
Energy Transformer

Benjamin Hoover*
IBM Research
Georgia Tech
benjamin.hoover@ibm.com

Yuchen Liang*
Department of CS
RPI
liangy7@rpi.edu

Bao Pham*
Department of CS
RPI
phamb@rpi.edu

Rameswar Panda
MIT-IBM Watson AI Lab
IBM Research
rpanda@ibm.com

Hendrik Strobelt
MIT-IBM Watson AI Lab
IBM Research
hendrik.strobelt@ibm.com

Duen Horng Chau
College of Computing
Georgia Tech
polo@gatech.edu

Mohammed J. Zaki
Department of CS
RPI
zaki@cs.rpi.edu

Dmitry Krotov
MIT-IBM Watson AI Lab
IBM Research
krotov@ibm.com

Abstract

Our work combines aspects of three promising paradigms in machine learning, namely, attention mechanism, energy-based models, and associative memory. Attention is the power-house driving modern deep learning successes, but it lacks clear theoretical foundations. Energy-based models allow a principled approach to discriminative and generative tasks, but the design of the energy functional is not straightforward. At the same time, Dense Associative Memory models or Modern Hopfield Networks have a well-established theoretical foundation, and allow an intuitive design of the energy function. We propose a novel architecture, called the Energy Transformer (or ET for short), that uses a sequence of attention layers that are purposely designed to minimize a specifically engineered energy function, which is responsible for representing the relationships between the tokens. In this work, we introduce the theoretical foundations of ET, explore its empirical capabilities using the image completion task, and obtain strong quantitative results on the graph anomaly detection and graph classification tasks.

1 Introduction

Transformers have become pervasive models in various domains of machine learning, including language, vision, and audio processing. Every transformer block uses four fundamental operations: attention, feed-forward multi-layer perceptron (MLP), residual connection, and layer normalization. Different variations of transformers result from combining these four operations in various ways. For instance, [2] propose to frontload additional attention operations and backload additional MLP layers in a sandwich-like instead of interleaved way, [3] prepend an MLP layer before the attention in each transformer block, [4] use neural architecture search methods to evolve even more sophisticated transformer blocks, and so on. Various methods exist to approximate the attention operation, multiple modifications of the norm operation, and connectivity of the block; see, for example, [5] for a taxonomy of different models. At present, however, the search for new transformer architectures

*B.Hoover, Y.Liang, and B.Pham equally contributed to this work. See [1] for the official conference version of this work.

is driven mostly by empirical evaluations, and the theoretical principles behind this growing list of architectural variations is missing.

In a seemingly unrelated line of work, Associative Memory models, also known as Hopfield Networks [6, 7], have been gaining popularity in the machine learning community thanks to theoretical advancements pertaining to their memory storage capacity and novel architectural modifications. Specifically, it has been shown that increasing the sharpness of the activation functions can lead to super-linear [8] and even exponential [9] memory storage capacity for these models, which is important for machine learning applications. Most importantly, the attention mechanism in transformers is demonstrated to be closely related to a special model of this family with the softmax activation function [10].

As a result, we use the recent theoretical advancements and architectural developments in Dense Associative Memories to design an energy function tailored to route the information between the tokens. The goal of this energy function is to represent the relationships between the semantic contents of tokens describing a given data point (e.g., the relationships between the contents of the image patches in the vision domain, or relationships between the nodes’ attributes in the graph domain). The core mathematical idea of our approach is that the sequential application of these unusual transformer blocks, which we call the Energy Transformer (ET), minimizes this global energy function. Thus, the sequence of conventional transformer blocks is replaced with a single ET block, which iterates the token representations until they converge to a fixed point attractor state. In the image domain, this fixed point corresponds to the completed image with masked tokens replaced by plausible auto-completions of the occluded image patches. In the graph domain, the fixed point reveals the anomaly status of a node given its neighbors, or the graph label. The energy function in our ET block is designed with the goal to describe the *relationships between the tokens* (e.g., continuity of straight lines in images, interaction between node features and label on graphs).

The core mathematical principle of the ET block – the existence of the global energy function – dictates strong constraints on the possible operations inside the block, the order in which these operations are executed in the inference pass, and the symmetries of the weights in the network. As a corollary of this theoretical principle, the attention mechanism of ET is different from the attention mechanism commonly used in feed-forward transformers [11]; and though the attention energy of ET is related to recent works on Modern Hopfield Networks [8, 9, 10], the overall formulation of the energy function of ET is unique (see Appendix C for further explanation). Lastly, our network may be viewed as an example of a broader class of Energy-Based Models [12] frequently discussed in the AI community. The proposed model is defined through the specific choice of the energy function, which on one hand is suitable for the computational task of interest, and on the other hand results in optimization equations closely related to the forward pass in feed-forward transformers.

2 Overview of ET

The overall pipeline of ET (shown in Figure 1) follows closely to that of the original transformer and its variants (such as Vision Transformer [13]). The input to ET can be seen as a series of physical *particles*, where they may either represent *patches* of an image, *words* in text, or *nodes* in a graph. Each particle is first to pass through an encoder along with adding positional information afterwards —yielding a token representation $\mathbf{x}_A \in \mathbb{R}^D$, where $A = 1 \dots N$ enumerates the number of particles. Given a token \mathbf{x}_A and prior to passing through the ET block, we apply a layer-normalize operation to yield the normalized token representation \mathbf{g}_A , see equation 6. This particular operation can be viewed as an activation for neurons and defined as a partial derivative of the Lagrangian function 7 (as discussed in [14, 15, 16]). See Appendix A for details on mathematical notation and Appendix B for full details on the design of ET.

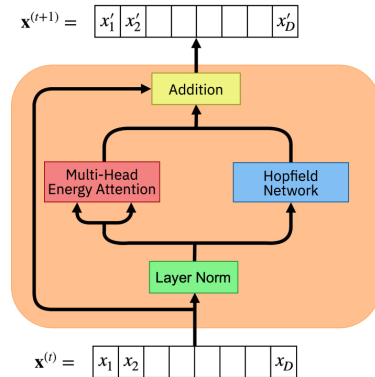


Figure 1: Inside the ET block. Input token \mathbf{x} passes through a sequence of operations and gets updated to produce the output token \mathbf{x}' . The operations inside the ET block are carefully engineered so that the entire network has a global energy function, which decreases with time and is bounded from below. The ET-based analogs of the attention module and the feed-forward MLP module are applied in parallel as opposed to consecutively as in conventional transformers.

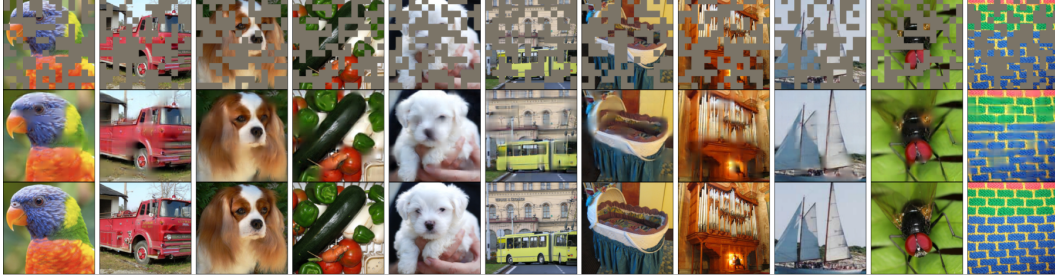


Figure 2: Reconstruction examples of our Energy Transformer using images from the ImageNet-1k validation set. *Top row*: input images where 50% of the patches are masked with the learned MASK token. *Middle row*: output reconstructions after 12 time steps. *Bottom row*: original images.

The first contribution to the ET’s energy function is responsible for exchanging information between particles. Similarly to the regular attention mechanism, each token generates a pair of queries and keys; where ET does not have a value matrix but instead the value matrix is a function of keys and queries. The goal of this energy function is to evolve the tokens in a meaningful way where the keys of the “open” particles are aligned with the queries of the “masked” particles in the internal space of the attention operation. The energy is thus described as:

$$E^{\text{ATT}} = -\frac{1}{\beta} \sum_h \sum_C \log \left(\sum_{B \neq C} \exp(\beta A_{hBC}) \right) \quad (1)$$

where $h = 1 \dots H$ index the heads and $B, C = 1 \dots N$ index the tokens. Meanwhile, A_{hBC} represents the attention matrix computed from query and key tensors obtained from learnable parameters $\mathbf{W}^K \in \mathbb{R}^{Y \times H \times D}$ and $\mathbf{W}^Q \in \mathbb{R}^{Y \times H \times D}$, where Y represents the dimensionality of each head h . The energy is minimal when all queries are aligned with the keys of a small number of other particles in the attention map, such that different heads (at index h) contribute to the energy additively.

The second contribution to ET, the Hopfield Network (HN), is responsible for ensuring that the token representations are consistent with what one expects to see in the original input space. The energy of this component is defined as:

$$E^{\text{HN}} = -\sum_B \sum_{\mu} G \left(\sum_j \xi_{\mu j} g_{jB} \right), \quad \xi \in \mathbb{R}^{K \times D} \quad (2)$$

where $\mu = 1 \dots K$ index the memories and $j = 1 \dots D$ index the token dimension. $\xi_{\mu j}$ is therefore a set of learnable weights or memories in the network and $G(\cdot)$ is an integral of a chosen activation function $r(\cdot)$, where $G(\cdot)' = r(\cdot)$, as described in [14, 16]. For example, if $r(\cdot)$ is chosen to be the softmax function, $G(\cdot)$ is then the log-sum-exp function. The energy contribution of this block is minimal when the token represents are aligned with some rows (memories) of the matrix ξ .

