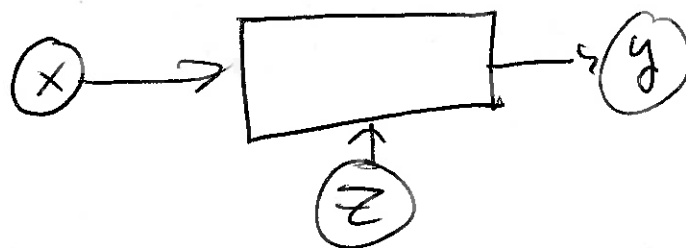


- Attention
 - Code example
 - Some more examples of Hamiltonian systems
-

We point to ENM 531 and Murphy's textbook for discussion about architectures, but for final projects I want to make sure everyone has access to a modern transformer architecture

Consider the conditional neural field



$$y = f(x | z; \theta)$$

Idea z modulates the x, y input/output relationship

ex

y - finite difference stencil
 x - grid function
 z - material properties / microscopy

Many options : - deep onet (Lu Lu)
- PCA / FNO (Stewart)
⊛ - cross Attention

Attention as soft dictionary lookup (see Murphy 15.4.1)

- Consider $(k_i, v_i)_{i=1}^N$ key-value pairs
describing input/output labels
- Consider a query q where model is
evaluated

def
$$\text{Attn}(q, \{k_i, v_i\}) = \sum_{i=1}^N \alpha_i(q, \vec{k}) v_i$$

for attn weights

$$0 \leq \alpha_i \leq 1$$

$$\sum \alpha_i = 1$$

introducing an attention score $a(q, k_i)$

$$\alpha_i = \text{softmax}(a(q, \vec{k})) = \frac{\exp[a(q, k_i)]}{\sum_j \exp[a(q, k_j)]}$$

Remark This is a data driven basis

e.g. Taking $\alpha_i(z, \vec{K}) = \phi_i(z)$

$\phi_i \in \text{CPWL}$

satisfies the same properties

$$f(x) = \sum_i \alpha_i(x, x_n) y_i$$

$$z \rightarrow x$$

$$K \rightarrow x_n \quad (\text{nodal values})$$

$$y_i \rightarrow \text{nodal basis functions}$$

Self-Attention

$$z, K, v = x$$

Multi-head Attention

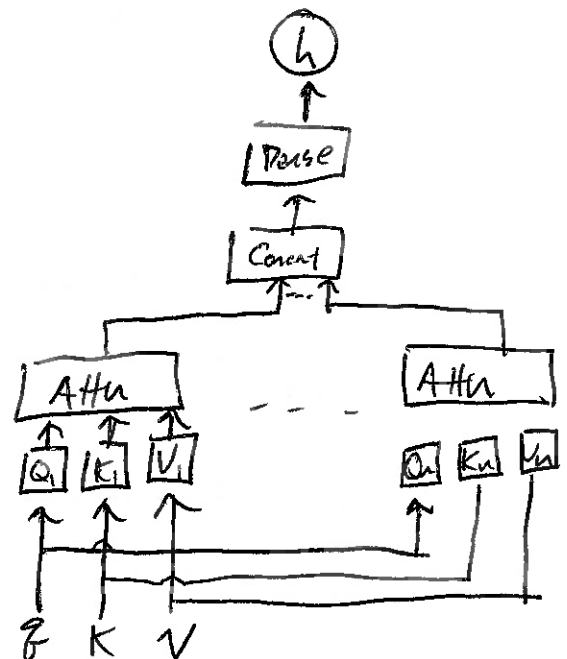
- Introduce Dense MLP Blocks $i=1, \dots, M$

