

Energy Transformer

Benjamin Hoover* Yuchen Liang* Bao Pham* Rameswar Panda Hendrik Strobelt Polo Chau Mohammed J. Zaki Dmitry Krotov

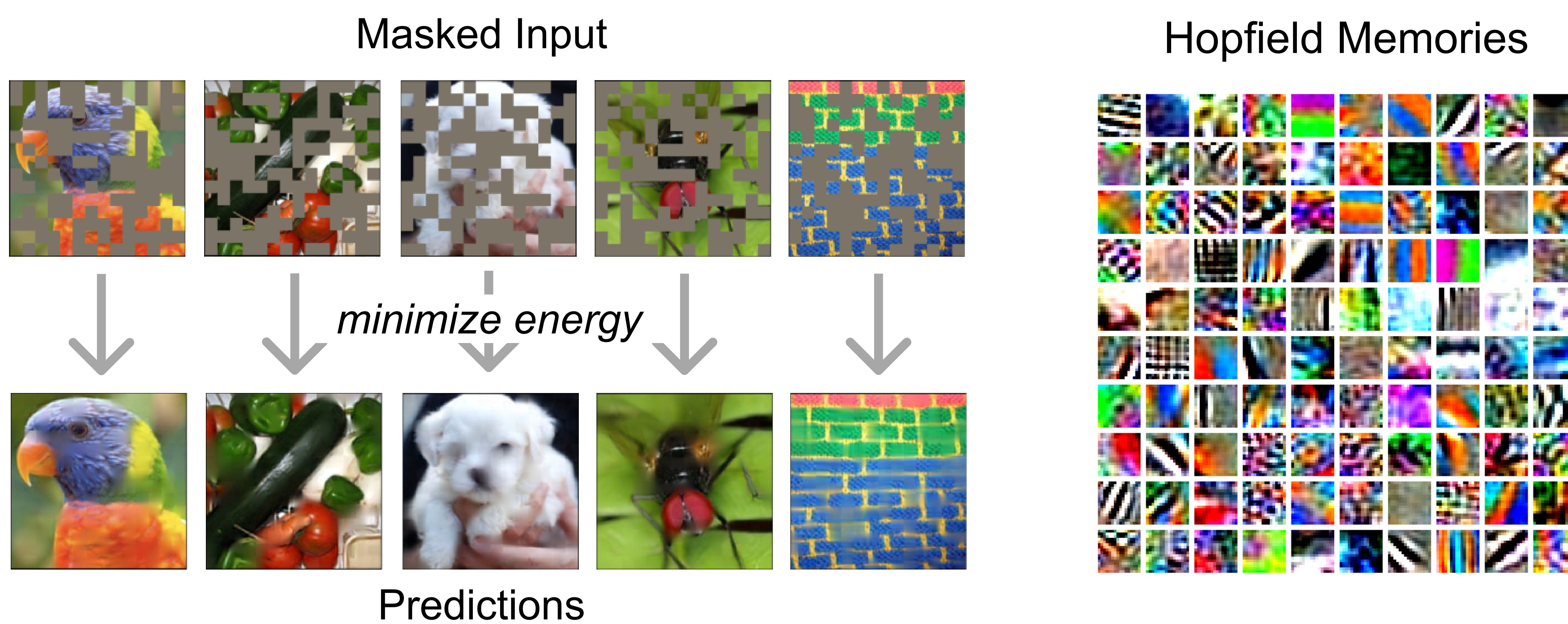
*Authors contributed equally

Energy Transformer (**E.T.**) is a novel architecture that is an **Energy-Based Model**, an **Associative Memory**, and a **Transformer**. Specifically, **E.T.** looks like a recurrent Transformer block that defines an attractor system (an *O.D.E. with guaranteed fixed points convergence*) that performs error correction via energy descent. **E.T.** shows excellent performance on MASKed image in-painting, graph anomaly detection, and graph classification.



