

Today:

Graph Neural Nets

Architecture Order In Class:

MLPs \rightarrow CNNs \rightarrow Graph NN \rightarrow RNN/State-space \rightarrow Transformers

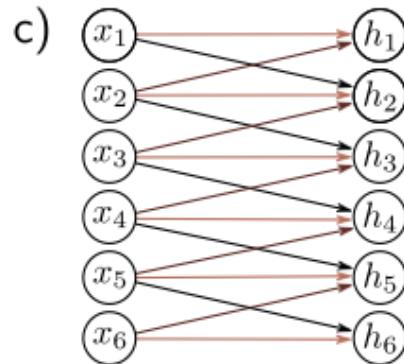
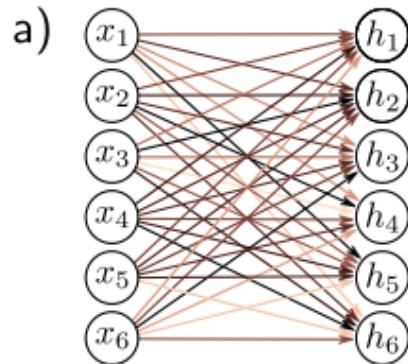
Reading: Prince through Ch 11 + Ch 13

Other Resources: distill.pub/2021/gnn-intro

Packages: dgl.ai, PyG, Spektral

More info: STANFORD CS 224W

Fri TA OH S31CORY 2-3 PM



b)

	x_1	x_2	x_3	x_4	x_5	x_6
h_1	Dark Brown	Medium Brown	Light Brown	Dark Brown	Medium Brown	Light Brown
h_2	Medium Brown	Black	Medium Brown	Dark Brown	Medium Brown	Light Brown
h_3	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Black
h_4	Light Brown	Black	Medium Brown	Dark Brown	Medium Brown	Light Brown
h_5	Black	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Light Brown
h_6	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Dark Brown

d)

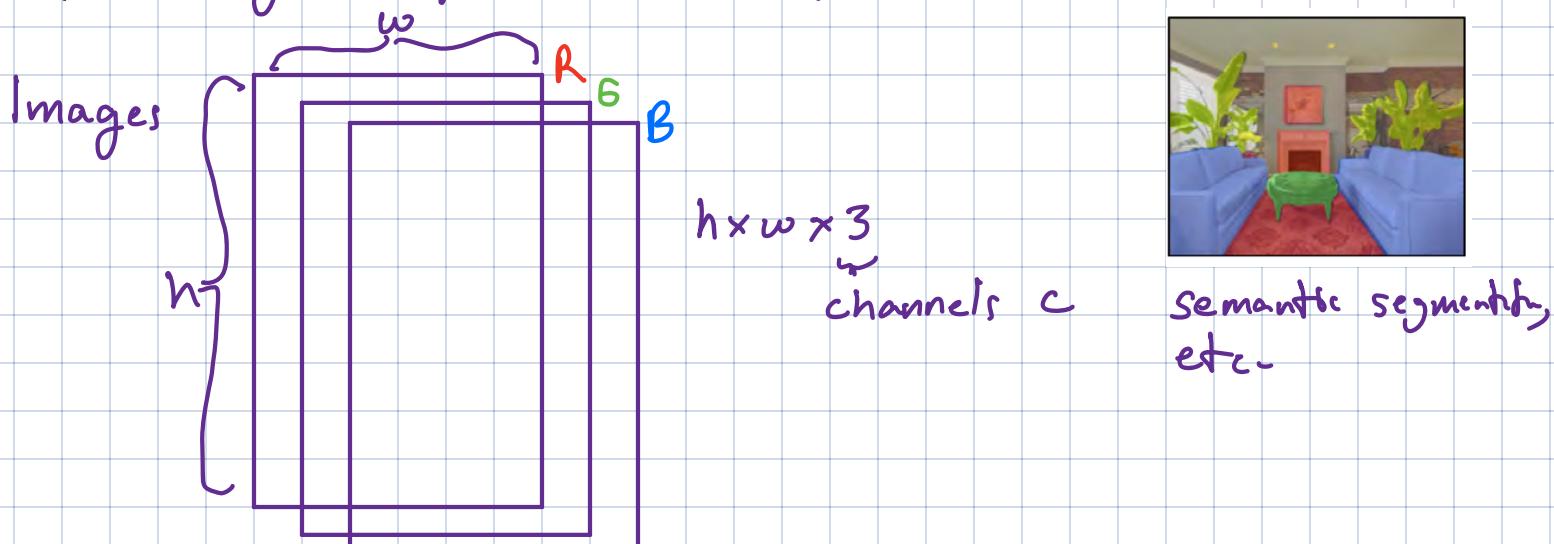
	x_1	x_2	x_3	x_4	x_5	x_6
h_1	Dark Brown	Medium Brown	Light Brown	Dark Brown	Medium Brown	Light Brown
h_2	Medium Brown	Black	Medium Brown	Dark Brown	Medium Brown	Light Brown
h_3	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Black
h_4	Light Brown	Black	Medium Brown	Dark Brown	Medium Brown	Light Brown
h_5	Black	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Light Brown
h_6	Light Brown	Medium Brown	Dark Brown	Dark Brown	Medium Brown	Dark Brown

