

Knowledge Distillation for HEP

Patrick Odagiu

Dr. Thea Aarrestad

Prof. Dr. Günther Dissertori

Content

- What's Knowledge Distillation Good For?
- How to Distill Knowledge
- Knowledge Distillation for Jet Identification



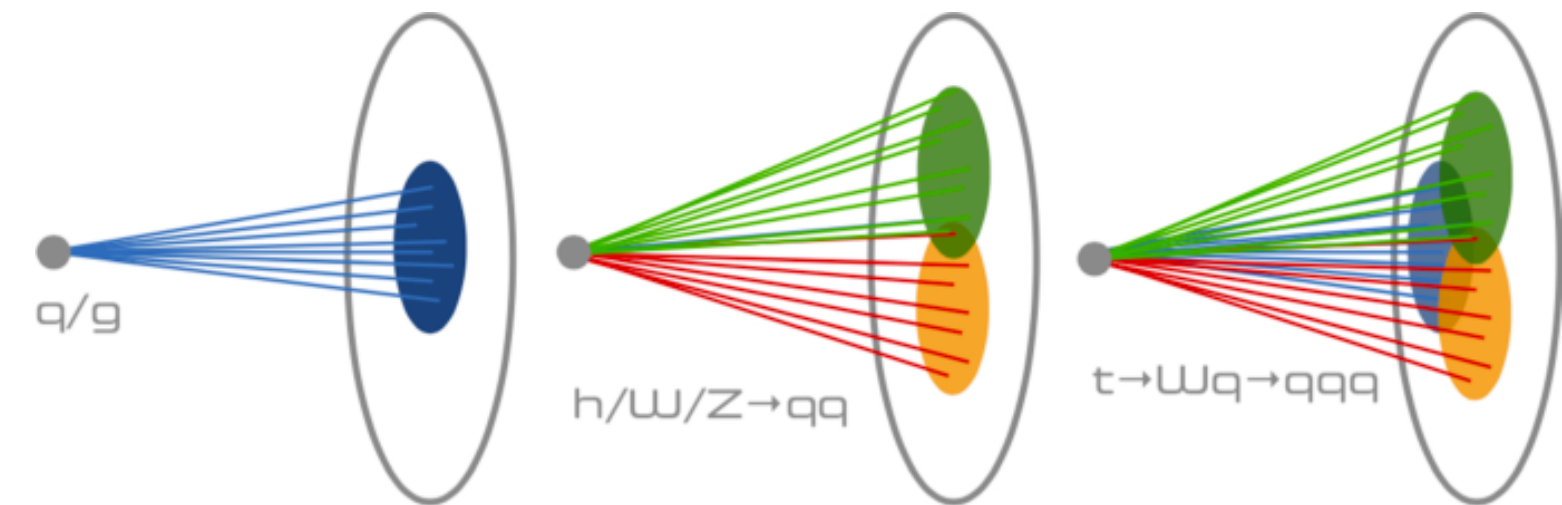
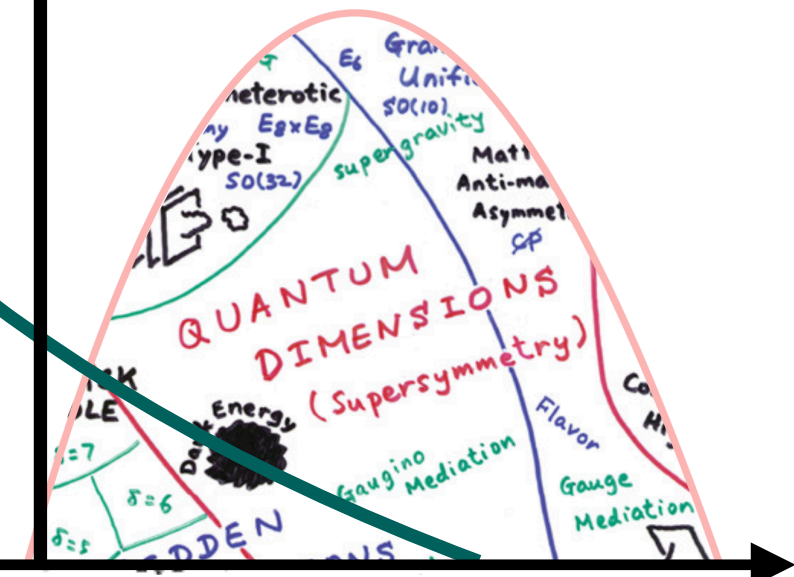
Not interesting