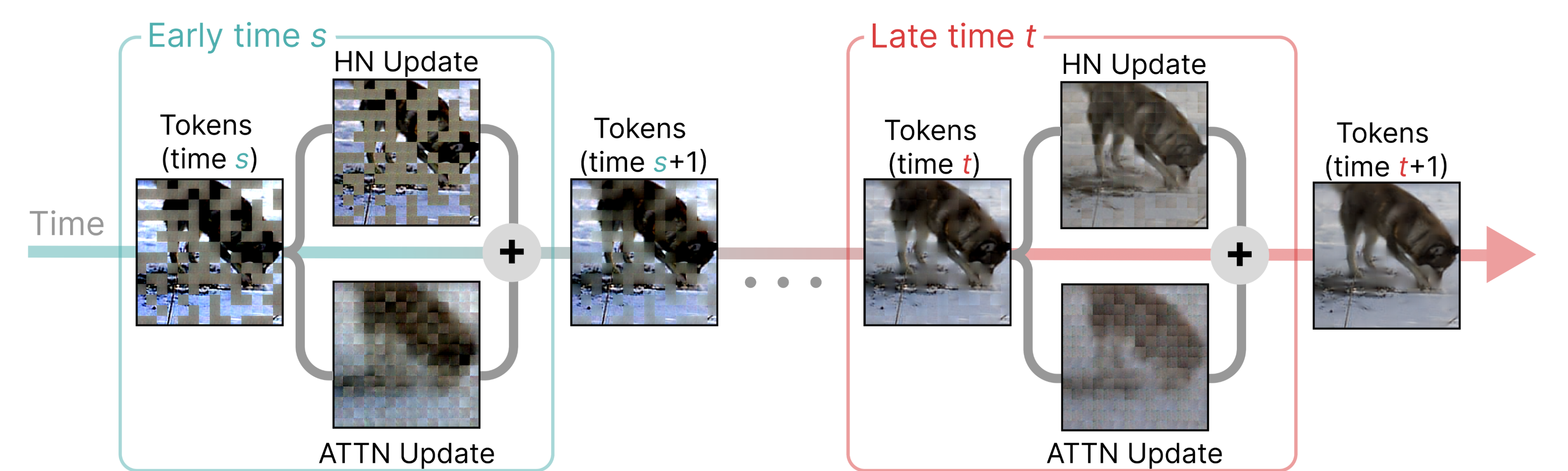
E.T. is Energy-Based Associative Memory

Minimizing energy (inference) performs error correction on corrupted data.



