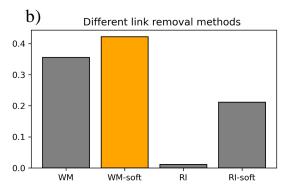


Figure 5: The ablation test of the β of the bipartite small world (BSW) model and the removal methods in CHTs. a) evaluates the influence of the rewiring rate β on the model performance when initialized with the BSW model. b) assesses the influence of link removal selecting from the weight magnitude (WM), weight magnitude soft (WM-soft), relative importance (RI), and relative importance soft (RI-soft). We utilize the win rate of the compared factors under the same setting across each realization of 3 seeds for all experiment combinations on MLP. The factor with the highest win rate is highlighted in orange.

diagonal. This score is given by:

$$S_{i,j} = d_{ij}^{\frac{1-r}{r}}, (5)$$

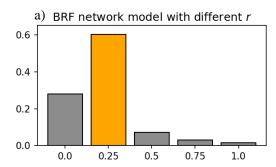
where

$$d_{ij} = min\{|i - j|, |(i - M) - j|, |i - (j - N)|\}$$
(6)

is the distance between the input and output neurons. Therefore, $S_{i,j}$ represents the distance from the diagonal, raised to the power of $\frac{1-r}{r}$. The parameter r controls how structured or random the adjacency matrix is. As $r \to 0$, the scoring function becomes more deterministic, with high scores assigned to entries near the diagonal and low scores to entries farther away. Conversely, as $r \to 1$, all scores $S_{i,j}$ become more uniform, leading to a more random, less structured adjacency matrix. The next step is to determine the degree distribution for the output nodes. This can either be fixed, assigning the same degree to all output nodes, or uniform, where the degrees are randomly sampled from a uniform distribution. Hence, we propose two variations of the BRF model: the Bipartite Receptive Field with fixed sampling (BRFf), in which the degrees of output nodes are fixed, and the Bipartite Receptive Field with uniform sampling (BRFu), where the degrees of the output nodes follow a uniform distribution. This represents an additional enhancement to the WS scheme, which offers no way to control how connections are allocated among the output nodes. In conclusion, to run the BRF model, the user should input an output degree distribution and a spatial dependent distance randomness.

Appendix E. Equal Partition and Neuron Resorting to enhance bipartite scale-free network initialization

As indicated in SET and CHT Mocanu et al. (2018); Zhang et al. (2024b), trained sparse models typically converge to a scale-free network. This suggests that initiating the network



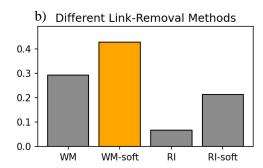


Figure 6: The ablation test of the parameter r in the bipartite receptive field (BRF) model and the removal methods in CHTs using the BRF initialization technique. a) evaluates the influence of the parameter r on the model performance when initialized with the BRF model. b) assesses the influence of link removal in the CHTs model with BRF initialization. We utilize the win rate of the compared factors under the same setting across each realization of 3 seeds for all experiment combinations on MLP. The factor with the highest win rate is highlighted in orange.

with a scale-free structure might initially enhance performance. However, starting directly with a Bipartite Scale-Free model (BSF, power-law exponent $\gamma = 2.76$) does not yield effective results. Upon deeper examination, two potential reasons emerge:

- The BSF model generates hub nodes randomly. However, this random assignment of hub nodes to less significant inputs leads to a less effective initialization, which is particularly detrimental in CHT, which merely utilizes the topology information to regrow new links.
- As demonstrated in CHT, in the final network, the hub nodes of one layer's output should correspond to the input layer of the subsequent layer, which means the hub nodes should have a high degree on both sides of the layer. However, the BSF model's random selection disrupts this correspondence, significantly reducing the number of Credit Assignment Paths (CAP) Zhang et al. (2024b) in the model. CAP is defined as the chain of transformation from input to output, which counts the number of links that go through the hub nodes in the middle layers.

