

---

# How Attentive are Graph Attention Networks?

---

**Shaked Brody**

Technion, Israel

shakedbr@cs.technion.ac.il

**Uri Alon**

Technion, Israel

urialon@cs.technion.ac.il

**Eran Yahav**

Technion, Israel

yahave@cs.technion.ac.il

## GRAPH ATTENTION NETWORKS

**Petar Veličković\***

Department of Computer Science and Technology  
University of Cambridge  
petar.velickovic@cst.cam.ac.uk

**Guillem Cucurull\***

Centre de Visió per Computador, UAB  
gcucurull@gmail.com

**Arantxa Casanova\***

Centre de Visió per Computador, UAB  
ar.casanova.8@gmail.com

**Adriana Romero**

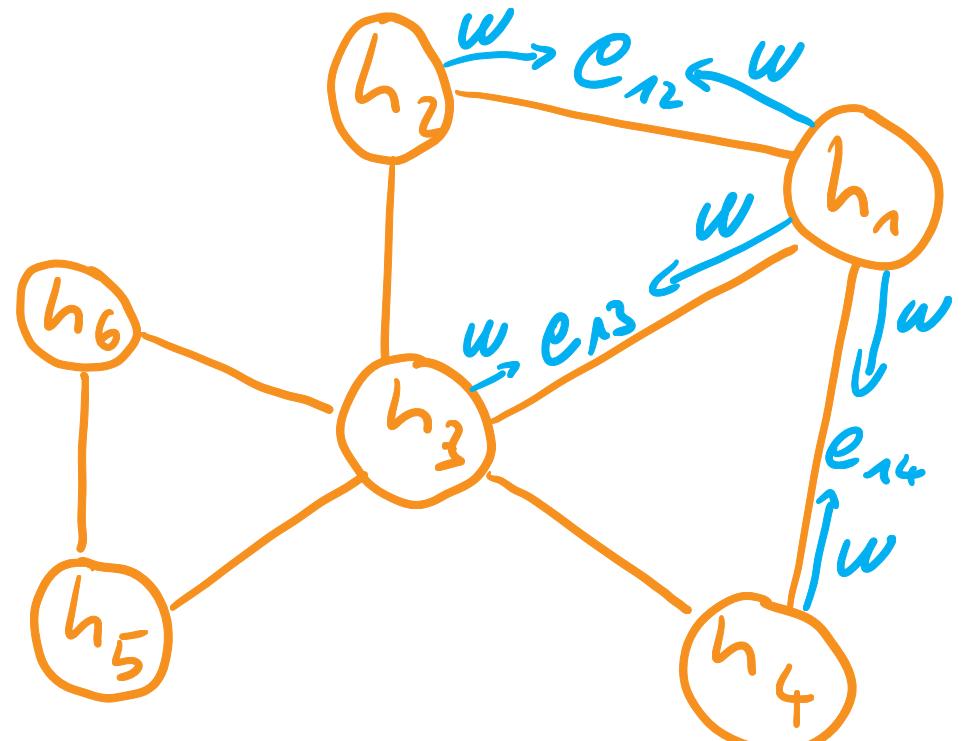
Montréal Institute for Learning Algorithms  
adriana.romero.soriano@umontreal.ca

**Pietro Liò**

Department of Computer Science and Technology  
University of Cambridge  
pietro.lio@cst.cam.ac.uk