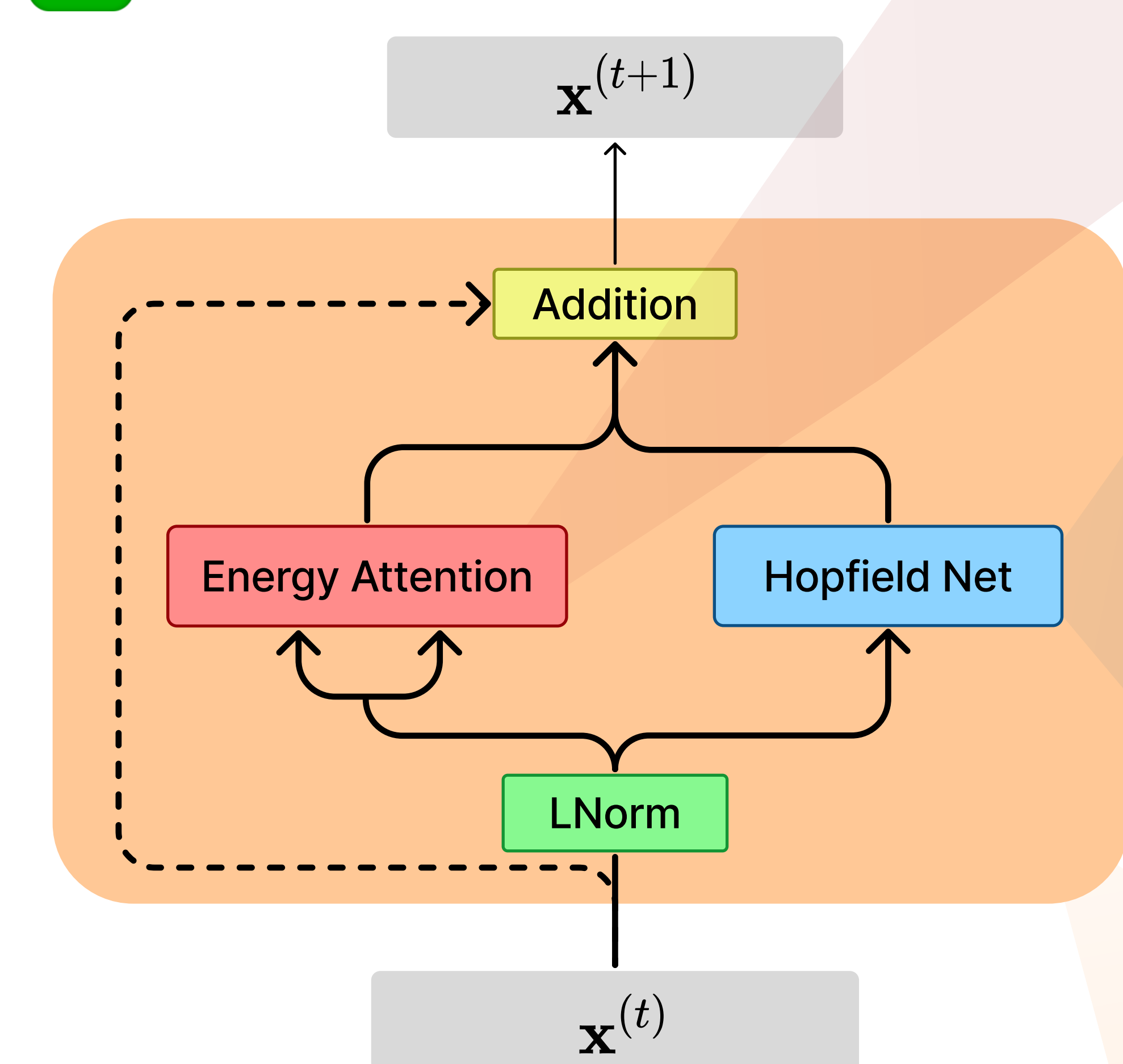
E.T. is Interpretable by Design

Visualize dynamic states (tokens and updates) at any time



E.T. is a Transformer

Energy Transformer Block



Energy Attention makes Queries align with Keys in the latent space (and vice versa!)

Hopfield Network makes tokens look like memories

Energy Transformer minimizes the sum of these energies