MLP

CNN

Fig 10.4
in Prince

Inspired by computer vision problems: classification,



Recall Concepts/Ideas from CNN module:

Weight-sharing to respect symmetries/invariances in domain

Convolutions - local processing

Combined with additional channels and depth/nonlinearities for expressivity

Data Augmentations

Dropout as an augmentation for MLP-like Elements in hidden layers

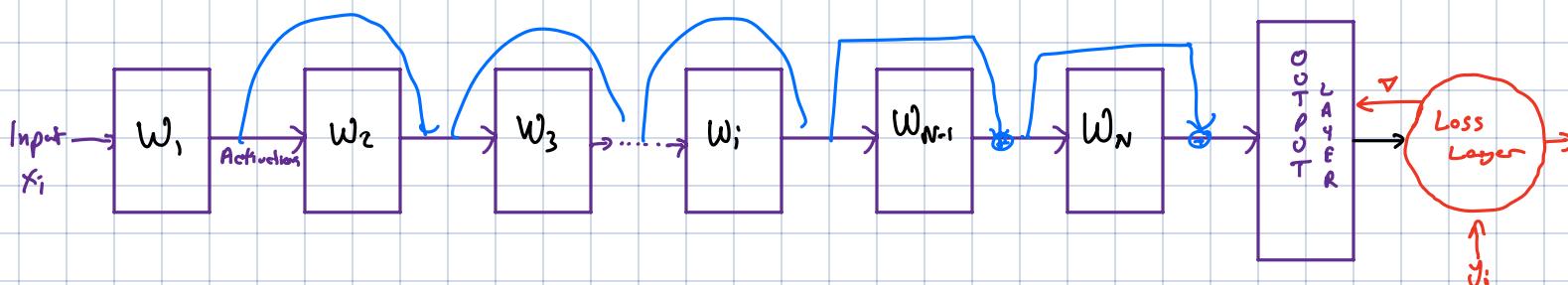
Normalizations for avoiding explosions in activations or gradients