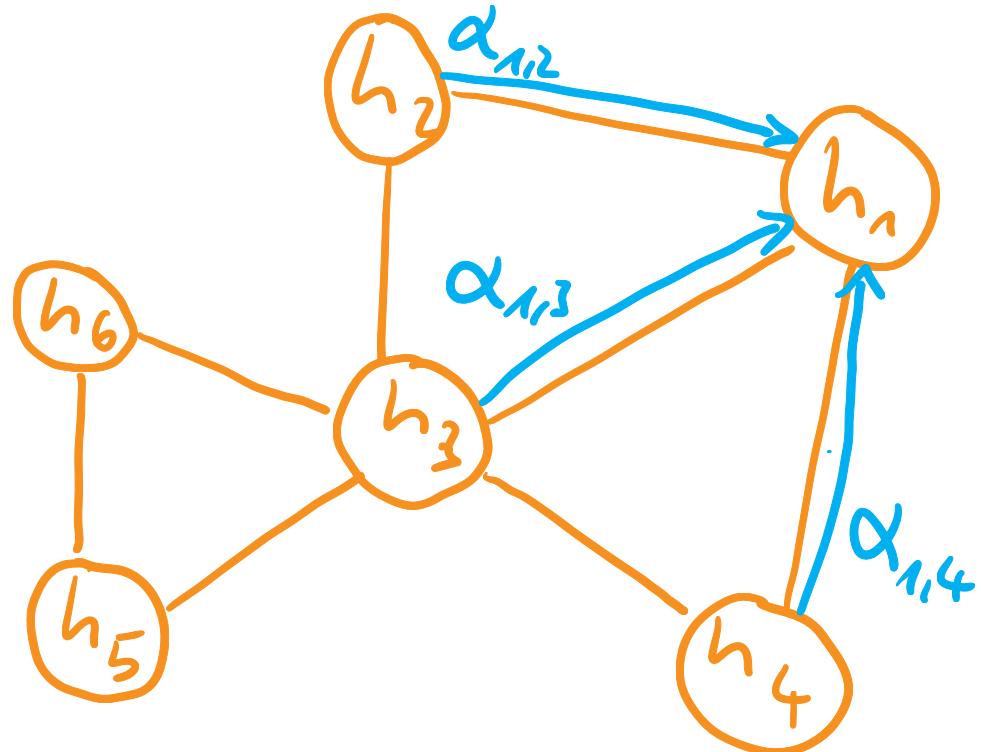
**Yoshua Bengio**

Montréal Institute for Learning Algorithms  
yoshua.umontreal@gmail.com



1. coefficients

$$e_{ij} = \text{LeakyReLU}(a^\top \cdot [Wh_i \| Wh_j])$$

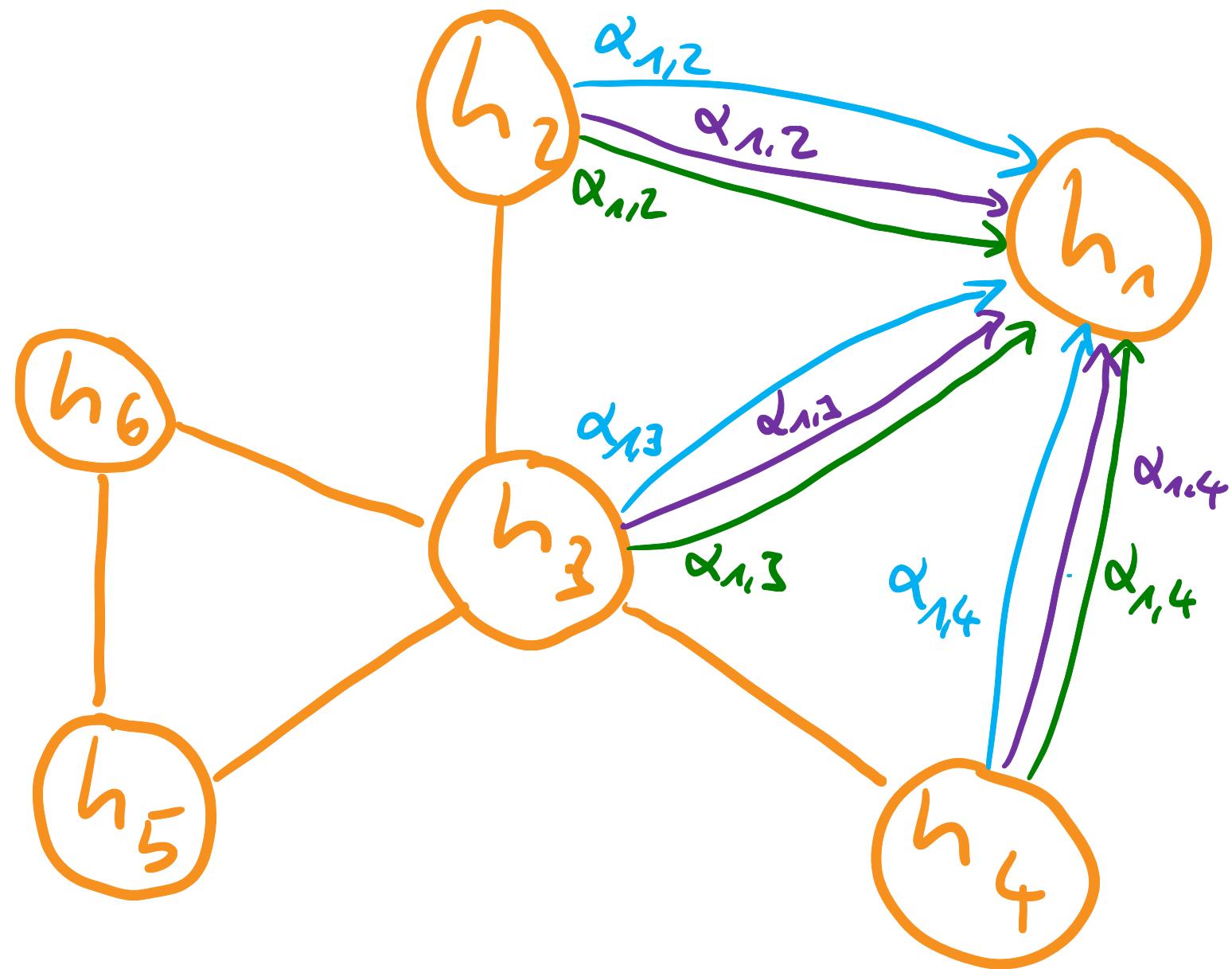


2. normalize:

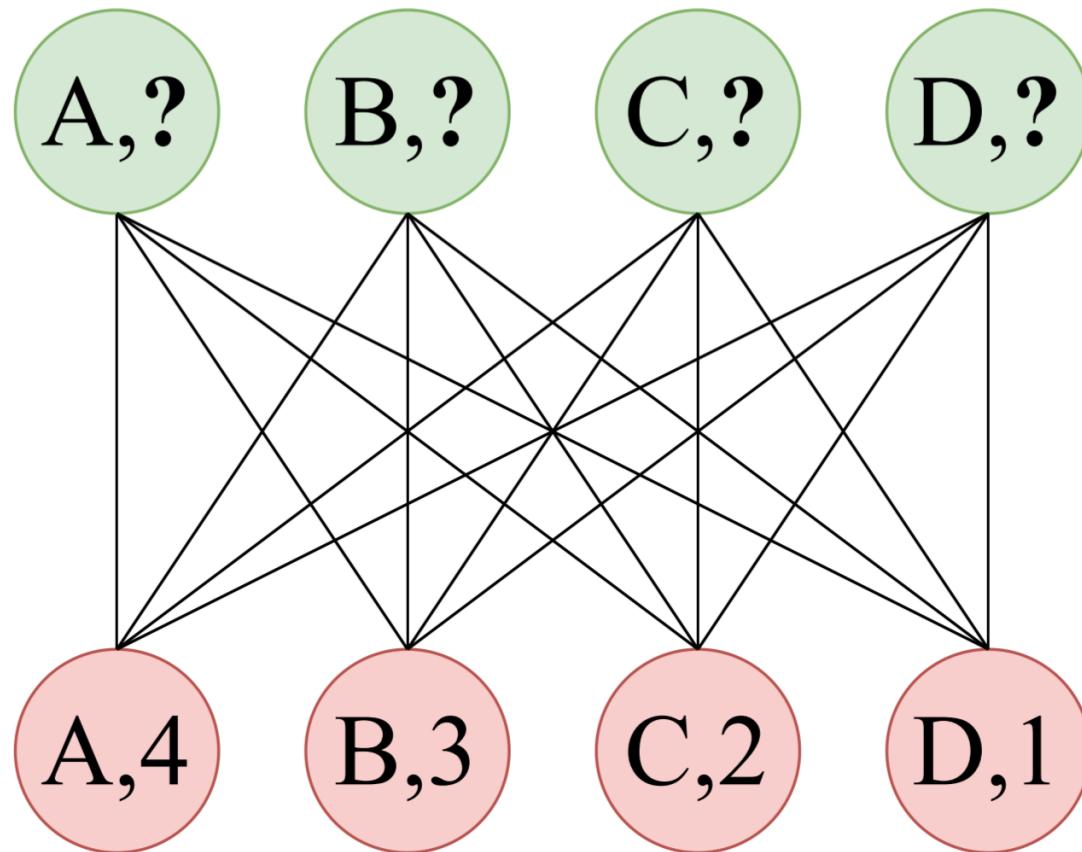
$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}_i} \exp(e_{ik})}.$$

3. aggregate

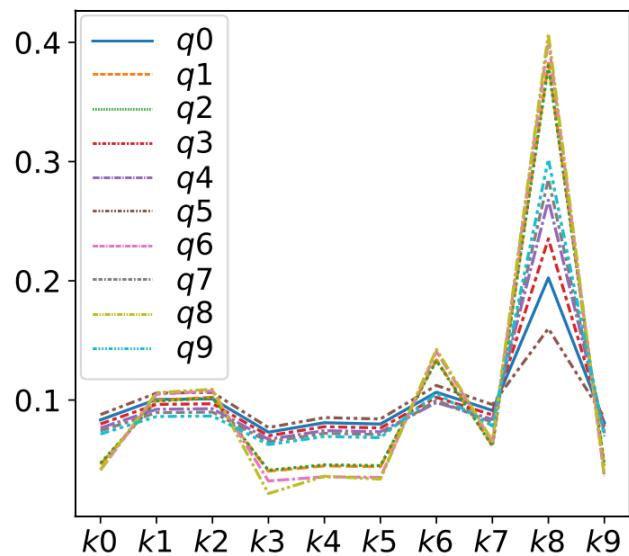
$$h'_i = \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij} \cdot Wh_j \right)$$