The inference pass of the ET network is described by the continuous time differential equation, which minimizes the sum of the two energies described above:

$$\tau \frac{dx_{iA}}{dt} = -\frac{\partial E}{\partial g_{iA}}, \quad \text{where } E = E^{\text{ATT}} + E^{\text{HN}} \quad (3)$$

Here, x_{iA} is the token representation (input and output from ET block), and g_{iA} is its layer-normalized version. The dynamical system finds a trade-off between the alignment of each token’s queries to the keys of its neighbors and the consistency of the token to the expectation of the identity of the particle in the original space, as described in E^{ATT} and E^{HN} . Most importantly, this dynamical system is guaranteed to minimize the energy, see equation 12.

3 ET for Masked Image Completion on ImageNet

We trained* the ET network on the masked image completion task using ImageNet-1k dataset [17]. Each image was broken into non-overlapping patches of 16x16 RGB pixels, which were projected

*The code is available: <https://github.com/bhoov/energy-transformer-jax>.

with a single affine encoder into the token space. Half of these tokens were “masked”, e.g., by replacing them with a learnable MASK token. Our ET block then processes all tokens recurrently for T steps. The updated token representations after the dynamics are passed to a simple linear decoder. The loss function is the standard MSE loss on the occluded patches.

Examples of occluded/reconstructed images (unseen during training) are shown in Figure 2. In general, our model learns to perform the task very well, capturing the texture in dog fur (column 3) and understanding meaningful boundaries of objects. However, we observe that our single ET block struggles to understand some global structure, e.g., failing to capture both eyes of the white dog (column 5) and completing irregular brick patterns in the name of extending the un-occluded borders (last column). See Appendix F for further details on training and model inspection.

4 Graph Anomaly Detection

Consider an undirected graph with N nodes. Every node has a vector of raw attributes $\mathbf{y}_A \in \mathbb{R}^F$, where F is the number of node’s features. Every node also has a binary label l_A indicating whether the node is benign or not. We focus on node anomaly and assume that all edges are trusted. The task is to predict the label of a node given the graph structure and the node’s features. Since there are far more benign nodes in the graph than anomalous ones, anomaly detection can be regarded as an imbalanced node classification task.

Similarly to the image domain, the feature vectors for every node are converted to a token representation via a linear embedding along with adding a learnable positional embedding, yielding $\mathbf{x}_A^{t=1}$, where $t = 1$ indicates the time of the update in the ET dynamics. This token is recursively applied through the ET for T iterations. Once T steps are completed, a layer-norm operation is applied onto $\mathbf{x}_A^{t=T}$, such the output of entire model is:

$$\mathbf{g}_A^{\text{final}} = \mathbf{g}_A^{t=1} \parallel \mathbf{g}_A^{t=T} \quad (4)$$

where \parallel denotes concatenation and $\mathbf{g}_A^{t=1}$ is the normalized version of $\mathbf{x}_A^{t=1}$. Following [18], the node representation $\mathbf{g}_A^{\text{final}}$ is fed into an MLP with a sigmoid activation function to compute the anomaly probabilities p_A . The weighted cross entropy

$$\text{Loss} = \sum_A \left[\omega l_A \log(p_A) + (1 - l_A) \log(1 - p_A) \right] \quad (5)$$

is used to train the network. Above, ω is the ratio of the benign labels ($l_A = 0$) to anomalous ones ($l_A = 1$). We compare ET with state-of-the-art approaches for graph anomaly detection* in Table 2, shown alongside full experimental details in Appendix D. Results on graph classification are included in Appendix E.

5 Discussion and Conclusions

A lot of recent research has been dedicated to understanding the striking analogy between Hopfield Networks and the attention mechanism in transformers. At a high level, the main message of our work is that the *entirety* of the transformer block can be viewed as a single large Hopfield Network, not just attention alone. At a deeper level, with the recent advances in the field of Hopfield Networks, we design a novel energy function that is tailored for dynamical information routing between the tokens, and representation of a large number of relationships between those tokens. We evaluated the ET network qualitatively on the image completion task, and quantitatively on node anomaly detection and graph classification. The qualitative investigation reveals the perfect alignment between the theoretical design principles of our network and its empirical computation. Meanwhile, the quantitative evaluation demonstrates strong results, which stand in line or exceed the methods recently developed specifically for these tasks. We believe that the proposed network will be useful for other tasks and domains (e.g., NLP, audio, and video), which provides future directions for a comprehensive investigation.

*The code is available: <https://github.com/zhuergou/Energy-Transformer-for-Graph-Anomaly-Detection/>.

References

- [1] Benjamin Hoover, Yuchen Liang, Bao Pham, Rameswar Panda, Hendrik Strobelt, Duen Horng Chau, Mohammed J Zaki, and Dmitry Krotov. Energy transformer. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [2] Ofir Press, Noah A Smith, and Omer Levy. Improving transformer models by reordering their sublayers. *arXiv preprint arXiv:1911.03864*, 2019.
- [3] Yiping Lu, Zhuohan Li, Di He, Zhiqing Sun, Bin Dong, Tao Qin, Liwei Wang, and Tie-Yan Liu. Understanding and improving transformer from a multi-particle dynamic system point of view. *arXiv preprint arXiv:1906.02762*, 2019.
- [4] David So, Quoc Le, and Chen Liang. The evolved transformer. In *International Conference on Machine Learning*, pages 5877–5886. PMLR, 2019.
- [5] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. *arXiv preprint arXiv:2106.04554*, 2021.
- [6] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.
- [7] John J Hopfield. Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the national academy of sciences*, 81(10):3088–3092, 1984.
- [8] Dmitry Krotov and John J. Hopfield. Dense associative memory for pattern recognition. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [9] Mete Demircigil, Judith Heusel, Matthias Löwe, Sven Upgang, and Franck Vermet. On a model of associative memory with huge storage capacity. *Journal of Statistical Physics*, 168(2):288–299, May 2017.
- [10] Hubert Ramsauer, Bernhard Schäfl, Johannes Lehner, Philipp Seidl, Michael Widrich, Thomas Adler, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [12] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, and Fugie Huang. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- [14] Dmitry Krotov and John J. Hopfield. Large associative memory problem in neurobiology and machine learning. In *International Conference on Learning Representations*, 2021.
- [15] Fei Tang and Michael Kopp. A remark on a paper of krotov and hopfield [arxiv: 2008.06996]. *arXiv preprint arXiv:2105.15034*, 2021.
- [16] Dmitry Krotov. Hierarchical associative memory. *arXiv preprint arXiv:2107.06446*, 2021.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [18] Jianheng Tang, Jiajin Li, Ziqi Gao, and Jia Li. Rethinking graph neural networks for anomaly detection. *arXiv preprint arXiv:2205.15508*, 2022.

- [19] Yongyi Yang, Zengfeng Huang, and David Wipf. Transformers from an optimization perspective. *arXiv preprint arXiv:2205.13891*, 2022.
- [20] Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. Deep anomaly detection on attributed networks. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages 594–602. SIAM, 2019.
- [21] Zhen Peng, Minnan Luo, Jundong Li, Luguoxue, and Qinghua Zheng. A deep multi-view framework for anomaly detection on attributed networks. *IEEE Transactions on Knowledge and Data Engineering*, 2020.
- [22] Tong Zhao, Chuchen Deng, Kaifeng Yu, Tianwen Jiang, Daheng Wang, and Meng Jiang. Error-bounded graph anomaly loss for gnn. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 1873–1882, 2020.
- [23] Kaize Ding, Jundong Li, Nitin Agarwal, and Huan Liu. Inductive anomaly detection on attributed networks. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 1288–1294, 2021.
- [24] Yuchen Liang, Dmitry Krotov, and Mohammed J Zaki. Modern hopfield networks for graph embedding. *Frontiers in big Data*, 5, 2022.
- [25] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [26] Zhiwei Liu, Yingdong Dou, Philip S Yu, Yutong Deng, and Hao Peng. Alleviating the inconsistency problem of applying graph neural network to fraud detection. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, pages 1569–1572, 2020.
- [27] Yingdong Dou, Zhiwei Liu, Li Sun, Yutong Deng, Hao Peng, and Philip S Yu. Enhancing graph neural network-based fraud detectors against camouflaged fraudsters. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 315–324, 2020.
- [28] Yang Liu, Xiang Ao, Zidi Qin, Jianfeng Chi, Jinghua Feng, Hao Yang, and Qing He. Pick and choose: a gnn-based imbalanced learning approach for fraud detection. In *Proceedings of the Web Conference 2021*, pages 3168–3177, 2021.
- [29] Shebuti Rayana and Leman Akoglu. Collective opinion spam detection: Bridging review networks and metadata. In *Proceedings of the 21th acm sigkdd international conference on knowledge discovery and data mining*, pages 985–994, 2015.
- [30] Julian John McAuley and Jure Leskovec. From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews. In *Proceedings of the 22nd international conference on World Wide Web*, pages 897–908, 2013.
- [31] Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs. *arXiv preprint arXiv:2012.09699*, 2020.
- [32] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, pages 233–240, 2006.
- [33] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks, 2022.
- [34] Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. Tudataset: A collection of benchmark datasets for learning with graphs. *CoRR*, abs/2007.08663, 2020.
- [35] Qi Zhao and Yusu Wang. Learning metrics for persistence-based summaries and applications for graph classification. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