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

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

$$E^{\text{HN}} = -\sum_{B=1}^N \sum_{\mu=1}^K G \left(\sum_{j=1}^D \xi_{\mu j} g_{jB} \right)$$

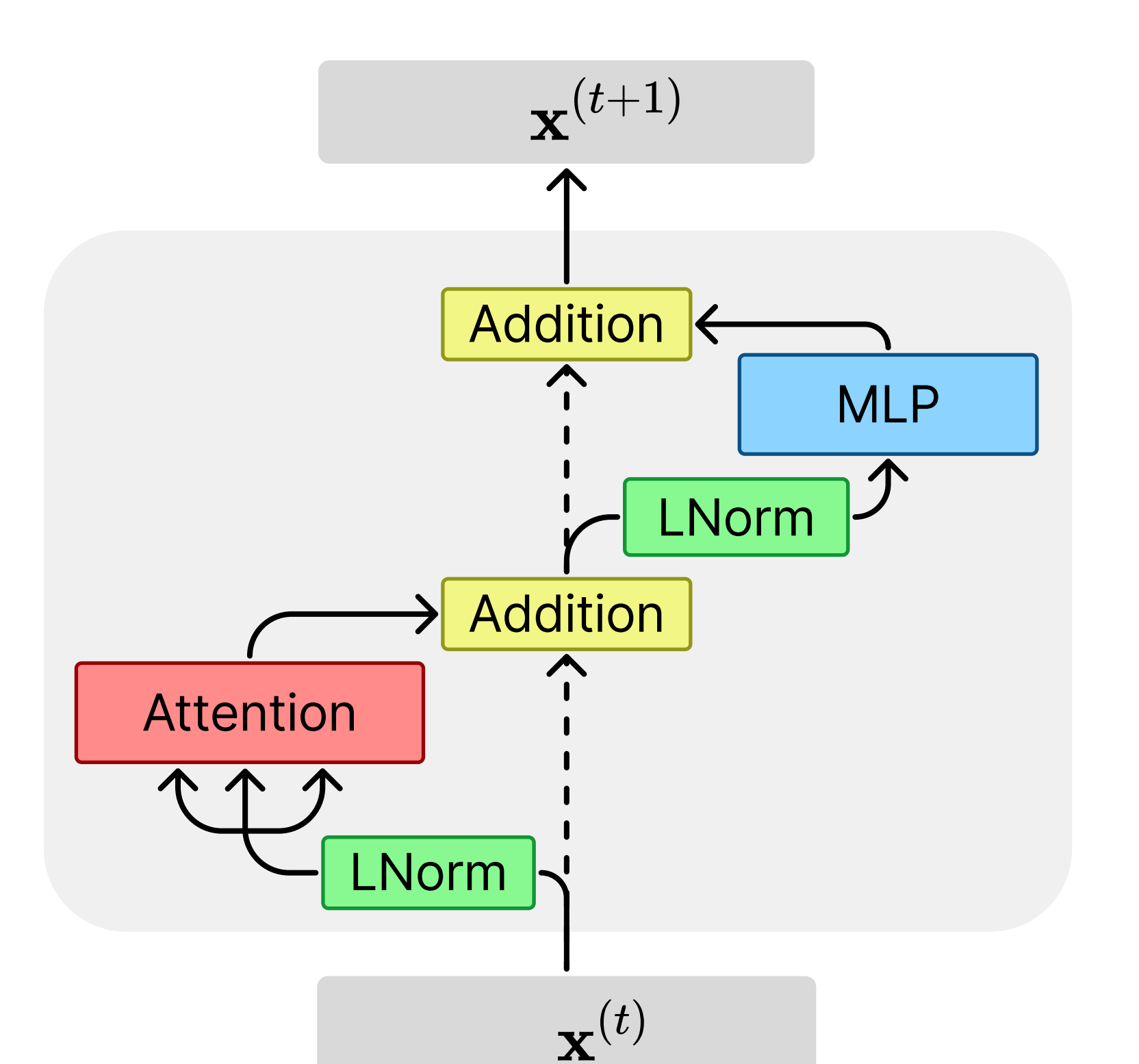
$$\text{Update: } -\frac{\partial E^{\text{HN}}}{\partial g_{iA}} = \sum_{\mu=1}^K \xi_{\mu i} G' \left(\sum_{j=1}^D \xi_{\mu j} g_{jA} \right)$$