To address these issues, we propose two solutions:

• Equal Partitioning of the First Layer: We begin by generating a BSF model, then rewire the connections from the input layer to the first hidden layer. While keeping the out-degrees of the output neurons fixed, we randomly sample new connections to the input neurons until each of the input neurons' in-degrees reaches the input layer's average in-degree. This approach ensures all input neurons are assigned equal importance while maintaining the power-law degree distribution of output neurons.

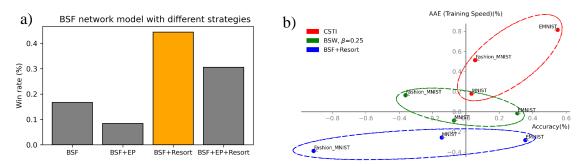


Figure 7: The Performance of the bipartite scale-free model and two enhanced techniques. a) shows the win rate of the Bipartite Scale-Free network model (BSF) with the different techniques. EP stands for equal partition of the first layer, and Resort refers to reordering the neurons based on their degree. b) assesses the comparison between Correlated Sparse Topological Initialization (CSTI), the Bipartite Scale-Free (BSF) model with the best solution from a), and the Bipartite Small-World (BSW) model with $\beta = 0.25$.

• Resorting Middle Layer Neurons: Given the mismatch in hub nodes between consecutive layers, we suggest permuting the neurons between the output of one layer and the input of the next, based on their degree. A higher degree in an output neuron increases the likelihood of connecting to a high-degree input neuron in the subsequent layer, thus enhancing the number of CAPs.

As illustrated in Figure 7, while the two techniques enhance the performance of the BSF initialization, they remain inferior to the BSW initialization. As noted in the main text, achieving scale-freeness is more effective when the model is allowed to learn and adapt dynamically rather than being directly initialized as a predefined structure.

Appendix F. Epitopological Local Minima

Let A_t be the set of existing links in the network at the training step t. Let \tilde{A}_t be the set of removal links and A_t^* be the set of regrown links. The overlap set between removed and regrown links at step t can be quantified as $O_t = \tilde{A}_t \cap A_t^*$. An ELM occurs if the size of O_t at step t is significantly large compared to the size of A_t^* , indicating a high probability of the same links being removed and regrown repeatedly throughout the subsequent training steps. This can be formally represented as $\frac{|O_t|}{|A_t^*|} \geq \theta$, where θ is a predefined threshold close to 1, indicating strong overlap. This definition is essential for the understanding of CHT, as evidenced by the article Zhang et al. (2024b) indicating that the overlap rate between removed and regrown links becomes significantly high within just a few epochs, leading to rapid topological convergence towards the ELM. Previously, CHT implements a topological early stop strategy to avoid predicting the same links iteratively. However, it will stop the topological exploration very fast and potentially trap the model within the ELM.

Appendix G. Soft link removal alternating from RI and Weight magnitude

We illustrate the link removal part of CHTs in Figure 1b1) and b2). We employ two methods, Weight Magnitude (WM) |**W**| and Relative Importance (RI) Zhang et al. (2024a), to remove the connections during dynamic sparse training.

$$\mathbf{RI}_{ij} = \frac{|\mathbf{W}_{ij}|}{\sum |\mathbf{W}_{*j}|} + \frac{|\mathbf{W}_{ij}|}{\sum |\mathbf{W}_{i*}|}$$
(7)

As illustrated in Equation 7, RI assesses connections by normalizing the absolute weight of links that share the same input or output neurons. This method does not require calibration data and can perform comparably to the baseline post-training pruning methods like sparsegpt Frantar and Alistarh (2023) and wanda Sun et al. (2023). Generally, WM and RI are straightforward, effective, and quick to implement in DST for link removal but give different directions for network percolation. WM prioritizes links with higher weight magnitudes, leading to rapid network percolation, whereas RI inherently values links connected to lower-degree nodes, thus maintaining a higher active neuron post-percolation (ANP) rate. The ANP rate is the ratio of the number of active neurons after training compared to the original number of neurons before training. These methods are equally valid but cater to different scenarios. For instance, using RI significantly improves results on the Fashion MNIST dataset compared to WM, whereas WM performs better on the MNIST and EMNIST datasets.