- [36] Muhammet Balcilar, Guillaume Renton, Pierre H eroux, Benoit Ga uz ere, S ebastien Adam, and Paul Honeine. Bridging the gap between spectral and spatial domains in graph neural networks. *CoRR*, abs/2003.11702, 2020.
- [37] Zhen Zhang, Jiajun Bu, Martin Ester, Jianfeng Zhang, Zhao Li, Chengwei Yao, Dai Huifen, Zhi Yu, and Can Wang. Hierarchical multi-view graph pooling with structure learning. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2021.
- [38] Dai Quoc Nguyen, Tu Dinh Nguyen, and Dinh Phung. Universal graph transformer self-attention networks. In *Companion Proceedings of the Web Conference 2022, WWW ’22*, page 193–196, New York, NY, USA, 2022. Association for Computing Machinery.
- [39] Miguel Domingue, Rohan Dhamdhere, Naga Durga Harish Kanamarlapudi, Sunand Raghupathi, and Raymond Ptucha. Evolution of graph classifiers. In *2019 IEEE Western New York Image and Signal Processing Workshop (WNYISPW)*, pages 1–5, 2019.
- [40] Mingqi Yang, Yanming Shen, Rui Li, Heng Qi, Qiang Zhang, and Baocai Yin. A new perspective on the effects of spectrum in graph neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 25261–25279. PMLR, 17–23 Jul 2022.
- [41] Francesco Orsini, Paolo Frasconi, and Luc De Raedt. Graph invariant kernels. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI’15*, page 3756–3762. AAAI Press, 2015.
- [42] Filippo Maria Bianchi, Daniele Grattarola, Lorenzo Livi, and Cesare Alippi. Hierarchical representation learning in graph neural networks with node decimation pooling. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [43] Ekagra Ranjan, Soumya Sanyal, and Partha Pratim Talukdar. ASAP: Adaptive structure aware pooling for learning hierarchical graph representations. *arXiv preprint arXiv:1911.07979*, 2019.
- [44] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019.
- [45] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: Composable transformations of Python+NumPy programs, 2018.
- [46] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020.
- [47] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*, 2020.
- [48] Alain Hor e and Djemel Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th International Conference on Pattern Recognition*, pages 2366–2369, 2010.

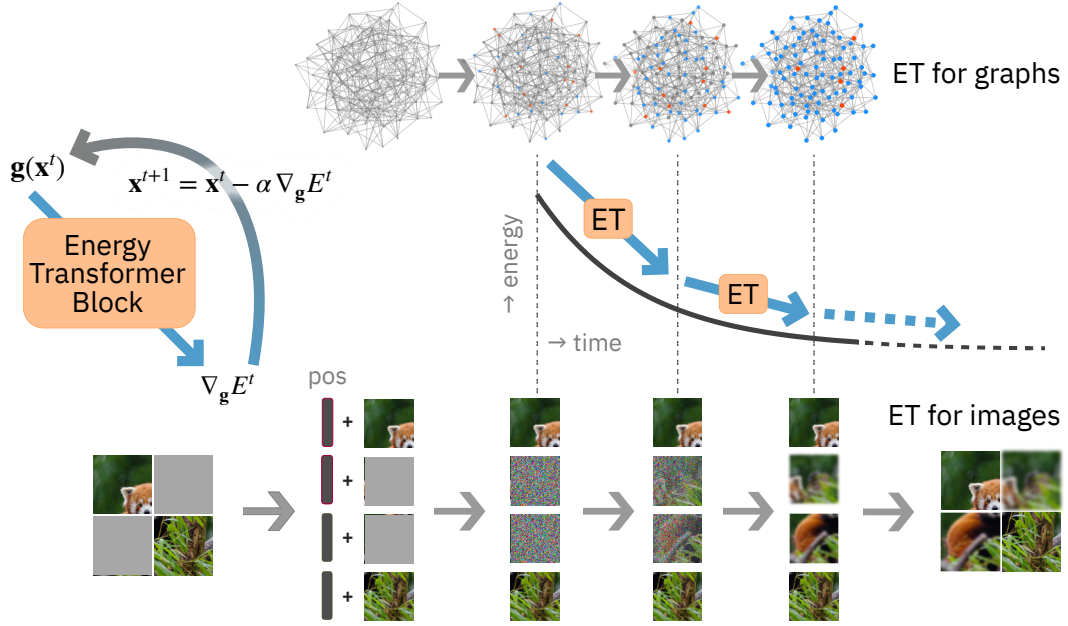


Figure 3: Overview of the Energy Transformer (ET). Instead of a sequence of conventional transformer blocks, a single recurrent ET block is used. The operation of this block is dictated by the global energy function. The token representations are updated according to a continuous time differential equation with the time-discretized update step $\alpha = dt/\tau$. On the image domain, images are split into non-overlapping patches that are linearly encoded into tokens with added learnable positional embeddings (POS). Some patches are randomly masked. These tokens are recurrently passed through ET, and each iteration reduces the energy of the set of tokens. The token representations at or near the fixed point are then decoded using the decoder network to obtain the reconstructed image. The network is trained by minimizing the mean squared error loss between the reconstructed image and the original image. On the graph domain, the same general pipeline is used. Each token represents a node, and each node has its own positional encoding. The token representations at or near the fixed point are used for the prediction of the anomaly status of each node.

A Notations Used in the Main Text and Appendices

Table 1 lists all the notations used in this paper.

B Energy Transformer Block

We now introduce the theoretical framework of the ET network. For clarity of presentation, we use language associated with the image domain. For the graph domain, one should think about “image patches” as nodes on the graph.

The overall pipeline is similar to the Vision Transformer networks (ViTs) [13] and is shown in Figure 3. An input image is split into non-overlapping patches. After passing these patches through the encoder and adding the positional information, the semantic content of each patch and its position is encoded in the token x_{iA} . In the following the indices $i, j, k = 1 \dots D$ are used to denote the token vector’s elements, indices $A, B, C = 1 \dots N$ are used to enumerate the patches and their corresponding tokens. It is helpful to think about each image patch as a physical particle, which has a complicated internal state described by a D -dimensional vector \mathbf{x}_A . This internal state describes the identity of the particle (representing the pixels of each patch), and the particle’s positional embedding (the patch’s location within the image). The ET block is described by a continuous time differential equation, which describes interactions between these particles. Initially, at $t = 1$ the network is given a set containing two groups of particles corresponding to open and masked patches. The “open” particles know their identity and location in the image. The “masked” particles only know where in the image

Table 1: Notations used in the paper.

Notation	Description
F	dimension of node’s feature space
D	dimension of token space
N	number of tokens
Y	number of hidden dimensions in the attention
K	number of hidden dimensions in the Hopfield Network module
T	number of recurrent time steps
H	number of heads
k_h	height of each image patch
k_w	width of each image patch
P	number of pixels per image patch ($3 \times k_h \times k_w$)
\mathbf{y}_A	input feature vector of node A
\mathbf{x}_A	vector representation of token A
x_{iA}	each element of vector representation of token A
\mathbf{g}_A	vector representation of token A after layernorm
g_{iA}	each element of vector representation of token A after layernorm
\mathbf{K}	key tensor
\mathbf{Q}	query tensor
$K_{\alpha h B}$	each element of the key tensor \mathbf{K}
$Q_{\alpha h C}$	each element of the query tensor \mathbf{Q}
A_{hBC}	each element of the attention tensor
l_A	label of node A on graph
$r(\cdot)$	activation function of the Hopfield Network module
$G(\cdot)$	energy contribution of the Hopfield Network module

they are located, but are not provided the information about what image patch they represent. The goal of ET’s non-linear dynamics is to allow the masked particles to find an identity consistent with their locations and the identities of open particles. This dynamical evolution is designed so that it minimizes a global energy function, and is guaranteed to arrive at a fixed point attractor state. The identities of the masked particles are considered to be revealed when the dynamical trajectory reaches the fixed point. Thus, the central question is: how can we design the energy function that accurately captures the task that the Energy Transformer needs to solve?

The masked particles’ search for identity is guided by two pieces of information: identities of the open particles, and the general knowledge about what patches are in principle possible in the space of all possible images. These two pieces of information are described by two contributions to the ET’s energy function: the energy based attention and the Hopfield Network, respectively, for reasons that will become clear in the next sections. Below we define each element of the ET block in the order they appear in Figure 4.

Layer Norm

Each token, or a particle, is represented by a vector $\mathbf{x} \in \mathbb{R}^D$. At the same time, most of the operations inside the ET block are defined using a layer-normalized token representation

$$g_i = \gamma \frac{x_i - \bar{x}}{\sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon}} + \delta_i, \quad \text{where} \quad \bar{x} = \frac{1}{D} \sum_{k=1}^D x_k \quad (6)$$

The scalar γ and the vector elements δ_i are learnable parameters, ε is a small regularization constant. Importantly, this operation can be viewed as an activation function for the neurons and can be defined as a partial derivative of the Lagrangian function (see [14, 15, 16] for the discussion of this property)

$$L = D\gamma \sqrt{\frac{1}{D} \sum_j (x_j - \bar{x})^2 + \varepsilon} + \sum_j \delta_j x_j, \quad \text{so that} \quad g_i = \frac{\partial L}{\partial x_i} \quad (7)$$

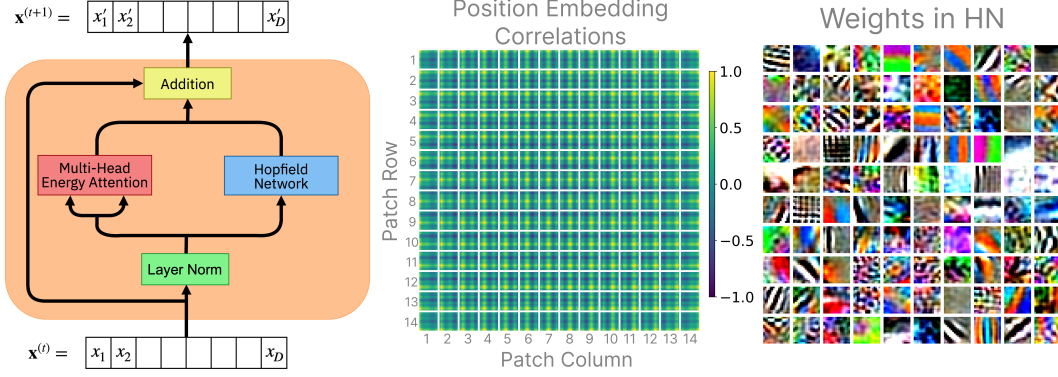


Figure 4: **Left:** Inside the ET block. The input token $\mathbf{x}^{(t)}$ passes through a sequence of operations and gets updated to produce the output token $\mathbf{x}^{(t+1)}$. The operations inside the ET block are carefully engineered so that the entire network has a global energy function, which decreases with time and is bounded from below. In contrast to conventional transformers, the ET-based analogs of the attention module and the feed-forward MLP module are applied in parallel as opposed to consecutively. **Center:** The cosine similarity between the learned position embedding of each patch and every other patch. In each cell, the brightest patch indicates the cell of consideration. **Right:** 100 selected memories stored in the HN memory matrix, visualized by the decoder as 16x16 RGB image patches. This visualization is unique to our model, as traditional Transformers cannot guarantee image representations in the learned weights.