$$\hat{z}_i = Q_i(z)$$

$$\hat{K}_i = K_i(K)$$

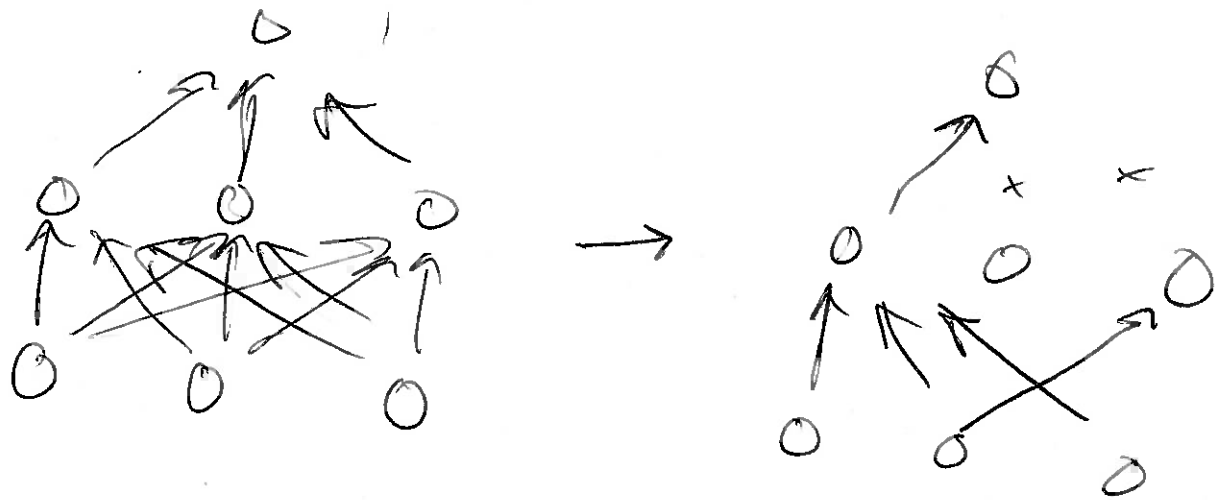
$$\hat{v}_i = V_i(v)$$

$$h = \text{MHA}(z, \{K, v\})$$



Dropout

With probability p , at each forward pass the output of a given neuron will be dropped



Replace weights

$$\theta_{lji} = w_{lji} \epsilon_{li}$$

$$\epsilon_{li} \sim \text{Ber}(1-p)$$

Graph Attention

Steps

1. Input features

~~✗~~

Each node $n_i \in N$ has feature vector $h_i^n \in \mathbb{R}^F$

2. Linear Transform

$$h_i' = W_i h_i^n, \quad W_i \in \mathbb{R}^{F \times F}$$

3. Preattention coeffs

$$e_{ij} = \sigma \left(\underset{\substack{\uparrow \\ \text{leaky ReLU}}}{a} \cdot \underset{\substack{\uparrow \\ \text{trainable}}}{[} W h_i' \parallel W h_j'] \underset{\substack{\uparrow \\ \text{concat}}}{]} \right)$$

4. Attention mech.

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in N} \exp(e_{ik})}$$

5. Aggregate

$$h^{n+1} = \sigma \left(\sum_{j \in N} \alpha_{ij} W h_j^n \right)$$

Physics-Inspired Arch.

GRAND

$$h_i^{n+1} = h_i^n + \sigma\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} (h_j^n - h_i^n)\right)$$

or in graph calculus

$$h_i^{n+1} = h_i^n + \sigma(\delta \star \delta h^n)$$

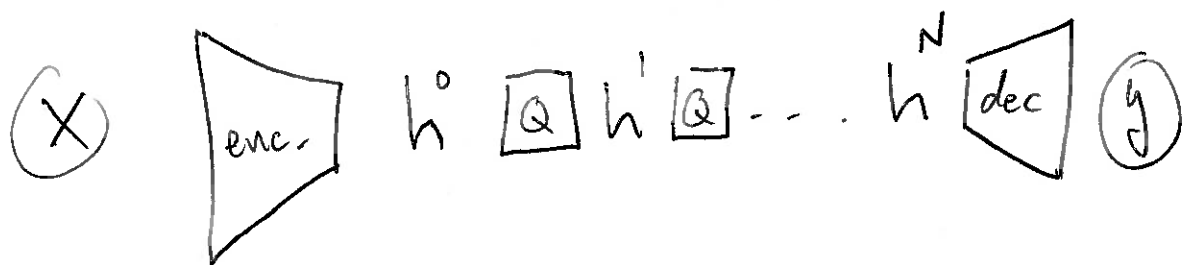
$$\langle \delta, \delta \star g \rangle$$