Soft link removal. In the early stages of training, both WM and RI are not reliable due to the model's underdevelopment. Therefore, rather than strictly selecting top values based on WM and RI, we also sample links from a multinomial distribution using an importance score calculated by the removal metrics. The final formula for link removal is defined in Equation 8.

$$\mathbf{S}_{ij} = \left(\frac{|\mathbf{W}_{ij}|/2}{\alpha + (1 - \alpha)\sum |\mathbf{W}_{i*}|} + \frac{|\mathbf{W}_{ij}|/2}{\alpha + (1 - \alpha)\sum |\mathbf{W}_{*j}|}\right)^{\frac{\delta}{1 - \delta}}$$
(8)

Here, α determines the removal strategy, shifting from weight magnitude ($\alpha=1$) to relative importance ($\alpha=0$). In all experiments, we only evaluate these two α values. δ adjusts the softness of the sampling process. As training progresses and weights become more reliable, we adaptively increase δ from 0.5 to 0.75 to refine the sampling strategy and improve model performance. These settings are constant for all the experiments in this article.

Appendix H. Experimental Setup

We evaluate the performance of CHTs using MLPs for image classification tasks on the MNIST LeCun et al. (1998), Fashion MNIST Xiao et al. (2017), and EMNIST Cohen et al. (2017) datasets. To further validate our approach, we apply the sigmoid gradual density decay strategy to Transformers for machine translation tasks on the Multi30k en-de Elliott et al. (2016), IWSLT14 en-de Cettolo et al. (2014), and WMT17 en-de Bojar et al. (2017) datasets. Additionally, we conduct language modeling experiments using the OpenWebText dataset Gokaslan and Cohen (2019). For MLP training, we sparsify all layers except the final layer, as ultra-sparsity in the output layer may lead to disconnected neurons, and the

connections in the final layer are relatively minor compared to the previous layers. For Transformers and LLaMA-130M, we apply dynamic sparse training (DST) to all linear layers, excluding the embedding and final generator layer. Detailed hyperparameter settings for each experiment are provided in Tables 1, 2, and 3.

Appendix I. Density Decay Strategies

In GraNet, the network evolution process consists of three steps: pruning, link removal, and link regrowth. The method first prunes the network to reduce the density, followed by removing and regrowing an equivalent number of links under the updated density. The density decrease in GraNet follows the same approach as Gradual Magnitude Pruning (GMP) Zhu and Gupta (2017), which adheres to a cubic function:

$$s_t = s_f + (s_i - s_f) \left(1 - \frac{t - t_0}{n\Delta t} \right)^3,$$
 (9)

where $t \in \{t_0, t_0 + \Delta t, \dots, t_0 + n\Delta t\}$, s_i is the initial sparsity, s_f is the target sparsity, t_0 is the starting epoch of gradual pruning, t_f is the end epoch of gradual pruning, and Δt is the pruning frequency.

However, this density decay scheduler exhibits a sharp decline in the initial stages of training, which risks pruning a substantial fraction of weights before the model has sufficiently learned. To mitigate this issue, we propose a sigmoid-based gradual density decrease strategy, defined as Equation 2 in the main text. We set k=6 for all the experiments in this article. This strategy ensures a smoother initial pruning phase, allowing the model to warm up and stabilize before undergoing significant pruning, thereby enhancing training stability and performance.

Since our work focuses on MLP, Transformer, and LLMs, where FLOPs are linearly related to the density of the linear layers, the FLOPs of the whole training process are linearly related to the integral of the density function across the training time. the The integral of the GraNet decrease function from t_0 to t_f is:

$$\int_{t_0}^{t_f} (s_i - s_f) \left(1 - \frac{t - t_0}{n\Delta t} \right)^3 dt$$

$$= \frac{1}{4} (s_i - s_f)(t_f - t_0).$$
(10)

For the sigmoid decrease, the integral is:

$$\int_{t_0}^{t_f'} (s_i' - s_f) \left(\frac{1}{1 + e^{-k\left(t - \frac{t_f' + t_0}{2}\right)}} \right) dt$$

$$= \frac{(s_i' - s_f)(t_f' - t_0)}{2}.$$
(11)

To maintain consistency in the computational cost (FLOPs) during training compared to the cubic decay strategy, we reduce the number of steps in the sigmoid-based gradual density decrease by half.