Residual Connections for successful gradient flow to all layers

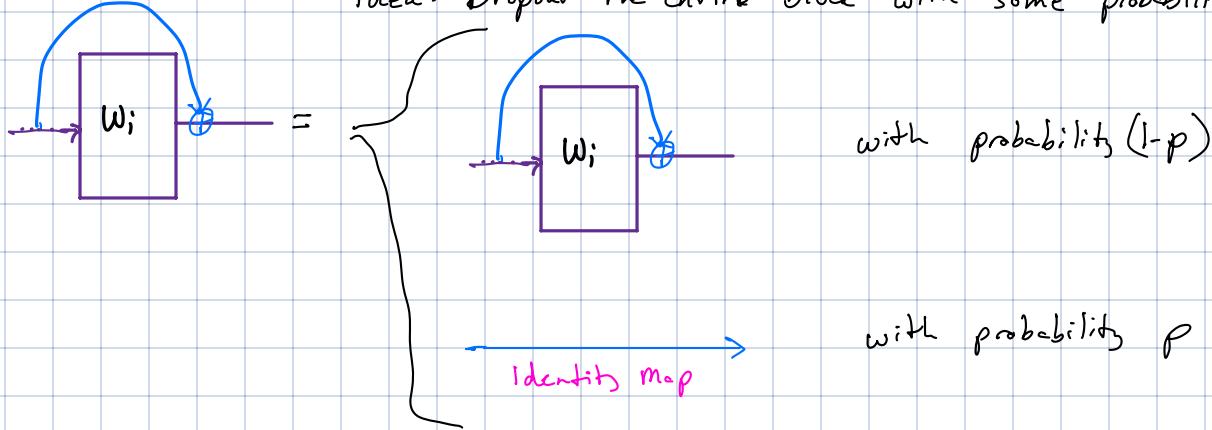
Downsampling/Pooling to get the big picture faster

Upsampling and U-Nets to restore appropriate detail at original scale

One loose end: stochastic depth regularization



Idea: Dropout the entire block with some probability



Like Dropout, helps with overfitting and creating functional redundancy (across blocks)

Note: Unlike Dropout, don't need to adjust anything at inference/test time. ← Added in O.H. (Because of Identity Map & Presence of Normalizations inside blocks)

Generalizing CNNs: Images \rightarrow General Graphs

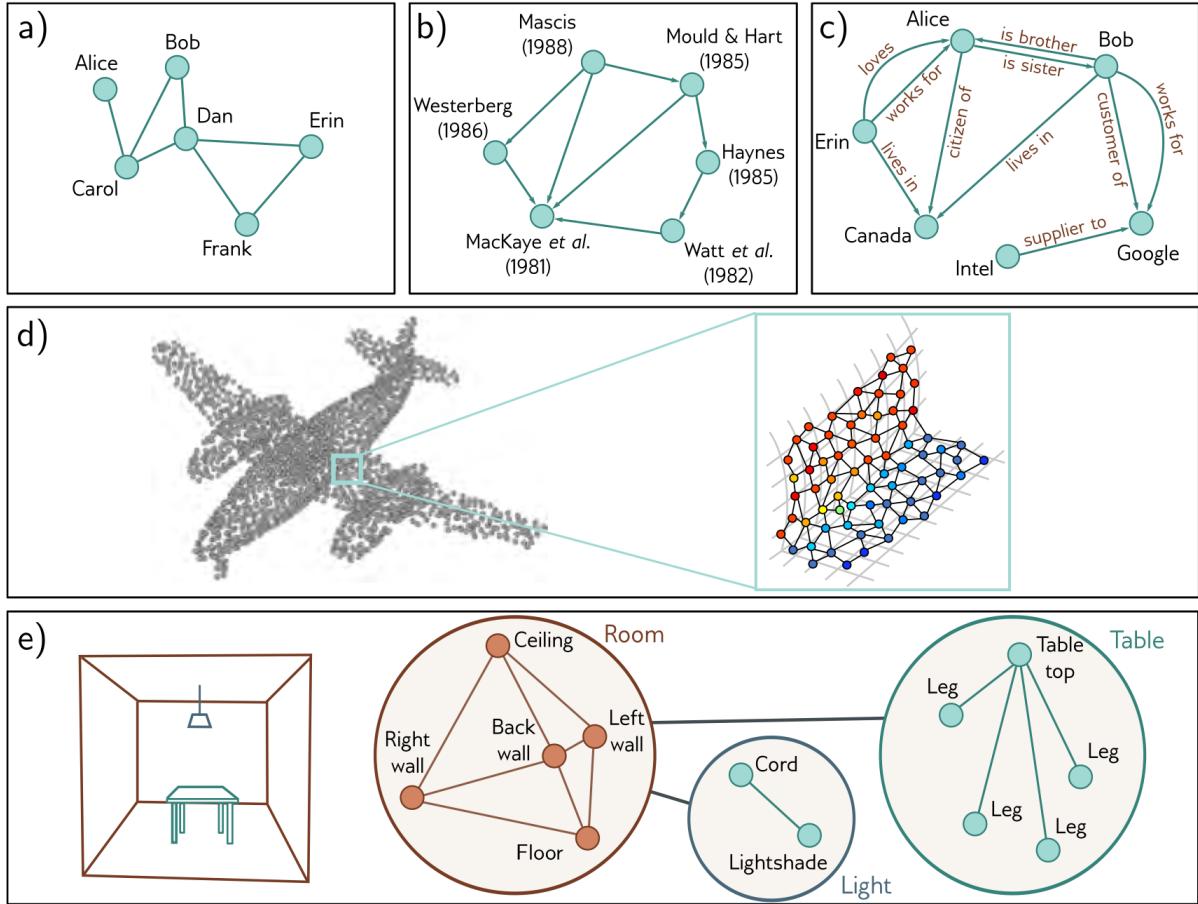
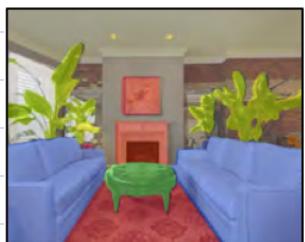