Multi-Head Energy Attention

The first contribution to the ET’s energy function is responsible for exchanging information between the particles (tokens). Similarly to the conventional attention mechanism, each token generates a pair of queries and keys (ET does not have a separate value matrix; instead the value matrix is a function of keys and queries). The goal of the energy based attention is to evolve the tokens in such a way that the keys of the open patches are aligned with the queries of the masked patches in the internal space of the attention operation. Below we use index $\alpha = 1 \dots Y$ to denote elements of this internal space, and index $h = 1 \dots H$ to denote different heads of this operation. With these notations the energy-based attention operation is described by the following energy function:

$$E^{\text{ATT}} = -\frac{1}{\beta} \sum_h \sum_C \log \left(\sum_{B \neq C} \exp(\beta A_{hBC}) \right) \quad (8)$$

where the attention matrix A_{hBC} is computed from query and key tensors as follows:

$$\begin{aligned} A_{hBC} &= \sum_{\alpha} K_{\alpha h B} Q_{\alpha h C}, & \mathbf{A} &\in \mathbb{R}^{H \times N \times N} \\ K_{\alpha h B} &= \sum_j W_{\alpha h j}^K g_{j B}, & \mathbf{K} &\in \mathbb{R}^{Y \times H \times N} \\ Q_{\alpha h C} &= \sum_j W_{\alpha h j}^Q g_{j C}, & \mathbf{Q} &\in \mathbb{R}^{Y \times H \times N} \end{aligned} \quad (9)$$

and the tensors $\mathbf{W}^K \in \mathbb{R}^{Y \times H \times D}$ and $\mathbf{W}^Q \in \mathbb{R}^{Y \times H \times D}$ are learnable parameters.

From the computational perspective each patch generates two representations: query (given the position of the patch and its current content, where in the image should it look for the prompts on how to evolve in time?), and key (given the current content of the patch and its position, what should be the contents of the patches that attend to it?). The log-sum-exp energy function (8) is minimal when for every patch in the image its queries are aligned with the keys of a small number of other patches connected by the attention map. Different heads (index h) contribute to the energy additively.

Hopfield Network Module

The next step of the ET block, which we call the Hopfield Network (HN), is responsible for ensuring that the token representations are consistent with what one expects to see in realistic images. The energy of this sub-block is defined as:

$$E^{\text{HN}} = - \sum_B \sum_{\mu} G \left(\sum_j \xi_{\mu j} g_{jB} \right), \quad \xi \in \mathbb{R}^{K \times D} \quad (10)$$

where $\xi_{\mu j}$ is a set of learnable weights (memories in the Hopfield Network), and $G(\cdot)$ is an integral of the activation function $r(\cdot)$, so that $G(\cdot)' = r(\cdot)$. Depending on the choice of the activation function this step can be viewed either as a classical continuous Hopfield Network [7] if the activation function grows slowly (e.g., $r(\cdot) = \text{ReLU}$), or as a modern continuous Hopfield Network [8, 10, 14] if the activation function is sharply peaked around the memories (e.g., $r(\cdot) = \text{power}$, or softmax). The HN sub-block is analogous to the feed-forward MLP step in the conventional transformer block but requires that the weights of the projection from the token space to the hidden neurons' space to be the same (transposed matrix) as the weights of the subsequent projection from the hidden space to the token space. Thus, the HN module here is an MLP with shared weights that is *applied recurrently*. The energy contribution of this block is low when the token representations are aligned with some rows of the matrix ξ , which represent memories, and high otherwise.

Dynamics of Token Updates

The inference pass of the ET network is described by the continuous time differential equation, which minimizes the sum of the two energies described above

$$\tau \frac{dx_{iA}}{dt} = - \frac{\partial E}{\partial g_{iA}}, \quad \text{where} \quad E = E^{\text{ATT}} + E^{\text{HN}} \quad (11)$$

Here x_{iA} is the token representation (input and output from the ET block), and g_{iA} is its layer-normalized version. The first energy is low when each patch's queries are aligned with the keys of its neighbors. The second energy is low when each patch has content consistent with the general expectations about what an image patch should look like (memory slots of the matrix ξ). The dynamical system (11) finds a trade-off between these two desirable properties of each token's representation. For numerical evaluations, equation (11) is discretized in time.

To demonstrate that the dynamical system (11) minimizes the energy, consider the temporal derivative

$$\frac{dE}{dt} = \sum_{i,j,A} \frac{\partial E}{\partial g_{iA}} \frac{\partial g_{iA}}{\partial x_{jA}} \frac{dx_{jA}}{dt} = - \frac{1}{\tau} \sum_{i,j,A} \frac{\partial E}{\partial g_{iA}} M_{ij}^A \frac{\partial E}{\partial g_{jA}} \leq 0 \quad (12)$$

The last inequality sign holds if the symmetric part of the matrix

$$M_{ij}^A = \frac{\partial g_{iA}}{\partial x_{jA}} = \frac{\partial^2 L}{\partial x_{iA} \partial x_{jA}} \quad (13)$$

is positive semi-definite (for each value of index A). The Lagrangian (7) satisfies this condition.

C Relationship to Modern Hopfield Networks and Conventional Attention

One of the theoretical contributions of our work is the design of the energy attention mechanism and the corresponding energy function (8). Although heavily inspired by prior work [8, 9, 10] on Modern Hopfield Networks, our approach is fundamentally different from it. Specifically, our energy function (8) may look somewhat similar to the energy function of a continuous Hopfield Network with the softmax activation function. The main difference, however, is that in order to use Modern Hopfield Networks *recurrently* (as opposed to applying their update rule only once), the keys must be constant parameters (called *memories* in the Hopfield language). In contrast, in our energy attention network the keys are *dynamical variables* that evolve in time with the queries.

To emphasize this further, it is instructive to write explicitly the ET attention contribution to the update dynamics (11). It is given by (for clarity, assume only one head of attention):

$$- \frac{\partial E^{\text{ATT}}}{\partial g_{iA}} = \sum_{C \neq A} \sum_{\alpha} W_{\alpha i}^Q K_{\alpha C} \text{softmax}_C \left(\beta \sum_{\gamma} K_{\gamma C} Q_{\gamma A} \right) + W_{\alpha i}^K Q_{\alpha C} \text{softmax}_A \left(\beta \sum_{\gamma} K_{\gamma A} Q_{\gamma C} \right)$$

In both terms the softmax normalization is done over the token index of the keys, which is indicated by the subscript in the equation. The first term in this formula is the conventional attention mechanism [11] with the value matrix equal to $\mathbf{V} = (\mathbf{W}^Q)^T \mathbf{K} = \sum_{\alpha} W_{\alpha i}^Q K_{\alpha C}$. The second term is the brand new contribution that is missing in the original attention mechanism. The presence of this second term is crucial to make sure that the dynamical system (11) minimizes the energy function if applied recurrently. This second term is the main difference of our approach compared to the Modern Hopfield Networks. The same difference applies compared to the other recent proposals [19]. Furthermore, the convention we utilize to formulate our energy function of the ET block follows closely to [16], which attempts to emulate the biological interaction between layers of neurons.

Lastly, we want to emphasize that our ET block contains two different kinds of Hopfield Networks acting in parallel, see Figure 4. The first one is the energy attention module, which is *inspired* by, but not identical to, Modern Hopfield Networks, as described above. The second one is the ‘‘Hopfield Network’’ module, which can be either a classical or Modern Hopfield Network. These two should not be confused.

D Details of ET training on Anomaly Detection Task

Graph anomaly detection refers to the process of detecting outliers that deviate significantly from the majority of the samples. Neural network based methods are very popular due to their capability of learning sophisticated data representations. DOMINANT [20] utilizes an auto-encoder framework, using a GCN as an encoder and two decoders for structural reconstruction and attribute reconstruction. ALARM [21] aggregates the encoder information from multiple view of the node attributes. Another study [22], propose a novel loss function to train graph neural networks for anomaly-detectable node representations. In [23] generative adversarial learning is used to detect anomaly nodes where a novel layer is designed to learn the anomaly-aware node representation. Recently, [18] pointed out that anomalies can lead to the ‘‘rightshift’’ of the spectral energy distribution – the spectral energy concentrates more on the high frequencies. They designed a filter that can better handle this phenomenon. We propose a new anomaly detection model from the perspective of Associative Memory (pattern matching), which does not have the over-smoothing problem often faced by GCNs, and has better model interpretability (outliers should be far from the common pattern). We also notice that Modern Hopfield Networks have been used before for node classification, link prediction, and graph coarsening tasks [24].