In addition to refining the decay function, we replace the weight magnitude criterion used in the original GMP and GraNet processes with relative importance (RI). This adjustment is motivated by prior work Zhang et al. (2024a), which has shown that RI provides a significant performance advantage over weight magnitude, particularly when pruning models initialized with high density.

Appendix J. Network percolation and extension to Transformer.

We have adapted network percolation Li et al. (2021); Zhang et al. (2024b) to suit the architecture of the Transformer after link removal. The core idea is to identify inactive neurons, which are characterized by having no connections on either one or both sides within a layer of neurons. Such neurons disrupt the flow of information during forward propagation or backpropagation. In addition, Layer-wise computation of the CH link prediction score further implies that neurons without connections on one side are unlikely to form connections in the future. Therefore, network percolation becomes essential to optimize the use of remaining links.

As shown in Figure 1, network percolation encompasses two primary processes: c1) inactive neuron removal to remove the neurons that lack connections on one or both sides; c2) incomplete path adjustment to remove the incomplete paths where links connect to the inactive neurons after c1). Typically applied in simpler continuous layers like those in an MLP, network percolation requires modification for more complex structures. For example, within the Transformer's self-attention module, the outputs of the query and key layers undergo a dot product operation. It necessitates percolation in these layers to examine the activity of the neurons in both output layers at the same position. Similar interventions are necessary in the up_proj and gate_proj layers of the MLP module in the LLaMA model family Touvron et al. (2023a,b).

Appendix K. Baseline Methods

K.1. Fixed Density Dynamic Sparse Training Methods

SET Mocanu et al. (2018): Removes connections based on weight magnitude and randomly regrows new links.

RigL Evci et al. (2020): Removes connections based on weight magnitude and regrows links using gradient information, gradually reducing the proportion of updated connections over time.

CHT Zhang et al. (2024b): A state-of-the-art (SOTA) gradient-free method that removes links with weight magnitude and regrows links based on CH3-L3 scores. CHT is often applied with early stopping to mitigate its computational complexity when working with large models.

K.2. Gradual Density Decrease Dynamic Sparse Training Methods

GMP Han et al. (2015); Zhu and Gupta (2017): Prunes the network with weight magnitude and gradually decreases the density based on Equation 10. Although originally a pruning

method, GMP is treated as a dynamic sparse training method in their implementation Zhu and Gupta (2017), as it stores historical weights and allows pruned weights to reappear during training, since, during training, the pruning threshold might change.

 $\mathbf{MEST}_{EM\&S}$ Yuan et al. (2021): Implements a two-stage density decrease strategy as described in the original work. It removes links based on the combination of weight magnitude and 0.01*gradient and regrows new links randomly.

GraNet Liu et al. (2021): Gradually decreases density using Equation 10. Similar to RigL, GraNet removes links based on the weight magnitude and regrows new links with the gradient of the existing links.

Table 9: Float32 Precision Comparison on LLaMA-130M. Bold values denote the best performance among DST methods. Lower perplexity corresponds to better model performance. s_i represents the initial sparsity for DST methods employing a density decay strategy.

Method	Sparsity		
	70%	80%	
FC	17.07		
RigL	18.34	19.64	
CHTs	17.99	19.25	
GraNet $(s_i = 0.5)$	17.92	18.79	
CHTss $(s_i = 0.5)$	17.76	18.69	

Appendix L. Ablation and Sensitivity Tests

An overall ablation test To fully assess each component's effectiveness, we conduct several ablation and sensitivity tests that help us understand how to select a sparse topological initialization and identify the best link removal and regrowth methods. We first made a global test for all the components in Table 10, which shows the effectiveness of each element introduced by this article. The node-based and path-based link regrowth methods have comparable performance, but the node-based versions are much faster.