Fig 13.2
in Prince

Image Classification



Semantic Segmentation

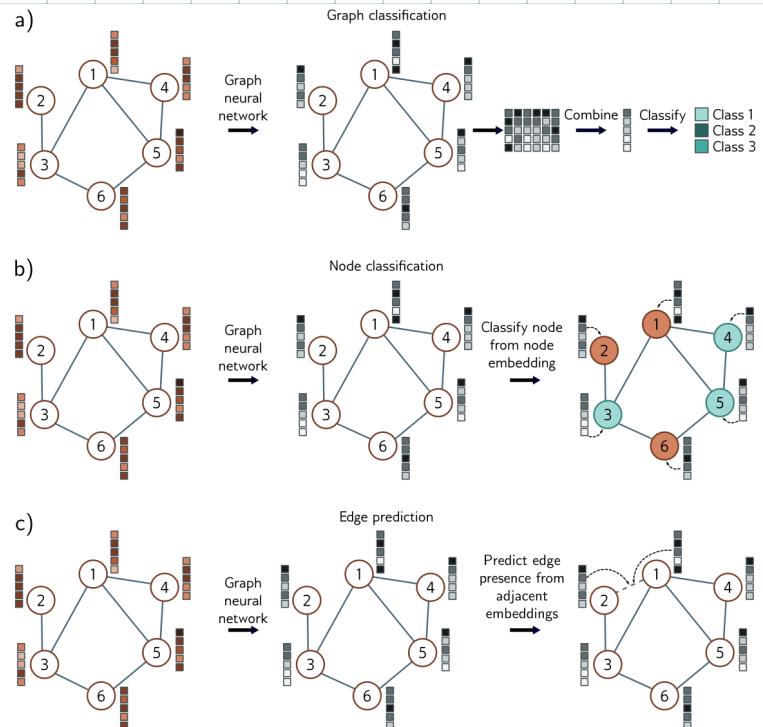
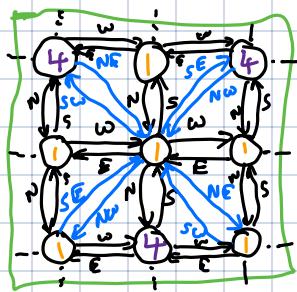