D.1 Detailed Model Structure for the Graph Anomaly Detection

First, we compute the features that are passed to our energy-based transformer. Each node’s features $\mathbf{y}_A \in R^F$ are mapped into the token space $\mathbf{x}_A \in R^D$, using a linear projection \mathbf{E} . Learnable positional embeddings λ_A are added to this token at $t = 1$,

$$\mathbf{x}_A^{t=1} = \mathbf{E}\mathbf{y}_A + \lambda_A \quad (14)$$

At each time step the input to the ET-block is layer normalized:

$$\mathbf{g}_A^t = \text{LayerNorm}(\mathbf{x}_A^t) \quad (15)$$

Let $\mathbf{W}^Q \in R^{Y \times H \times D}$ and $\mathbf{W}^K \in R^{Y \times H \times D}$ be the query and key weight matrices, respectively. Here Y is the projection dimension in the attention operation, H is the number of heads. We define

$$\begin{aligned} K_{\alpha h B} &= \sum_j W_{\alpha h j}^K g_{j B} \\ Q_{\alpha h C} &= \sum_j W_{\alpha h j}^Q g_{j C} \end{aligned} \quad (16)$$

If we let h indicate the index of the head, we have

$$\Delta x_{iA}^t = \sum_{C \in \mathcal{N}_A} \sum_{h, \alpha} \left[W_{\alpha h i}^Q K_{\alpha h C} \omega_{CA} + W_{\alpha h i}^K Q_{\alpha h C} \omega_{AC} \right] + \sum_{\mu} \xi_{\mu i} r \left(\sum_j \xi_{\mu j} g_{j A} \right) \quad (17)$$

where

$$\omega_{CA} = \text{softmax}_C \left(\beta \sum_{\gamma} K_{\gamma h C} Q_{\gamma h A} \right) \quad (18)$$

Table 2: Performance on Yelp, Amazon, T-Finance, and T-Social datasets with different training ratios. Following [18], mean and standard deviation over 5 runs with different train/dev/test split are reported (standard deviations are only included if they are available in the prior work). Best results are in **bold**. Our model is state of the art or near state of the art on every category.

	Datasets	Split	GraphConsis	CAREGNN	PC-GNN	BWGNN	MLP	GT	ET (Ours)
Macro-F1	Yelp	1%	56.8 \pm 2.8	62.1 \pm 1.3	59.8 \pm 1.4	61.1 \pm 0.4	53.9 \pm 0.2	61.7 \pm 0.4	63.0 \pm 0.6
		40%	58.7 \pm 2.0	63.3 \pm 0.9	63.0 \pm 2.3	71.0 \pm 0.9	57.5 \pm 0.8	68.7 \pm 0.4	71.5 \pm 0.1
	Amazon	1%	68.5 \pm 3.4	68.7 \pm 1.6	79.8 \pm 5.6	90.9 \pm 0.7	74.6 \pm 1.2	88.6 \pm 0.5	89.3 \pm 0.7
		40%	75.1 \pm 3.2	86.3 \pm 1.7	89.5 \pm 0.7	92.2 \pm 0.4	79.1 \pm 1.2	91.7 \pm 0.8	92.8 \pm 0.3
	T-Finance	1%	71.7	73.3	62.0	84.8	61.0	81.5	85.1 \pm 1.0
		40%	73.4	77.5	63.1	86.8	70.5	83.6	88.2 \pm 1.0
	T-Social	1%	52.4	55.8	51.1	75.9	50.0	64.3	79.1 \pm 0.7
		40%	56.5	56.2	52.1	83.9	50.3	68.2	83.5 \pm 0.4
AUC	Yelp	1%	66.4 \pm 3.4	75.0 \pm 3.8	75.4 \pm 0.9	72.0 \pm 0.5	59.8 \pm 0.4	72.5 \pm 0.6	73.2 \pm 0.8
		40%	69.8 \pm 3.0	76.1 \pm 2.9	79.8 \pm 0.1	84.0 \pm 0.9	66.5 \pm 1.0	81.9 \pm 0.5	84.9 \pm 0.3
	Amazon	1%	74.1 \pm 3.5	88.6 \pm 3.5	90.4 \pm 2.0	89.4 \pm 0.3	83.6 \pm 1.7	89.0 \pm 1.2	91.9 \pm 1.0
		40%	87.4 \pm 3.3	90.5 \pm 1.6	95.8 \pm 0.1	98.0 \pm 0.4	89.8 \pm 1.0	95.4 \pm 0.6	97.3 \pm 0.4
	T-Finance	1%	90.2	90.5	90.7	91.1	82.9	90.0	92.8 \pm 1.1
		40%	91.4	92.1	91.2	94.3	87.1	88.2	95.0 \pm 3.0
	T-Social	1%	65.2	71.2	59.8	88.0	56.3	81.4	91.9 \pm 0.6
		40%	71.2	71.8	68.4	95.2	56.9	82.5	93.9 \pm 0.2

Here β controls the temperature of the softmax, \mathcal{N}_A stands for the neighbors of node A — a set of all the nodes connected to node A , r is the ReLU function. Restriction of the attention operation to the neighborhood of a given node is similar to that used in the Graph Attention Networks (GAT), see [25]. Finally, we have residual connection (which is a natural consequence of the discretized differential equation dynamics)

$$\mathbf{x}_A^{t+1} = \mathbf{x}_A^t + \Delta \mathbf{x}_A^t \quad (19)$$

Intuitively, the first term in Equation 17 describes the influence (attention score) of the neighbor nodes with respect to the target node, the second term describes the influence of the target node with respect to each of its neighbor, and the third term is the contribution of the Hopfield Network module. It can be shown that the forward pass of our energy-based transformer layer minimizes the following energy function:

$$E = -\frac{1}{\beta} \sum_C \sum_h \log \left(\sum_{B \in \mathcal{N}_C} \exp \left(\beta \sum_{\alpha} K_{\alpha h B} Q_{\alpha h C} \right) \right) - \frac{1}{2} \sum_{C, \mu} G \left(\sum_j \xi_{\mu j} g_{jC} \right) \quad (20)$$

This energy function will decrease as the forward pass progresses until it reaches a local minimum.

After T iterations when the retrieval is stable, we have the final representation for each node $\mathbf{g}_A^{\text{final}}$ as

$$\mathbf{g}_A^{\text{final}} = \mathbf{g}_A^{t=1} \parallel \mathbf{g}_A^{t=T} \quad (21)$$

where \parallel is the concatenation sign. Following [18], we treat anomaly detection as semi-supervised learning task in this work. The node representation $\mathbf{g}_A^{\text{final}}$ is fed to another MLP with the sigmoid function to compute the abnormal probability p_A , weighted log-likelihood is then used to train the network. The loss function is as follow:

$$\text{Loss} = \sum_A \left[\omega l_A \log(p_A) + (1 - l_A) \log(1 - p_A) \right] \quad (22)$$

where ω is the ratio of normal labels ($l_A = 0$) to anomaly labels ($l_A = 1$).

D.2 Experimental Details

We compare ET with state-of-the-art approaches for graph anomaly detection, which include GraphConsis [26], CAREGNN [27], PC-GNN [28] and BWGNN [18] on four datasets: YelpChi [29], Amazon [30], T-Finance, and T-Social [18]. Additionally, multi-layer perceptrons (MLP) and Graph Transformer (GT) [31] are included as baselines. Following previous work, macro-F1 score (unweighted mean of F1 score) and the Area Under the Curve (AUC) are used as the evaluation metrics on the test datasets [32]. These results are reported in Table 2.

We train all models for 100 epochs using the Adam optimizer with a learning rate of 0.001, and use the model with the best Macro-F1 on the validation set for reporting the final results on the test set. Following [18], we use training ratios 1% and 40% respectively (randomly select 1% and 40% nodes of the dataset to train the model, and use the remaining nodes for the validation and testing). These remaining nodes are split 1:2 for validation:testing. The statistics of the datasets are listed in Table 3. For the four datasets used in the experiments, Amazon and Yelp datasets can be obtained from the DGL library, T-Finance and T-Social can be obtained from [18]. We report the average performance

Dataset	$ V $	$ E $	Anomaly(%)	Features
Amazon	11944	4398392	6.87%	25
Yelp	45954	3846979	14.53%	32
T-Finance	39357	21222543	4.58%	10
T-Social	5781065	73105508	3.01%	10

Table 3: Summary of all the datasets.

of 5 runs on the test datasets. The hyperparameters of our model are tuned based on the validation set, selecting the best parameters within 100 epochs. To speedup the training process, for the large graph datasets T-Finance and T-Social, we sample a different subgraph to train for each epoch (subgraphs have 5% of the nodes with respect to the whole training data). The hyperparameters include the number of hidden dimensions in ET-attention Y , the number of neurons K in the hidden layer within the Hopfield Network Module, the number of time iterations T , and the number of heads H . The weights are learned via backpropagation, which includes embedding projection \mathbf{E} , positional embedding λ_A , softmax inverse temperature parameter β , ET-attention weight tensors \mathbf{W}^Q and \mathbf{W}^K . The optimal hyperparameters used in Table 2 are reported in Table 4. The last row in that table summarizes the range of the hyperparameter search that was performed in our experiments. In general, we have observed that for small datasets (Yelp, Amazon, T-Finance) a 1 or 2 applications of our network is sufficient for achieving strong results, for larger datasets (T-Social) more iterations (3) are necessary. For even bigger dataset (ImageNet) our network needs about 12 iterations.

Dataset	Y	K	T	H
Amazon (40%)	128	640	1	2
Amazon (1%)	64	128	1	1
Yelp (40%)	128	256	1	1
Yelp (1%)	128	256	1	1
T-Finance (40%)	128	256	1	3
T-Finance (1%)	128	256	1	1
T-Social (40%)	128	256	3	3
T-Social (1%)	128	256	3	3
Range of hyperparameters	{64, 128, 256}	{2Y, 3Y, 4Y, 5Y}	{1,2,3}	{1,2,3}

Table 4: Hyperparameters choice of our method on all the datasets.