$h'_1 \parallel h'_2 \parallel h'_3$

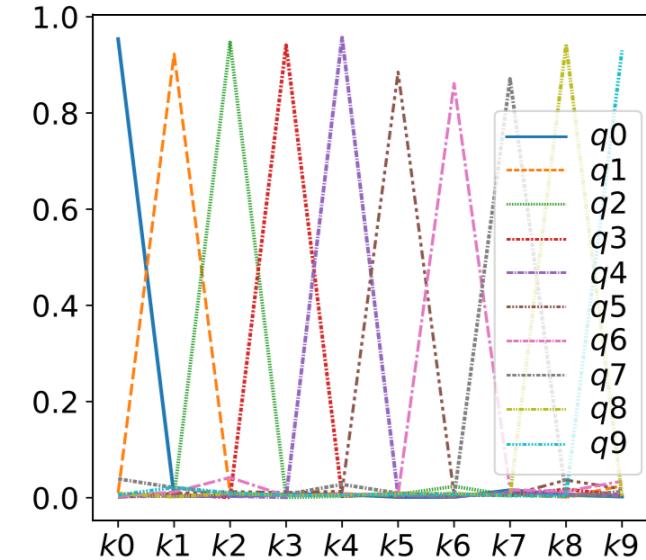


	$k0$	$k1$	$k2$	$k3$	$k4$	$k5$	$k6$	$k7$	$k8$	$k9$
$q0$	0.08	0.10	0.10	0.07	0.08	0.08	0.11	0.09	0.20	0.08
$q1$	-0.05	0.10	0.10	0.04	0.04	0.04	0.13	0.06	0.38	0.04
$q2$	-0.05	0.10	0.10	0.04	0.05	0.05	0.13	0.06	0.38	0.05
$q3$	-0.08	0.10	0.10	0.07	0.08	0.08	0.10	0.09	0.24	0.08
$q4$	-0.08	0.09	0.09	0.07	0.07	0.07	0.10	0.08	0.27	0.07
$q5$	-0.09	0.11	0.11	0.08	0.09	0.08	0.11	0.10	0.16	0.09
$q6$	-0.04	0.10	0.11	0.03	0.04	0.04	0.14	0.06	0.40	0.04
$q7$	-0.07	0.09	0.09	0.06	0.07	0.07	0.10	0.08	0.29	0.07
$q8$	-0.04	0.11	0.11	0.02	0.04	0.03	0.14	0.07	0.41	0.04
$q9$	-0.07	0.09	0.09	0.06	0.07	0.07	0.11	0.08	0.30	0.07



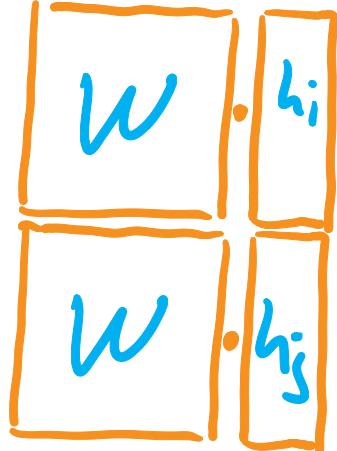
(a) Attention in standard GAT (Veličković et al. (2018))

	$k0$	$k1$	$k2$	$k3$	$k4$	$k5$	$k6$	$k7$	$k8$	$k9$
$q0$	-0.95	0.00	0.00	0.01	0.01	0.00	0.00	0.02	0.01	0.00
$q1$	-0.01	0.92	0.01	0.01	0.01	0.00	0.01	0.01	0.00	0.02
$q2$	-0.00	0.00	0.95	0.00	0.00	0.01	0.02	0.01	0.00	0.00
$q3$	-0.01	0.01	0.00	0.94	0.00	0.01	0.00	0.00	0.02	0.01
$q4$	-0.00	0.00	0.00	0.00	0.96	0.00	0.00	0.01	0.01	0.00
$q5$	-0.00	0.01	0.01	0.01	0.01	0.89	0.01	0.01	0.04	0.02
$q6$	-0.00	0.01	0.04	0.00	0.01	0.01	0.86	0.02	0.01	0.03
$q7$	-0.04	0.02	0.01	0.01	0.03	0.01	0.00	0.87	0.00	0.01
$q8$	-0.01	0.00	0.01	0.01	0.01	0.01	0.01	0.00	0.94	0.00
$q9$	-0.01	0.02	0.01	0.01	0.01	0.01	0.01	0.00	0.00	0.93