Fig 13.6
in Prince

First Conceptual Step: Image as a graph

Nodes = Pixels

1	1	1	1	1	1	1	1
1	1	4	1	4	1	1	1
1	1	1	1	1	1	1	1
1	1	1	4	1	1	1	1
1	4	1	1	4	1	1	1
1	1	4	4	4	1	1	1
1	1	1	1	1	1	1	1



Directed Graph

Edges labeled with one of
 $\{N, S, E, W, NE, SE, NW, SW\}$

(Always a special relation to self)

Recall Picture:

3x3 examples:

W_{NW}	W_N	W_{NE}
W_W	W_{me}	W_E
W_{SW}	W_S	W_{SE}

$\underbrace{\text{Cin}}_{\text{out}} \{ \overrightarrow{W_i}, \overrightarrow{b} \} \underbrace{\text{Cout}}$

$$\overrightarrow{h}_{out} = \overrightarrow{b} + \sum_i \overrightarrow{W_i} \overrightarrow{h}_{in,i}$$

A 3x3 Conv can be viewed as receiving messages from neighbors, each of whom has their own labeled edge to me, multiplying by an edge-label-specific weight matrix, and adding those up together with a special self-specific matrix acting on this node's info and a bias.

Note: Implicit Zero-padding and Same-Pad
 Weight-sharing across nodes
 1×1 Convs are local to a node

Matrices W_i and biases \overrightarrow{b}
 are specific to layer of processing

Taking Stock: We have, in principle, an immediate recipe to generalize CNNs to directed graphs with the same property: Every edge labeled from a discrete finite set.

What generalizes immediately: Weight-sharing (Nodes do the same things)

3x3 convs (As above, Matrices for edge-labels)

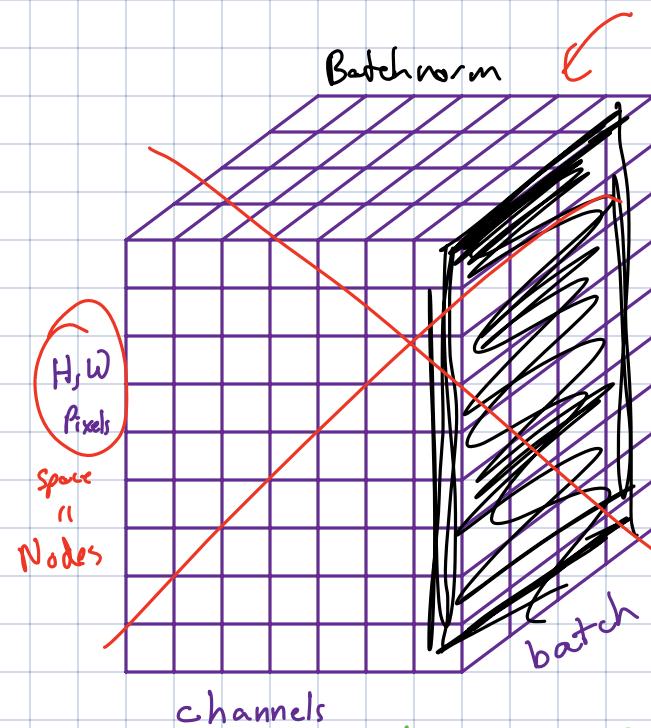
Additional Channels, Residual Connections, Dropout for MLPs, Stochastic RNN,