Sparse topological initialization. For sparse topological initialization, we compare BRF, BSW, BSF, and CSTI Zhang et al. (2024b) across three image classification datasets, as shown in Figure 7b. The results indicate that when the inputs can directly access task-relevant information, CSTI consistently achieves the best performance. In general, BRF and BSW perform similarly under these conditions, but outperform the BSF initialization.

To further validate our findings, we evaluate BRF and BSW network initializations on machine translation tasks using Transformer models. Figure 8 and Figure 9 present the performance comparisons between BSW and BRF on the Multi30k and IWSLT datasets, while Figure 10 shows the win-rate analysis. These comparisons demonstrate that BRF consistently outperforms BSW across most cases. Additionally, Figure 6a analyzes the impact of the receptive field range r on BRF initialization for MNIST, Fashion MNIST, and EMNIST tasks using MLPs, with results indicating that r=0.25 yields the best performance.

Table 10: Ablation results of Transformer on Multi30K and IWSLT datasets at 90% sparsity. The scores indicate BLEU scores, the higher the better. Bold values denote the best performance among DST methods.

Variant	Multi30K (90% sparsity)	IWSLT (90% sparsity)
a. CHT	28.38	19.91
b. CHTss without node-based implementation	32.68 (2.42 hours)	24.82 (18 hours)
c. CHTss without soft sampling	28.92	21.88
d. CHTss without sigmoid decay (= CHTs)	30.35	21.60
e. CHTss (full model)	32.79 (0.25 hours)	24.57 (1.5 hours)

Building on this prior knowledge, we further evaluate BRF on LLaMA-60M and LLaMA-130M models, testing r values in the range [0,0.3] and comparing two different degree distributions. The results, shown in Table 7 and Table 8, indicate that on LLaMA models, the choice of r and distribution has limited impact. While r=0.1 wins slightly more often, the improvements remain marginal. Finally, Table 6 reports the best performance combinations of r and degree distributions derived from these evaluations.

Table 11: Performance comparison of CHTs and CHTss at 90% sparsity across different removal methods. The tested dataset is Multi30K, and the reported metric is BLEU, which is the higher the better.

Remove Method	CHTs	CHTss
set	28.82	25.76
wm	28.17	31.15
${ m wm_soft}$	30.35	32.79
ri	28.91	32.20
ri_soft	27.86	31.86
MEST	28.70	32.07
snip	28.23	31.66
sensitivity	29.02	29.73
Rsensitivity	28.18	30.67

Link removal. We first conduct a simple evaluation of the link removal methods introduced in this article when changing the α and δ inside Figure 1b2) on Figure 5b) and Figure 6b). The removal methods are selected from Weight Magnitude (WM), Weight Magnitude soft (WMs), Relative Importance (RI), and Relative Importance soft (RIs). For WM we fix the hyperparameters $\alpha = 1$ and $\delta = 1$; for RI we fix the hyperparameters $\alpha = 0$ and $\delta = 0.5$; for

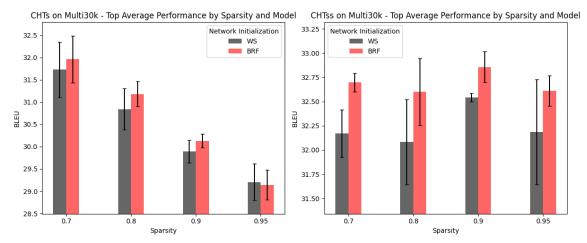


Figure 8: Top average BLEU for WS (grey) and BRF initialization methods on the Multi30k translation dataset of CHTs (left) and CHTss (right, with sigmoid decay) at different sparsity levels. Error bars denote the standard error across three seeds.

WMs we fix $\alpha = 1$ and we let δ increase linearly from 0.5 to 0.9; for RIs we fix $\alpha = 0$ and let δ increase linearly from 0.5 to 0.9. From the results, it can be observed that WMs performs the best in most cases. We compare these methods with those in Nowak et al. (2023) in Table 11 on two machine translation tasks. The results indicate that using WMs as a link removal method generally outperforms the alternatives.