(b) Attention in GATv2, our fixed version of GAT

$$e_{i,j} = \text{LeakyReLU}(\alpha^T \cdot Wh_i || Wh_j)$$

$$e_{i,j} = (\boxed{\alpha^T} \cdot \boxed{w \cdot h_i})$$


$$e_{i,j} = (\boxed{\alpha_1^T} \boxed{w \cdot h_i} + \boxed{\alpha_2^T} \boxed{w \cdot h_j})$$

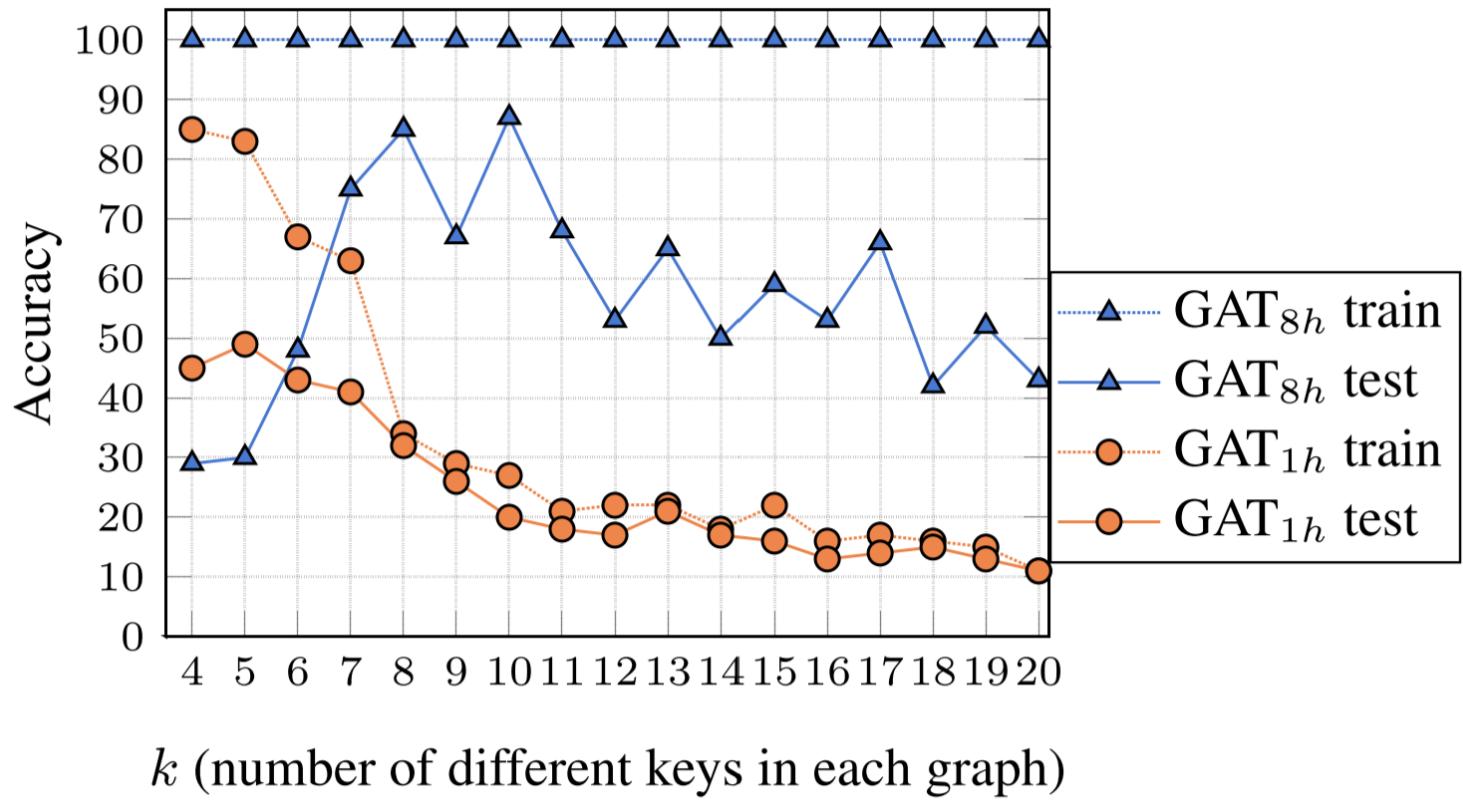
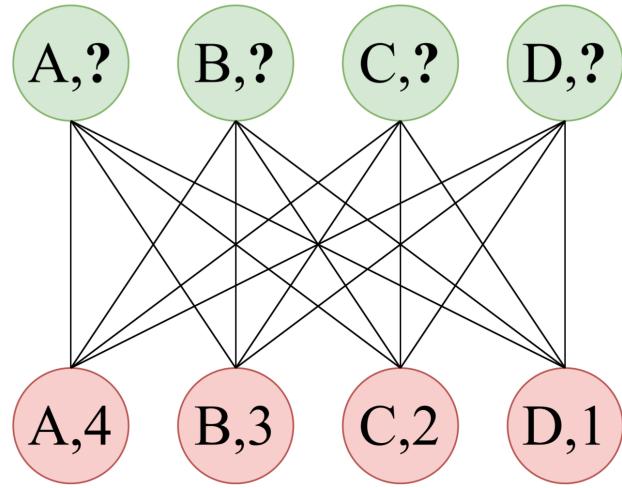
$$e_{i,j} = \text{LeakyReLU}(\alpha_1^T \cdot Wh_i + \alpha_2^T \cdot Wh_j)$$