What does not generalize: Pooling / Downsampling and Upsampling

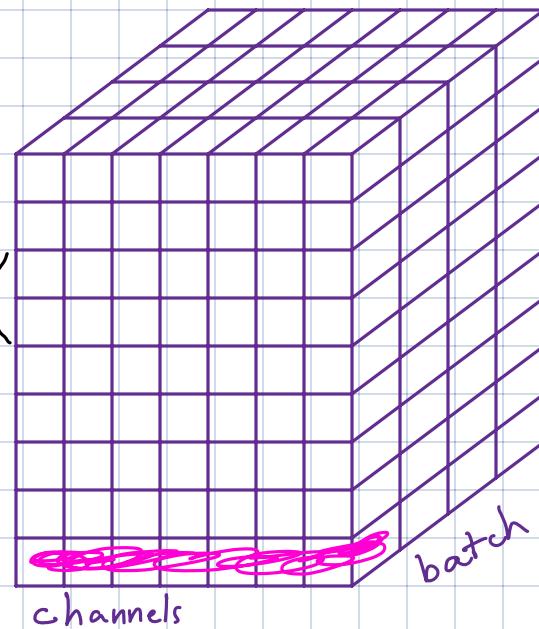
What needs discussion: Normalization.

Core Insight: The graph plays the role of the empty grid of pixels in CNNs.
 The input is a graph with vectors in it

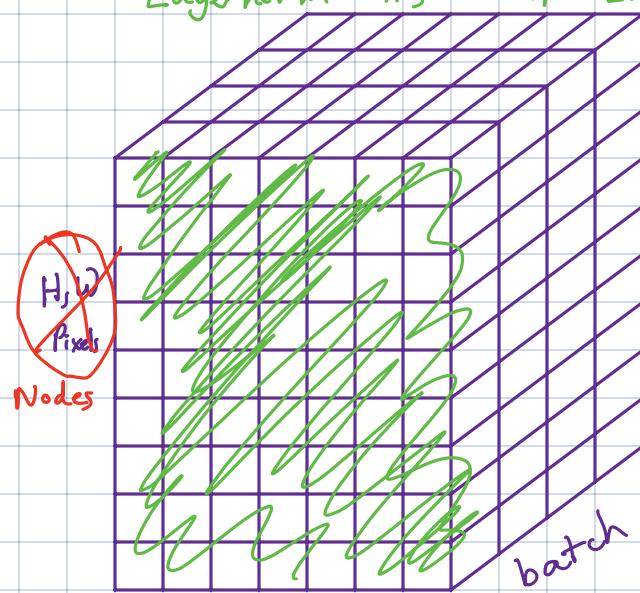
Normalization Layers in GNNs



Picture inaccurate
because
different
graphs
in
Batch
can
have
diff
node
count.



Easy choice: Just per-embedding normalization
Same applied at each node
Like weight-share any learned γ, β
Scale ↑ Shift ↑



Can do this

So are we done? No.	<u>Edge Labels (finite)</u>	Edge Labels (infinite)	None
Kinds of Graphs:	Directed	✓	
	Undirected		???

First Attempt: Reduce to already understood case.

- 1) Treat each undirected edge as both directions \rightarrow replaces —
- 2) Pick a single dummy label for every edge.

What happens? Conv: $\vec{b} + \sum w_i \vec{h}_i$
 $= \vec{b} + \vec{w} \sum \vec{h}_i$

Step back: What do we know?

- a) Me — my own node information is special h_0
- b) How many neighbors I have n
- c) A set $\{h_1, h_2, \dots, h_n\}$ of info from neighbors

Architectural Constraints: I can use n and h_0 specially
 But must be permutation-invariant to h_1, \dots, h_n

So for each output channel, I can compute:

e.g. $f_{W_1}(\vec{h}_0, \sum_i \vec{g}_{W_2}(\vec{h}_i))$ with learnable parameters W_1, W_2

Other choices for Σ :

\prod multiplication

max, min, median (General Order Statistics) + mode

Norm (right kind)

Average, Variance

Count Unique

Correlation with self

Least Squares \leftarrow Use F-matrix, thresh H's down.