$$\langle \delta \delta, g \rangle_\alpha$$

view attention
as learning
inner-product

Writing $h^{n+1} = Q^n h^n$

Arbitrary graph data $\mathcal{D} = \{X_i, y_i\}$



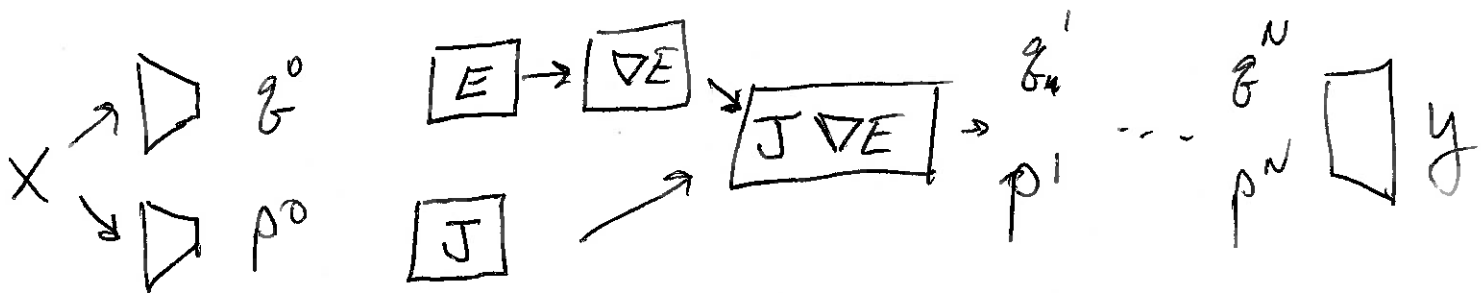
Let's revisit Hamiltonian mech.

$$\frac{d}{dt} \begin{pmatrix} q \\ p \end{pmatrix} = \begin{pmatrix} 0 & -S_0^* \\ S_0 & 0 \end{pmatrix} \begin{pmatrix} \partial_q E \\ \partial_p E \end{pmatrix}$$

$$E = NN(p, q)$$

Then $\dot{E} = 0$

Same as GRAND



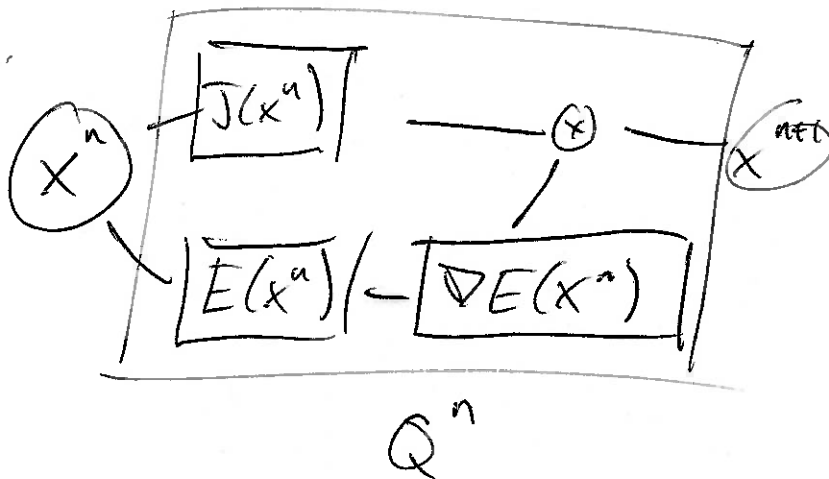
Double Bracket dynamics

$$\frac{dx}{dt} = J \nabla E + J^T J \nabla E$$

$$\frac{dE}{dt} = \frac{\partial E}{\partial x} \cdot \frac{dx}{dt}$$

$$= \frac{\partial E}{\partial x} \cancel{J} \frac{\partial E}{\partial x} + \underbrace{\frac{\partial E}{\partial x} J^T J \frac{\partial E}{\partial x}}_{\geq 0}$$

$= 0$



$$x \rightarrow h^0 \square Q^0 h^1 \square Q^1 \dots h^n \square y$$