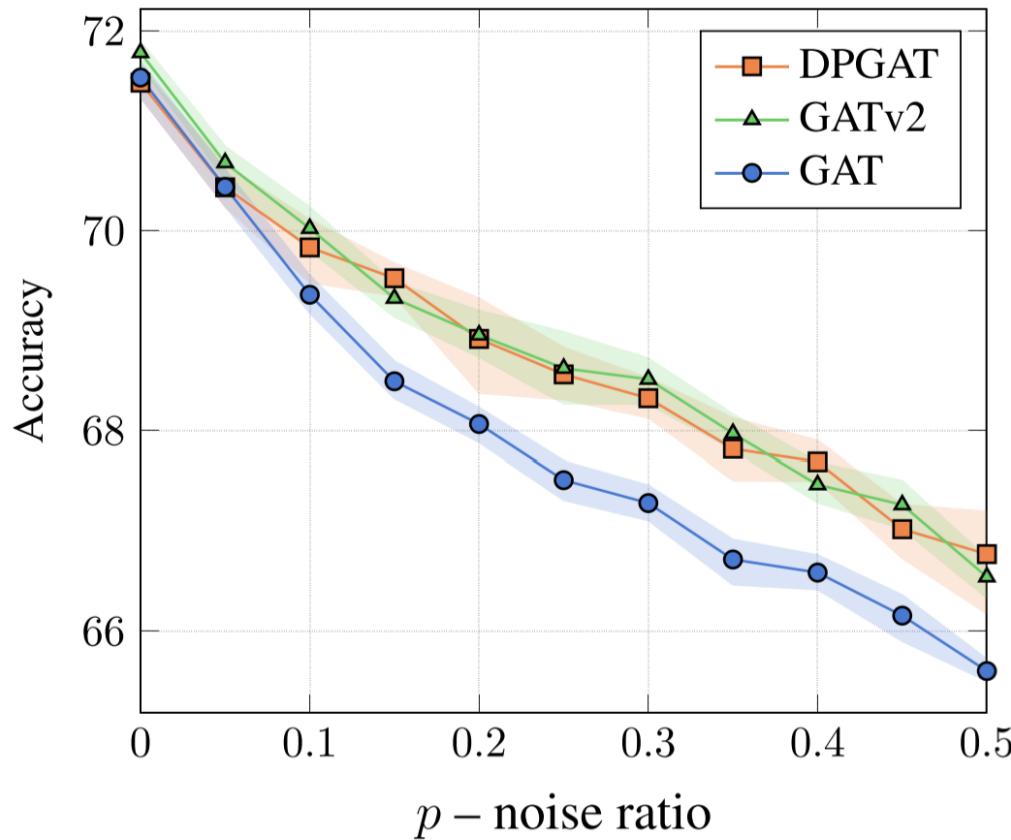
$$e_{i,j} = \text{LeakyReLU}(a^T \cdot Wh_i || Wh_j)$$

$$e_{i,j} = (a^T \cdot \begin{matrix} w \\ w \end{matrix} \cdot \begin{matrix} h_i \\ h_j \end{matrix})$$