Softmax

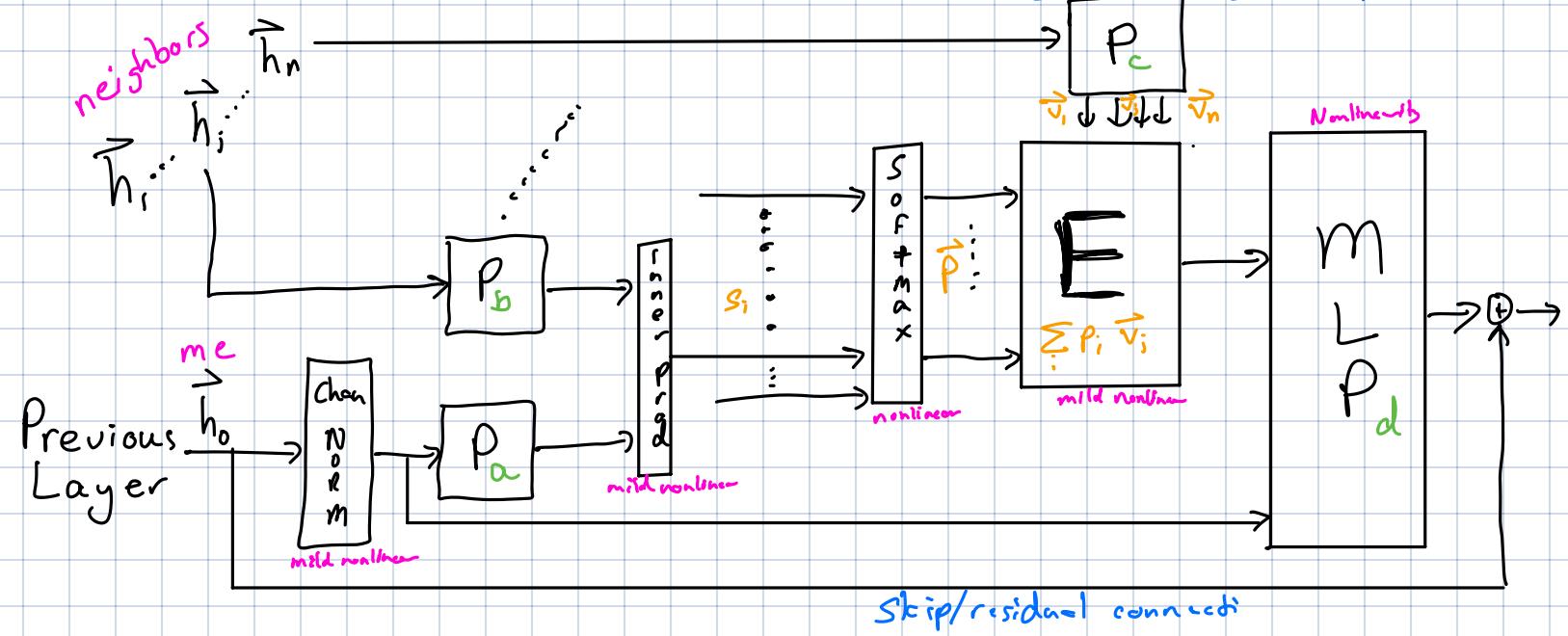
More generally: $f_{w_i}(\vec{h}_o, \sum \vec{g}_{w_i}(\vec{h}_o, \vec{h}_i))$

or with more structure $f_{w_i}(\vec{h}_o, \sum s_{w_2}(\vec{h}_o, \vec{h}_i) \vec{g}_{w_3}(\vec{h}_i))$

Should this be an expectation?

Natural Example: Design Choices:
GNN Building Block

Linear for g (Learned)
Inner Product for Similarity s (Learned)
Expectation-style for "sum" (Learned Linear Function)
MLP for combining f (Learned)



Question: What's nonlinear here?

Nonlinearity is a prereq for avoiding oversmoothing when we have no edge labels.

Note: This approach is related to attention mechanisms in Transformers but the GNN motivation is from minimal "what's possible" considerations.

In fact Transformers can be usefully thought of as GNNs on a complete graph.