We also evaluate how to define the softness in WMs. During sampling, we have a hyperparameter to decide the temperature of the scores that convert to the probability of being removed. We perform a test using a linear decay solution, since, generally, the weights in the model become more reliable as training progresses. Figure 11 shows the variation in BLEU scores as we change the starting and ending values of the δ parameter in the soft weight magnitude removal method on transformer models. Recalling that we define the temperature by $T=\frac{1}{1-\delta}$, we observe that for a simple benchmark like Multi30k, a high starting temperature produces better performance. This is motivated by the fact that loss decreases very fast through epochs, meaning that weights are learned quickly, and we can deterministically remove weights with high reliability. In more complex datasets, like IWSLT, low starting temperatures are preferred. This is because during the early stages of training, weights are learned slowly, meaning that a deterministic removal can be less reliable. To be more consistent, we select a start $\delta=0.5$ and end $\delta=0.9$ for all the tasks in the main article.

Appendix M. Historical weights

Inspired by GMP Han et al. (2015); Zhu and Gupta (2017), we incorporate historical weights into our CHTs and CHTss implementation. During training, we maintain a historical weight matrix that records previously learned weights throughout the training process. When CHTs and CHTss predict new links, we initialize them using their corresponding historical weights - specifically, the values they held before being pruned. In this way, CHTs and CHTss

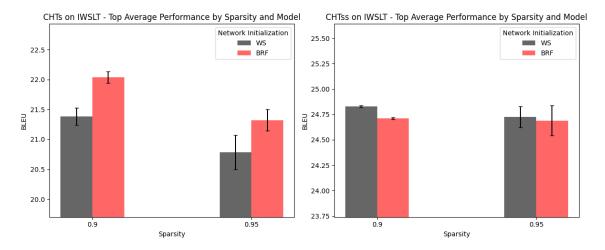


Figure 9: Top average BLEU for WS (grey) and BRF initialization methods on the IWSLT translation dataset of CHTs (left) and CHTss (right, with sigmoid decay) at different sparsity levels. Error bars denote the standard error across two seeds.

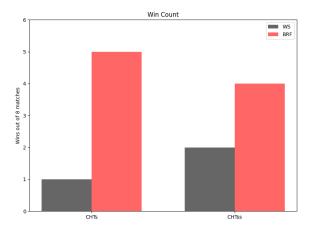


Figure 10: Win rates of BRF against WS over CHTs and CHTss models on different datasets (Multi30k and IWSLT) and different sparsities (0.9 and 0.95 for IWSLT and 0.7, 0.8, 0.9, 0.95 for Multi30k).

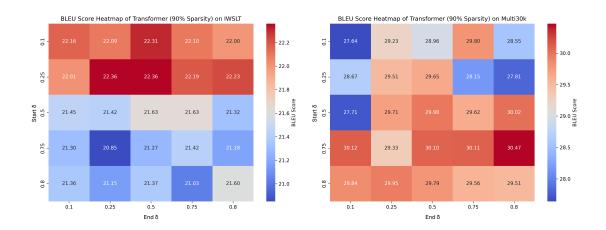


Figure 11: Investigating the level of randomness in link removal strategies. Top BLEU scores of the transformer model using CHTs with weight magnitude soft removal strategy, as the initial and final values of δ take values in $\{0.1, 0.25, 0.5, 0.75, 0.8\}$.

enable weight recovery with preserved memory, allowing the model to retain valuable prior information.

Appendix N. MLP for image classification

Ablation Test. Using MLP, we conduct an ablation study on each component proposed within the CHTs framework to determine the most effective implementation to apply next for the Transformer model. Figure 5a) compares the topologies initialized with the Bipartite Small-World (BSW) model at different values of β , clearly indicating that $\beta=0.25$ yields the best results. Figures 5b) assess the link removal methods, concluding that the weight magnitude soft (WM-soft) method outperforms all others. We consider the best settings showcased in these results to decide the CHTs strategy for training Transformers and LLaMA-130M.