E Graph Classification with ET

To fully explore the efficacy of ET, we further evaluate its performance on the graph classification problem*. Unlike the anomaly detection task, where we predict a label for every node, in the graph classification task, a single label is predicted for the entire graph.

Consider a graph G , where each node has a raw feature vector $\mathbf{y}_A \in \mathbb{R}^F$ with F feature-dimension. Each vector is projected to the token space yielding $\mathbf{x}_A \in \mathbb{R}^D$, where D is the token dimension. A learnable CLS token \mathbf{x}_{CLS} is concatenated to the set of tokens resulting in the token matrix $\mathbf{X} \in \mathbb{R}^{(N+1) \times D}$, and a learnable linear projection of the top k smallest eigen-vectors of the normalized Laplacian matrix, following [33], is added to \mathbf{X} . Graph structural information is provided to ET within the attention operation. We also use a stacked version of the Energy Transformer consisting of S ET blocks, where each block has the same number of temporal unfolding steps T and its own LayerNorm. The final representation of the CLS token $\mathbf{x}_{CLS}^{\ell=S, t=T}$ is projected via a linear

*The code is available: <https://github.com/Lemon-cmd/Energy-Transformer-For-Graph>.

embedding into the predictor space. A softmax over the number of classes is applied to this predictor representation, and the cross-entropy loss is used to train the network.

Table 5: Graph classification performance on eight datasets of TUDataset. Following [34], mean and standard deviation obtained from 100 runs of 10-fold cross validation are reported. For baselines standard deviations are only included if they are available in prior work. If the entry is unavailable in prior literature it is denoted by ‘-’; best results are in **bold**. The performance difference between non-baseline approaches (including ours) and the baselines (specified by their gray cell) is indicated by ∇ (decrease), \blacktriangle (increase), and ∇ (no change within the error bars) along with the value.

Method	Dataset							
	PROTEINS	NCI1	NCI109	DD	ENZYMES	MUTAG	MUTAGENICITY	FRANKENSTEIN
WKPI (kmeans)	78.5 \pm 0.4 ∇ (6.4)	87.5 \pm 0.5	85.9 \pm 0.4 ∇ (1.5)	82.0 \pm 0.5 ∇ (13.7)	-	85.8 \pm 2.5 ∇ (14.2)	-	-
WKPI (kcenters)	75.2 \pm 0.4 ∇ (9.7)	84.5 \pm 0.5 ∇ (3.0)	87.4 \pm 0.3	80.3 \pm 0.4 ∇ (15.4)	-	88.3 \pm 2.6 ∇ (11.7)	-	-
Spec-GN	-	84.8 \pm 1.6 ∇ (2.7)	83.6 \pm 0.8 ∇ (3.8)	-	72.5 \pm 5.8 ∇ (5.9)	-	-	-
Norm-GN	-	84.9 \pm 1.7 ∇ (2.6)	83.5 \pm 1.3 ∇ (3.9)	-	73.3 \pm 8.0 ∇ (5.1)	-	-	-
GWL-WL	75.8 \pm 0.6 ∇ (9.1)	-	-	-	71.3 \pm 1.1 ∇ (7.1)	-	-	78.9 \pm 0.3
HGP-SL	84.9 \pm 1.6	78.5 \pm 0.8 ∇ (9.1)	80.7 \pm 1.2 ∇ (6.7)	81.0 \pm 1.3 ∇ (14.7)	68.8 \pm 2.1 ∇ (9.6)	-	82.2 \pm 0.6	-
DSGCN	77.3 \pm 0.4 ∇ (7.6)	-	-	-	78.4 \pm 0.6	-	-	-
U2GNN	80.0 \pm 3.2 ∇ (4.9)	-	-	95.7 \pm 1.9	-	88.5 \pm 7.1 ∇ (11.5)	-	-
NDP	73.4 \pm 3.1 ∇ (11.5)	74.2 \pm 1.7 ∇ (13.3)	-	72.8 \pm 5.4 ∇ (22.9)	44.5 \pm 7.4 ∇ (34.9)	87.9 \pm 5.7 ∇ (12.1)	77.9 \pm 1.4 ∇ (4.3)	-
ASAP	74.2 \pm 0.8 ∇ (10.7)	71.5 \pm 0.4 ∇ (16.0)	70.1 \pm 0.6 ∇ (17.3)	76.9 \pm 0.7 ∇ (18.8)	-	-	-	66.3 \pm 0.5 ∇ (12.6)
EvoG	-	-	-	-	55.7 ∇ (22.7)	100.0	-	-
ET (Ours)	90.3\pm0.7 \blacktriangle(5.4)	90.1\pm0.1 \blacktriangle(2.6)	90.5\pm0.1 \blacktriangle(3.1)	95.9\pm0.8 \blacktriangle(0.2)	99.8 \blacktriangle(21.4)	96.6 \pm 0.2 ∇ (3.4)	98.7\pm0.1 \blacktriangle(16.5)	99.8\pm0.1 \blacktriangle(20.9)

E.1 Experimental Evaluation

Eight graph datasets from the TUDataset [34] collection are used for experimentation. NCI1, NCI109, MUTAG, MUTAGENICITY, and FRANKENSTEIN are a common class of graph datasets consisting of small molecules with class labels representing toxicity or biological activity determined in drug discovery projects [34]. Meanwhile, DD, ENZYMES, and PROTEINS represent macromolecules. The task for both DD and PROTEINS is to classify whether a protein is an enzyme. Lastly, for ENZYMES, the task is to assign enzymes to one of the six classes, which reflect the catalyzed chemical reaction [34]. See Table 6 for more details of the datasets.

We compare ET with the current state-of-the-art approaches for the mentioned datasets, which include WKPI-kmeans [35], WKPI-kcenters [35], DSGCN [36], HGP-SL [37], U2GNN [38], and EvoG [39]. Additionally, approaches [36, 40, 41, 42, 43], which are close to the baselines, are included to further contrast the performance of our model. Following the 10-fold cross validation process delineated in [34], accuracy score is used as the evaluation metric and reported in Table 5. In general, we have observed that the modified ET demonstrates strong performance across the eight datasets. Specifically, with the exception of MUTAG, ET beats other methods on all the considered baselines by a substantial margin.

Table 6: The statistics and properties of the eight datasets of TUDataset (additional node attributes are indicated by ‘+’ if exist).

Dataset	Graphs	Avg. Nodes	Avg. Edges	Node Attr	Classes
MUTAG	188	17.93	19.79	7	2
ENZYMES	600	32.63	62.14	18 + 3	6
PROTEINS	1113	39.06	72.82	0 + 4	2
DD	1178	284.32	715.66	89	2
NCI1	4110	29.87	32.30	37	2
NCI109	4127	29.68	32.13	38	2
MUTAGENICITY	4337	30.32	30.77	14	2
FRANKENSTEIN	4337	16.90	17.88	780	2

F Details of Training on ImageNet

We trained the ET network on a masked-image completion task on the ImageNet-1k (IN1K) dataset. We treat all images in IN1K as images of shape 224×224 that are normalized according to standard IN1K practices (mean 0, variance 1 on the channel dimension) and use data augmentations provided by the popular `timm` library [44] (See Table 7). Following the conventional ViT pipeline [13], we split these images into non-overlapping patches of 16×16 RGB pixels which are then projected with a single affine encoder into the token dimension D for a total of 196 encoded tokens per image. We proceed to randomly and uniformly assign 100 of these tokens as “occluded” which are the only tokens considered by the loss function. “Occluded” tokens are designated as follows: of the 100 tokens, 90 tokens are replaced with a learnable MASK token of dimension D and 10 we leave untouched (which we find important for the HN to learn meaningful patch representations). To all tokens we then add a distinct learnable position bias.

These tokens are then passed to our Energy Transformer block which we recur for T steps (the “depth” of the model in conventional Transformers). At each step, the feedback signal (the sum of the energy gradients from our attention block and HN block) is subtracted from our original token representation with a scalar step size $\alpha = \frac{dt}{\tau}$ which we treat as a non-learnable hyperparameter in our experiments. The token representations after T steps are passed to a simple linear decoder (consisting of a layer norm and an affine transformation) to project our representations back into the image plane. We then use the standard MSE Loss between the original pixels and reconstructed pixels for only the 100 occluded patches. We allow self attention as in the following formula for the energy of multiheaded attention

$$E^{\text{ATT}} = \sum_h -\frac{1}{\beta} \sum_C \log \left(\sum_B \exp \left(\beta \sum_\alpha K_{\alpha h B} Q_{\alpha h C} \right) \right) \quad (23)$$

We give details of our architectural choices in Table 7. In the main paper we present our Energy Transformer with a configuration similar to the standard base Transformer configuration (e.g., token dimension 768, 12 heads each with $Y = 64$, softmax’s $\beta = \frac{1}{\sqrt{Y}}, \dots$), with several considerations learned from the qualitative image evaluations:

- The $\frac{dt}{\tau}$ (step size) of 1 implicitly used in the traditional transformer noticeably degrades our ability to smoothly descend the energy function. We find that a step size of 0.1 provides a smoother descent down the energy function and benefits the image reconstruction quality.
- We observe that our MSE loss must include some subset of un-occluded patches in order for the HN to learn meaningful filters.
- Values of β in the energy attention that are too high prevent our model from training. This is possibly due to vanishing gradients in our attention operation from a softmax operation that is too spiky.
- Without gradient clipping, our model fails to train at the learning rates we tried higher than $1e-4$. We observe that gradient clipping not only helps our model train faster at the trainable learning rates, it also allows us to train at higher learning rates.