$$E = E^{\text{ATT}} + E^{\text{HN}}$$

$$\text{Continuous update: } \tau \frac{dx_{iA}}{dt} = -\frac{\partial E}{\partial g_{iA}}$$

$$\text{Discrete update: (stepsize } \alpha) \quad x_{iA}^{(t+1)} = x_{iA}^{(t)} - \alpha \left(\frac{\partial E^{\text{ATT}}}{\partial g_{iA}} + \frac{\partial E^{\text{HN}}}{\partial g_{iA}} \right)$$

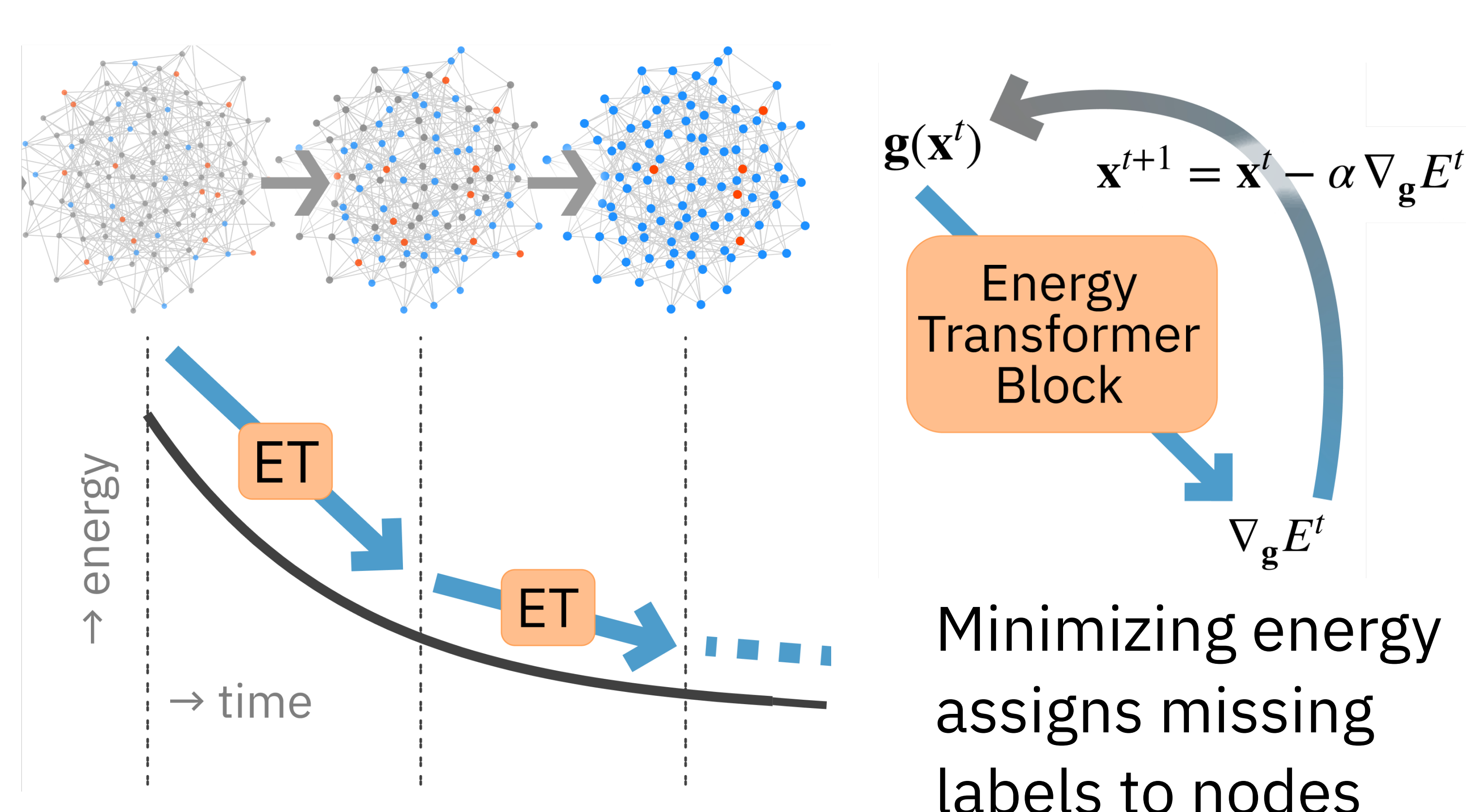
Standard Transformer Block



E.T. sets SOTA on Graphs

Node anomaly detection

| | Datasets | Split | Top Baseline Score | ET (Ours) |
|----------|-----------|----------------|---------------------------------|---------------------------------|
| Macro-F1 | Yelp | 1% | 62.1 \pm 1.3 | 63.0 \pm 0.6 \uparrow 0.9 |
| | | 40% | 71.0 \pm 0.9 | 71.5 \pm 0.1 \uparrow 0.5 |
| | Amazon | 1% | 90.9 \pm 0.7 | 89.3 \pm 0.7 \downarrow 0.7 |
| | | 40% | 92.2 \pm 0.4 | 92.8 \pm 0.3 \uparrow 0.6 |
| | T-Finance | 1% | 84.8 \pm 0.0 | 85.1 \pm 1.0 \uparrow 0.3 |
| | | 40% | 86.8 \pm 0.0 | 88.2 \pm 1.0 \uparrow 1.4 |
| T-Social | 1% | 75.9 \pm 0.0 | 79.1 \pm 0.7 \uparrow 3.2 | |
| | 40% | 83.9 \pm 0.0 | 83.5 \pm 0.4 \downarrow 0.4 | |
| AUC | Yelp | 1% | 75.4 \pm 0.9 | 73.2 \pm 0.8 \downarrow 2.2 |
| | | 40% | 84.0 \pm 0.9 | 84.9 \pm 0.3 \uparrow 0.9 |
| | Amazon | 1% | 90.4 \pm 2.0 | 91.9 \pm 1.0 \uparrow 1.5 |
| | | 40% | 98.0 \pm 0.4 | 97.3 \pm 0.4 \downarrow 0.7 |
| | T-Finance | 1% | 91.1 \pm 0.0 | 92.8 \pm 1.1 \uparrow 1.7 |
| | | 40% | 94.3 \pm 0.0 | 95.0 \pm 3.0 \uparrow 0.7 |
| T-Social | 1% | 88.0 \pm 0.0 | 91.9 \pm 0.6 \uparrow 3.9 | |
| | 40% | 95.2 \pm 0.0 | 93.9 \pm 0.2 \downarrow 1.3 | |



Graph Classification

| Datasets | Top Baseline Score | ET (Ours) |
|--------------|--------------------|---------------------------------|
| PROTEINS | 84.9 \pm 1.6 | 90.3 \pm 5.4 \uparrow 5.4 |
| NCI1 | 87.5 \pm 0.5 | 90.1 \pm 0.1 \uparrow 2.6 |
| NCI109 | 87.4 \pm 0.3 | 90.5 \pm 0.1 \uparrow 3.1 |
| DD | 95.7 \pm 1.9 | 95.9 \pm 0.8 \uparrow 0.2 |
| ENZYMES | 78.4 \pm 0.6 | 99.8 \pm 0.0 \uparrow 21.4 |
| MUTAG | 100.0 \pm 0.0 | 96.6 \pm 0.2 \downarrow 3.4 |
| MUTAGENICITY | 82.2 \pm 0.6 | 98.7 \pm 0.1 \uparrow 16.5 |
| FRANKENSTEIN | 78.9 \pm 0.3 | 99.8 \pm 0.1 \uparrow 20.9 |