Main Results. In the MLP evaluation, we aim to assess the fundamental capacity of DST methods to train the fully connected module, which is common across many ANNs. The sparse topological initialization of CHT and CHTs is CSTI Zhang et al. (2024b) since the input bipartite layer can directly receive information from the input pixels. Table 4 displays the performance of DST methods compared to their fully connected counterparts across three basic datasets of MNIST, Fashion MNIST, and EMNIST. The DST methods are tested at 99% sparsity. As shown in Table 4, both of the two regrowth methods of CHTs outperform the other fixed sparsity DST methods. Notably, the path-based CH3-L3p outperforms the fully connected one in all the datasets. The node-based CH2-L3n also achieves comparable performance on these basic datasets. However, considering the running time of CH3-L3p is

unacceptable, especially in large scale models, in the rest of the experiments of this article, we only use CH2-L3n as the representative method to regrow new links. Table 12 presents a comparison of fixed-sparsity dynamic sparse training (DST) methods against the fully connected (FC) baseline. Notably, CHTs outperforms all other DST methods and achieve an 11% improvement in accuracy over the fully connected model. In addition, we present the active neuron post-percolation rate (ANP) for each method in Table 4 and Table 12. It is evident that CHTs adaptively percolates the network more effectively while retaining performance.

Table 12: Performance comparison of different sparsity dynamic sparse training methods on the CIFAR10 dataset trained on an MLP at 99% sparsity. The density decay of GMP, GraNet, and CHTss starts with a sparsity of 50%. ACC represents accuracy, and ANP denotes the active neuron percolation rate, indicating the final size of the network. The lowest anp rate and the best dynamic sparse training method are highlighted in bold, and performances surpassing the fully connected model are marked with "*". The results present a standard error taken over three seeds of the experiments.

Method	ACC (%)	Comparison to F	C ANP
FC	62.85 ± 0.16	_	_
CHTs CHT RigL SET	$69.97 \pm 0.06^*$ 59.10 ± 0.06 $63.90 \pm 0.19^*$ 62.70 ± 0.11	+11.33% $-5.97%$ $+1.67%$ $-0.24%$	54% 96% 59% 100%
CHTss GraNet GMP		$+13.43\% \\ +10.28\% \\ +3.60\%$	63% 61% 75%

Appendix O. Extra results of LLaMA1b

Table 13: Validation perplexity of different dynamic sparse training (DST) methods on OpenWebText using LLaMA-1B across varying sparsity levels. Lower perplexity corresponds to better model performance. The performances that surpass the fully connected model are marked with "*".

Sparsity	0.7	0.9	0.95
FC		14.62	
CHTs CHTss	14.53* 15.15	17.14 15.62	10.00

Language modeling. We present a comparison of CHTs, CHTss, and fully connected network on language modeling tasks using the LLaMA-1B model on Table 13. The results clearly demonstrate that CHTs consistently outperform the fully connected (FC) baseline

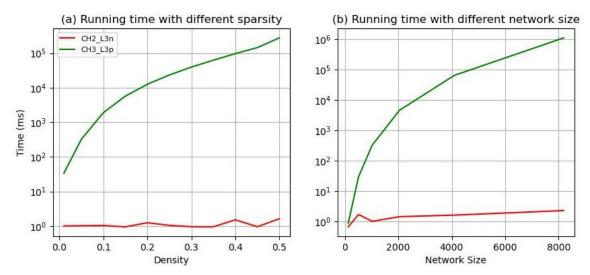


Figure 12: One-time Link Prediction Runtime Performance Evaluation of node-based and path-based methods across varying densities and network sizes. In (a), the network size is fixed at 1024×1024 , while in (b), the density is fixed at 5%.

at 70%, even at a high sparsity of 95%, CHTss achieves a perplexity of 16.51, which is remarkably close to the FC baseline.

Appendix P. Experiments compute resources

All experiments were conducted on NVIDIA A100 80GB GPUs. MLP and Transformer models were trained using a single GPU, while LLaMA models were trained using eight GPUs in parallel.