Our architecture and experiments for the image reconstruction task were written in JAX [45] using Flax [46]. This engineering choice means that our architecture definitions are quite lightweight, as we can define the desired energy function of the ET and use JAX’s autograd to automatically calculate the desired update.

F.1 Exploring the Hopfield Memories

A distinctive aspect of our network is that any variable that has a vector index i of tokens can be mapped into the data domain by applying the decoder network to this variable. This makes it possible to inspect all the weights in the model. For instance, the concept of “memories” is crucial to understanding how Hopfield networks function. The memories within the HN module represent the building blocks of *all possible image patches* in our data domain, where an encoded image patch is a superposition of a subset of memories. The complete set of memory vectors from the HN module is shown in Figure 5. The same analysis can be applied to the weights of the ET-attention module. In Figure 6, we show all the weights from this module mapped into the image plane.

Table 7: Hyperparameter, architecture, and data augmentation choices for ET model during ImageNet-1k masked training experiments. Data augmentations are listed as parameters passed to the equivalent `timm` dataloader functionality.

Training		Architecture		Data Augmentation	
batch_size	768	token_dim	768	random_erase	None
epochs	100	num_heads	12	horizontal_flip	0.5
lr	5e-4	head_dim	64	vertical_flip	0
warmup_epochs	2	β	1/8	color_jitter	0.4
start & end lr	5e-7	train_betas	No	scale	(0.08, 1)
b1, b2 (ADAM)	0.9, 0.99	step size α	0.1	ratio	(3/4, 4/3)
weight_decay	0.05	depth	12	auto_augment	None
grad_clipping	1.	hidden_dim HN	3072		
		bias in HN	None		
		bias in ATT	None		
		bias in LNORM	Yes		

F.2 Positional Bias Correlations

The relationships between our position embeddings exhibit similar behavior to the position correlations of the original ViT in that they are highly susceptible to choices of the hyperparameters (Figure 10 of [13]). In particular, we consider the effect of weight decay and the β parameter that serves as the inverse temperature of the attention operation (see Equation 8). The lower the temperature (i.e., the higher the value of β), the spikier the softmax distribution. By using a lower β , we encourage the attention energy to distribute its positional embeddings across a wider range of patches in the model.

F.3 Contribution of Attention and Hopfield Modules

We additionally visualize the gradients of the energy function (which are equal to the token updates, see Equation 11) of both ATTN block and the HN block, see Figure 8. Early in time, almost all signal to the masked tokens comes from the ATTN block, which routes information from the open patches to the masked ones; no meaningful signal comes from the HN block to the masked patch dynamics. Later in time we observe a different phenomenon: almost all signal to masked tokens comes from the HN module while ATTN contributes a blurry and uninformative signal. Thus, the attention layer is crucial early in the network dynamics, feeding signal to masked patches from the visible patches, whereas the HN is crucial later in the dynamics as the model approaches the final reconstruction, sharpening the masked patches. All the qualitative findings presented in this section are in accord with the core computational strategy of the ET block as it was designed theoretically in Appendix B.

G Parameter comparison

The energy function enforces symmetries in our model, which means ET has fewer parameters than its ViT counterparts. In particular, ET has no “Value Matrix” \mathbf{W}^V in the attention mechanism, and the HN module has only one of the two matrices in the standard MLP of the traditional Transformer block. We report these differences in Table 8. We take the ET configuration used in this paper, which has an architecture fully comparable to the original ViT-base [13] with `patch_size=16`, and report the number of parameters against ViT-base and an “ALBERT” version of ViT [47] where a single ViT block is shared across layers. We saw no benefit when including biases in ET, so we also exclude the biases from the total parameter count in the configuration of ViT and ALBERT-ViT. We report both the total number of parameters and the number of parameters per Transformer block. Furthermore, as a means for a fairer comparison, we reduce the parameters of ViT-Base such that the parameter count is similar to that of ET. We contrast ET and the reduced ViT on the image reconstruction task using PSNR (Peak-Signal-to-Noise-Ratio) and SSIM (Structural-SIMilarity) [48] as the evaluation metrics of the reconstructions. The results are demonstrated in Table 9, where the performance of ET is very close to that of the reduced ViT, which has a slight advantage in parameter count.

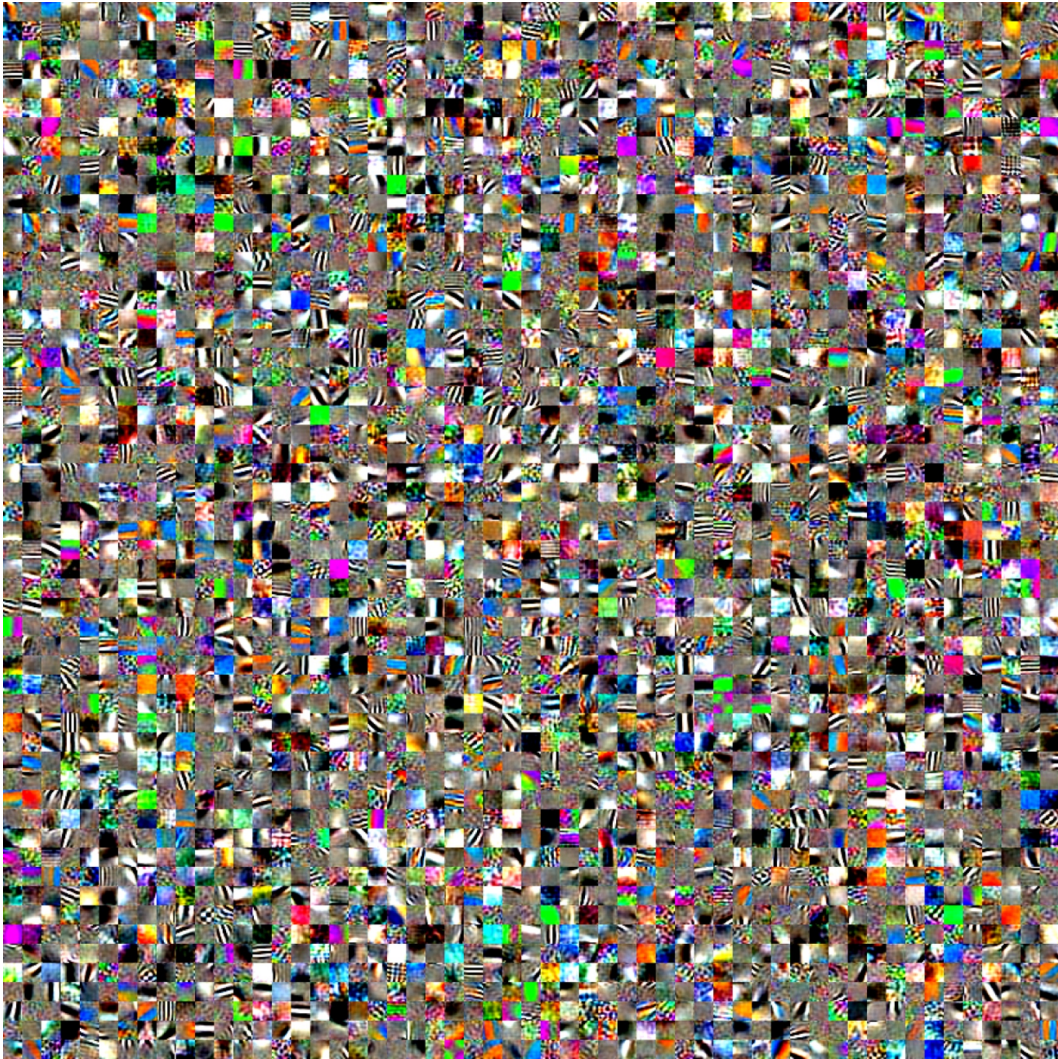


Figure 5: Visualizing a randomly selected 3025 patch memories of the 3072 learned by weight matrix in the Hopfield Network module (HN) of our model. These memories are vectors of the same dimensions D as the patch tokens, stored as rows in the weight matrix ξ . Each image patch is visualized using the model’s trained decoder.

Table 8: Comparison between the number of parameters in a standard ViT, an ALBERT version of ViT where standard Transformer blocks are shared across layers, and our ET. Comparison is done assuming no biases in any operation.

Model	NParams		NParams (per block)	
ViT-Base	86.28M	▼0.00%	7.08M	▼0.00%
ALBERT ViT-Base	8.41M	▼90.25%	7.08M	▼0.00%
ET	4.87M	▼94.36%	3.54M	▼50.02%

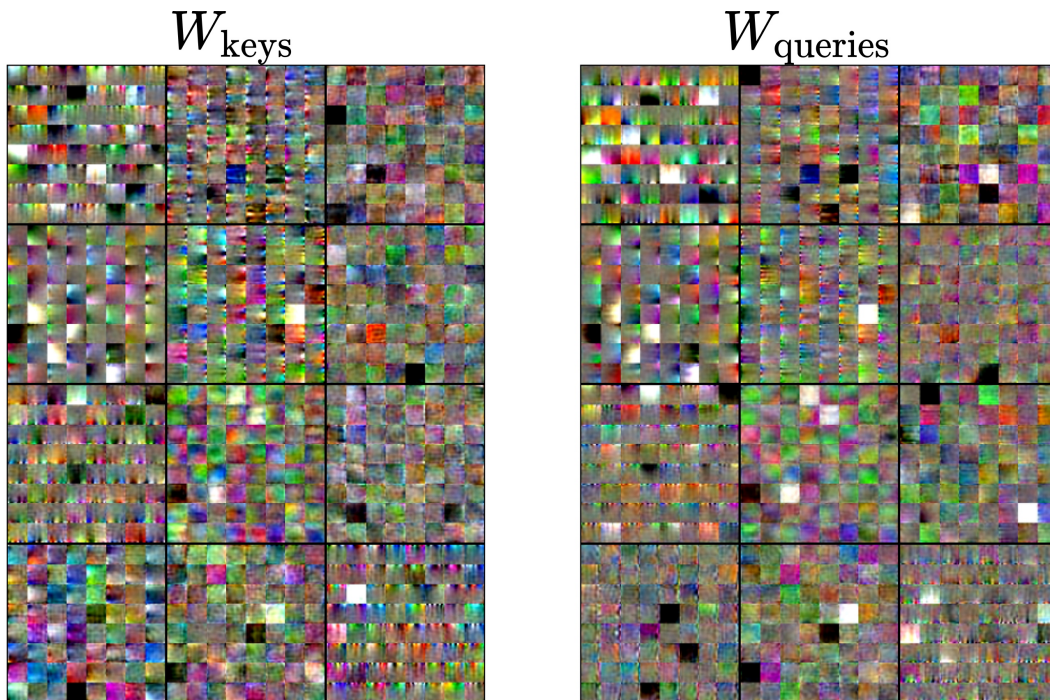


Figure 6: Visualizing the token dimension of the “key” and “query” matrices of the attention as image patches. Each head is represented as a cell on the 4×3 grid above. We use the trained decoder of our model to visualize each weight.