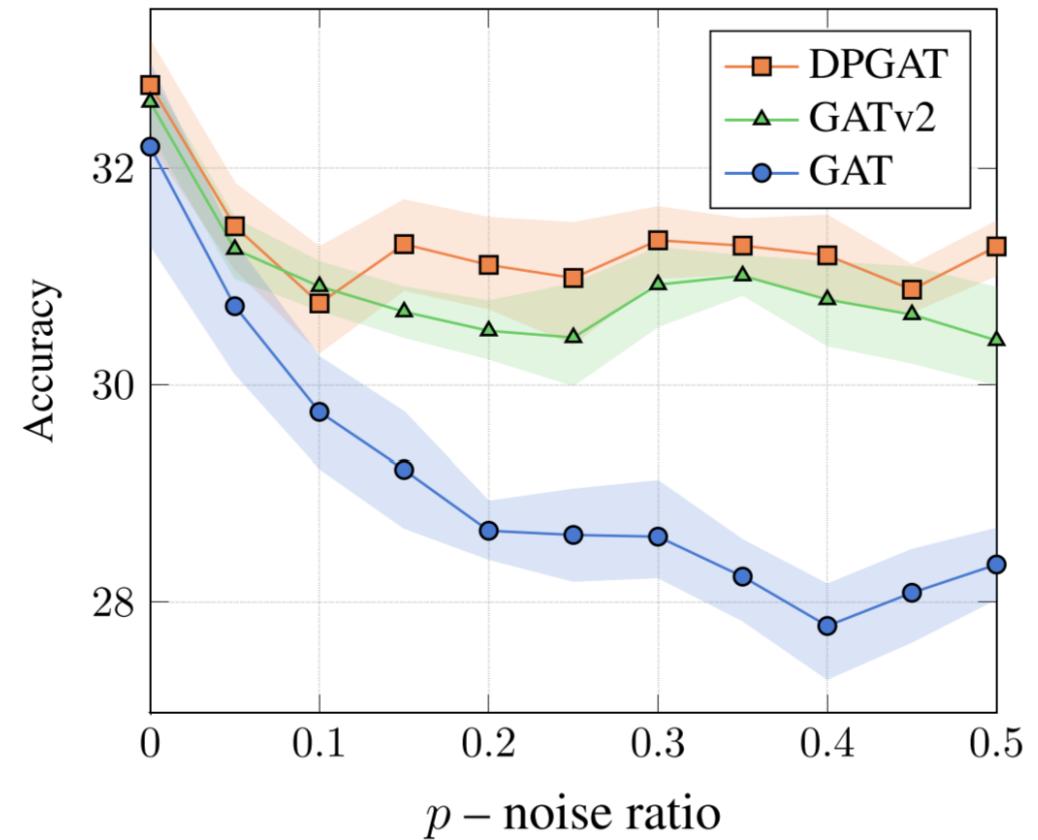
$$e_{i,j} = a^T \cdot \text{LeakyReLU}(w \cdot h_i || h_j)$$