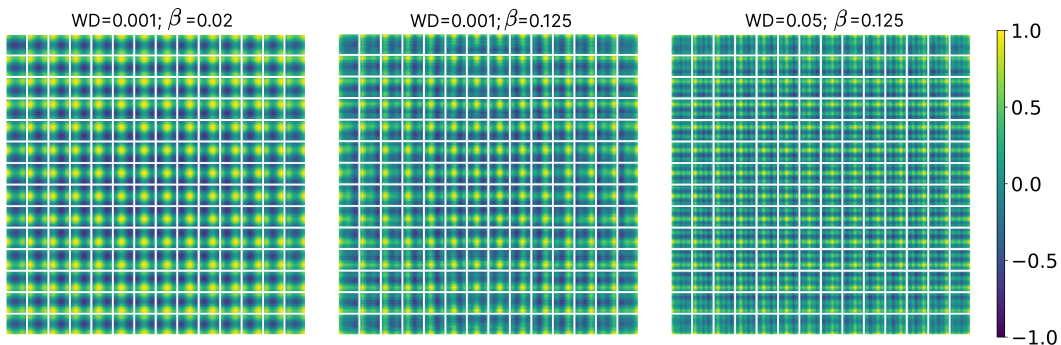


Figure 7: The cosine similarity between position biases of patches when the ET model is trained under different hyperparameter choices for β (inverse temperature of the attention energy) and weight decay. Our ET sees a trend where smoother correlations are observed with smaller β and weight decay.

Table 9: Comparison between ET and ‘comparable-size’ ViT on image reconstruction task. Given ViT-Base, we reduce its parameter count down to a number similar to that of ET for the image domain. The metrics, PSNR (Peak-Signal-to-Noise-Ratio) and SSIM (Structural-SIMilarity), are recorded for the image reconstruction evaluations.

Model	NParams	PSNR	SSIM
ViT	5.52M ∇ 0.00%	22.11 ∇ 0.00	0.715 ∇ 0.00
ET	4.28M ∇ 22.46%	21.51 ∇ 0.59	0.681 ∇ 0.03

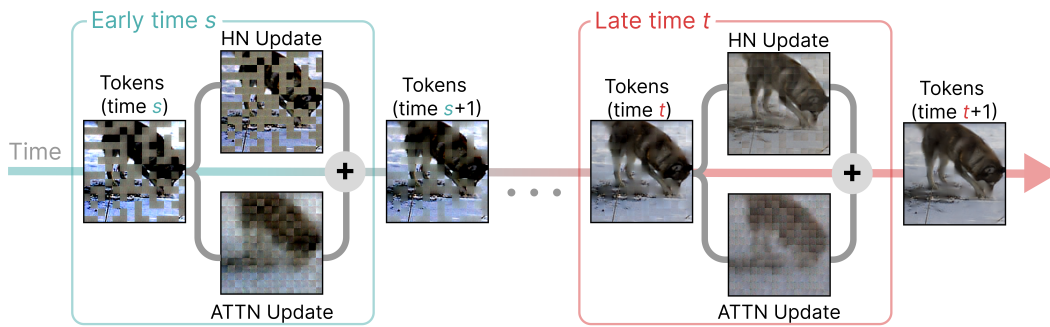


Figure 8: Token representations and gradients are visualized using the decoder at different times during the dynamics. The Energy Attention (ATTN) block contributes general structure information to the masked patches at *earlier* time steps, whereas the Hopfield Network (HN) significantly sharpens the quality of the masked patches at *later* time steps.

H Formal Algorithm for Training and Inference of ET

Algorithm 1: Training and inference pseudocode of ET for image reconstruction task

```

1 HyperParameters
2    $\alpha$ : Energy descent stepsize
3    $\epsilon$ : Learning rate
4    $p$ : Token mask probability
5    $b$ : batch size
6 Parameters
7    $\mathbf{W}^K \in \mathbb{R}^{Y \times H \times D}$ ,  $\mathbf{W}^Q \in \mathbb{R}^{Y \times H \times D}$ : Key, Query kernels of the Energy Attention
8    $\xi \in \mathbb{R}^{M \times D}$ : Kernel of Hopfield Network
9    $\gamma_{\text{norm}} \in \mathbb{R}$ ,  $\delta_{\text{norm}} \in \mathbb{R}^D$ : Scale, bias of LayerNorm
10  MASK  $\in \mathbb{R}^D$ : Mask token
11   $\delta_{\text{pos}} \in \mathbb{R}^{N \times D}$ : Position bias, added to each token
12   $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{P \times D}$ ,  $\delta_{\text{enc}} \in \mathbb{R}^D$ : Kernel, bias of affine Encoder
13   $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{D \times P}$ ,  $\delta_{\text{dec}} \in \mathbb{R}^D$ : Kernel, bias of affine Decoder
14 Infer
15 Inputs
16   Corrupted image tokens  $\tilde{X} \in \mathbb{R}^{N \times D}$ 
17   Add position biases:  $\tilde{X} \leftarrow \tilde{X} + \delta_{\text{pos}}$ ;
18   for timesteps  $t = 1, \dots, T$  do
19     Normalize each token:
20      $\tilde{g} \leftarrow \text{LayerNorm}(\tilde{X}; \gamma_{\text{norm}}, \delta_{\text{norm}})$ ;  $\tilde{g} \in \mathbb{R}^{N \times D}$ 
21     Calculate Energy of tokens:
22      $E \leftarrow \text{EnergyTransformer}(\tilde{g}; \mathbf{W}^K, \mathbf{W}^Q, \xi)$ ;  $E \in \mathbb{R}$ 
23      $\tilde{X} \leftarrow \tilde{X} - \alpha \nabla_{\tilde{g}} E$ ;
24   return  $\tilde{X}$ 
25 Train
26 Inputs
27   Dataset  $S_{\text{train}}$  with elements  $X \in \mathbb{R}^{\text{channels} \times \text{height} \times \text{width}}$ 
28 Initialize
29   Randomly initialize from  $\mathcal{N}(0, 0.02)$ :
30    $\mathbf{W}^K, \mathbf{W}^Q, \xi, \text{MASK}, \mathbf{W}_{\text{enc}}, \mathbf{W}_{\text{dec}}, \delta_{\text{pos}} \sim \mathcal{N}(0, 0.02)$ 
31   Set other biases to zero:  $\delta_{\text{enc}}, \delta_{\text{dec}}, \delta_{\text{norm}} \leftarrow 0$ 
32   Set LayerNorm scale to one:  $\gamma_{\text{norm}} \leftarrow 1$ 
33 for epoch  $n = 1, \dots, N_{\text{epoch}}$  do
34    $S_{\text{epoch}} \leftarrow S_{\text{train}}$ 
35   for batch  $B \subset S_{\text{epoch}}$   $B \in \mathbb{R}^{b \times \text{channels} \times \text{height} \times \text{width}}$ 
36     do
37       Convert image into non-overlapping patches:
38        $B_{\text{patch}} \leftarrow \text{Patchify}(B)$ ;  $B_{\text{patch}} \in \mathbb{R}^{b \times N \times P}$ 
39       Embed image patches into tokens:
40        $X \leftarrow \text{Encode}(B_{\text{patch}}; \mathbf{W}_{\text{enc}}, \delta_{\text{enc}})$ ;  $X \in \mathbb{R}^{b \times N \times D}$ 
41       Replace image tokens randomly by MASK:
42        $\tilde{X}, I_{\text{mask}} \leftarrow \text{Mask}(X; \text{MASK}, p)$   $\tilde{X} \in \mathbb{R}^{b \times N \times D}, I_{\text{mask}} \in \{0, 1\}^{b \times N}$ 
43       Reconstruct tokens with ET:
44        $\tilde{X} \leftarrow \text{Infer}(\tilde{X})$ 
45       Decode tokens:
46        $\hat{B}_{\text{patch}} \leftarrow \text{Decode}(\tilde{X}[I_{\text{mask}}]; \mathbf{W}_{\text{dec}}, \delta_{\text{dec}})$ ;  $\hat{B}_{\text{patch}} \in \mathbb{R}^{b \times N \times P}$ 
47       Calculate MSE loss on corrupted tokens:
48        $L \leftarrow \text{Mean}(|\hat{B}_{\text{patch}}[I_{\text{mask}}] - B_{\text{patch}}[I_{\text{mask}}]|^2)$   $L \in \mathbb{R}$ 
49        $\text{params} \leftarrow \text{params} - \epsilon \nabla_{\text{params}} L$ 
50        $S_{\text{epoch}} \leftarrow S_{\text{epoch}} \setminus B$ 
51 return  $\text{params}$ 

```