$$e_{i,j} = a^T \cdot (w \cdot \begin{matrix} h_i \\ h_j \end{matrix})$$

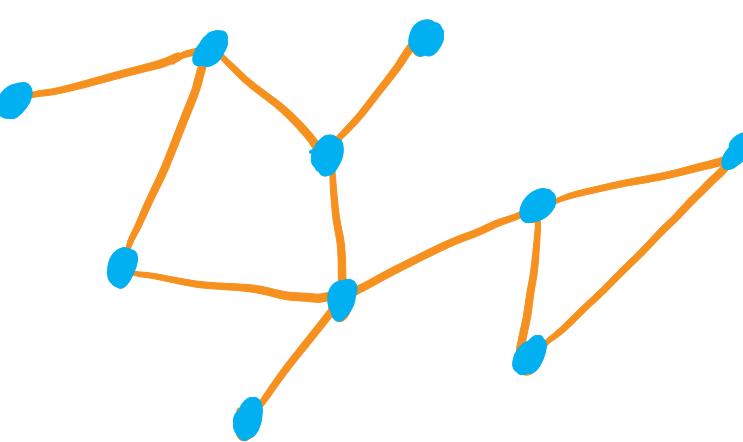




(a) **ogbn-arxiv**



(b) **ogbn-mag**



	Model	<b>ogbn-arxiv</b>	<b>ogbn-products</b>	<b>ogbn-mag</b>	<b>ogbn-proteins</b>
No-Attention	GCN <sup>†</sup>	$71.74 \pm 0.29$	$78.97 \pm 0.33$	$30.43 \pm 0.25$	$72.51 \pm 0.35$
	GraphSAGE <sup>†</sup>	$71.49 \pm 0.27$	$78.70 \pm 0.36$	$31.53 \pm 0.15$	$77.68 \pm 0.20$
Attention 1 head	GAT <sub>1h</sub>	$71.59 \pm 0.38$	$79.04 \pm 1.54$	$32.20 \pm 1.46$	$70.77 \pm 5.79$
	DPGAT <sub>1h</sub>	$71.52 \pm 0.17$	$76.49 \pm 0.78$	$32.77 \pm 0.80$	$63.47 \pm 2.79$
	GATv2 <sub>1h</sub>	$71.78 \pm 0.18$	<b><math>80.63 \pm 0.70</math></b>	$32.61 \pm 0.44$	$77.23 \pm 3.32$
Attention 8 heads	GAT <sub>8h</sub>	$71.54 \pm 0.30$	$77.23 \pm 2.37$	$31.75 \pm 1.60$	$78.63 \pm 1.62$
	DPGAT <sub>8h</sub>	$71.48 \pm 0.26$	$73.53 \pm 0.47$	$27.74 \pm 9.97$	$72.88 \pm 0.59$
	GATv2 <sub>8h</sub>	<b><math>71.87 \pm 0.25</math></b>	$78.46 \pm 2.45$	$32.52 \pm 0.39$	<b><math>79.52 \pm 0.55</math></b>

(a)

(b)

Table 2: Average accuracy (Table 2a) and ROC-AUC (Table 2b) in node-prediction datasets (10 runs $\pm$ std). In all datasets, GATv2 outperforms GAT. <sup>†</sup> – previously reported by Hu et al. (2020).

		<b>ogbl-collab</b>		<b>ogbl-citation2</b>
		w/o val edges	w/ val edges	
Model				
No- Attention	GCN <sup>†</sup>	44.75±1.07	47.14±1.45	80.04±0.25
	GraphSAGE <sup>†</sup>	48.10±0.81	54.63±1.12	<b>80.44</b> ±0.10
Attention 1 head	GAT <sub>1h</sub>	39.32±3.26	48.10±4.80	79.84±0.19
	DPGAT <sub>1h</sub>	44.15±1.59	49.61±3.16	80.31±0.21
	GATv2 <sub>1h</sub>	42.00±2.40	48.02±2.77	80.33±0.13
Attention 8 heads	GAT <sub>8h</sub>	42.37±2.99	46.63±2.80	75.95±1.31
	DPGAT <sub>8h</sub>	<b>48.99</b> ±0.70	<b>56.90</b> ±1.45	79.64±0.42
	GATv2 <sub>8h</sub>	42.85±2.64	49.70±3.08	80.14±0.71

(a)

(b)

Table 3: Average Hits@50 (Table 3a) and mean reciprocal rank (MRR) (Table 3b) in link-prediction benchmarks from OGB (10 runs±std). † – previously reported by Hu et al. (2020).

Model	Predicted Property													Rel. to $\text{GAT}_{8h}$
	1	2	3	4	5	6	7	8	9	10	11	12	13	
GCN <sup>†</sup>	3.21	<b>4.22</b>	1.45	1.62	2.42	16.38	17.40	7.82	8.24	9.05	7.00	3.93	<b>1.02</b>	-1.5%
GIN <sup>†</sup>	2.64	4.67	1.42	1.50	2.27	<b>15.63</b>	12.93	<b>5.88</b>	18.71	<b>5.62</b>	<b>5.38</b>	<b>3.53</b>	1.05	-2.3%
$\text{GAT}_{1h}$	3.08	7.82	1.79	3.96	3.58	35.43	116.5	28.10	20.80	15.80	10.80	5.37	3.11	+134.1%
$\text{DPGAT}_{1h}$	3.20	8.35	1.71	2.17	2.88	25.21	65.79	12.93	13.32	14.42	13.83	6.37	3.28	+77.9%
$\text{GATv2}_{1h}$	3.04	6.38	1.68	2.18	2.82	20.56	77.13	10.19	22.56	15.04	22.94	5.23	2.46	+91.6%
$\text{GAT}_{8h}^{\dagger}$	2.68	4.65	1.48	1.53	2.31	52.39	14.87	7.61	6.86	7.64	6.54	4.11	1.48	+0%
$\text{DPGAT}_{8h}$	<b>2.63</b>	4.37	1.44	<b>1.40</b>	<b>2.10</b>	32.59	<b>11.66</b>	6.95	7.09	7.30	6.52	3.76	1.18	-9.7%
$\text{GATv2}_{8h}$	2.65	4.28	<b>1.41</b>	1.47	2.29	16.37	14.03	6.07	<b>6.28</b>	6.60	5.97	3.57	1.59	<b>-11.5%</b>

Table 4: Average error rates (lower is better), 5 runs for each property, on the QM9 dataset. <sup>†</sup> was previously tuned and reported by [Brockschmidt \(2020\)](#). Standard deviation is reported in Appendix F.

$$e_{i,j} = \frac{h_i^T Q \cdot (h^T K)^T}{\sqrt{d_K}}$$