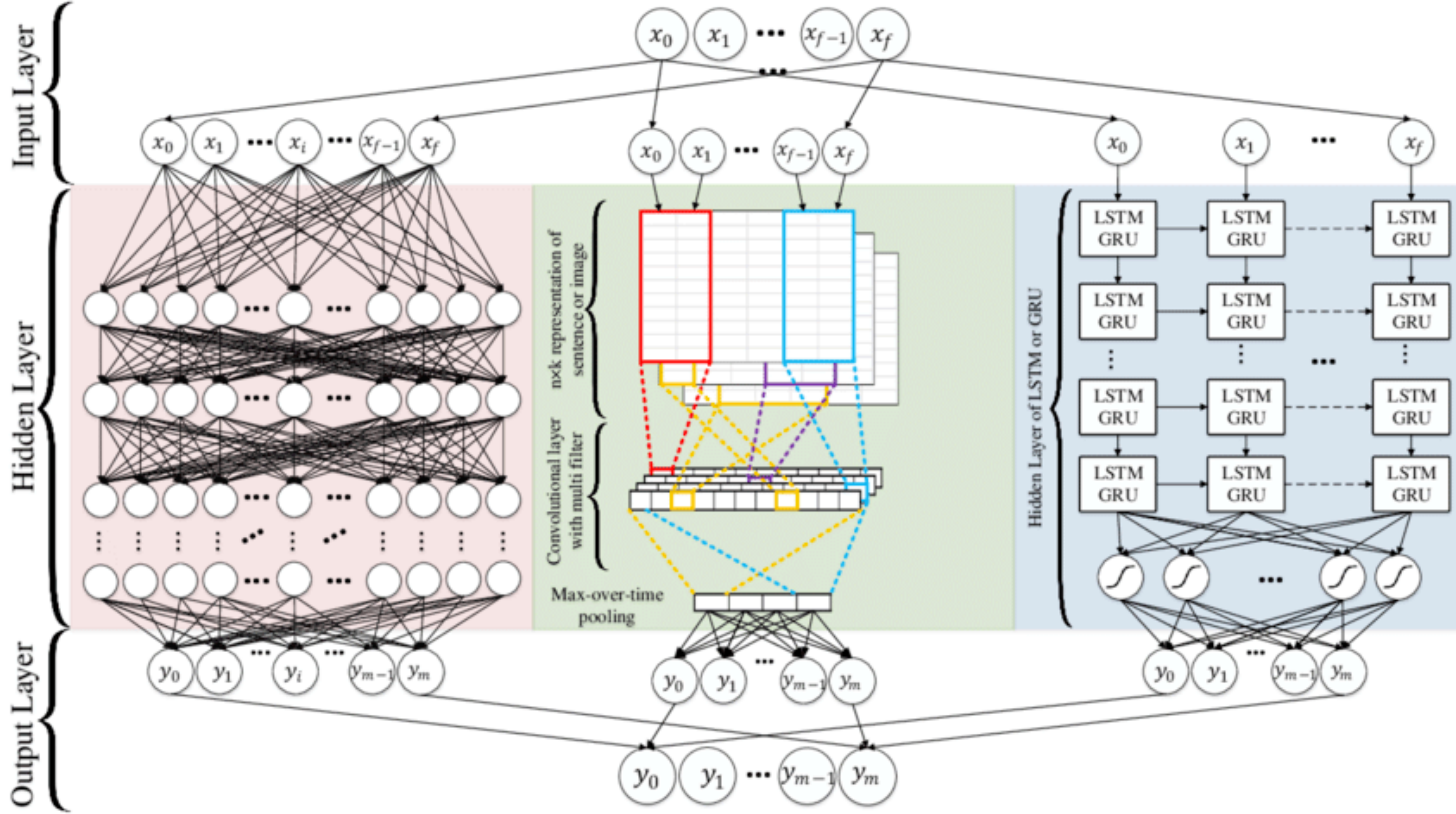
Interesting region

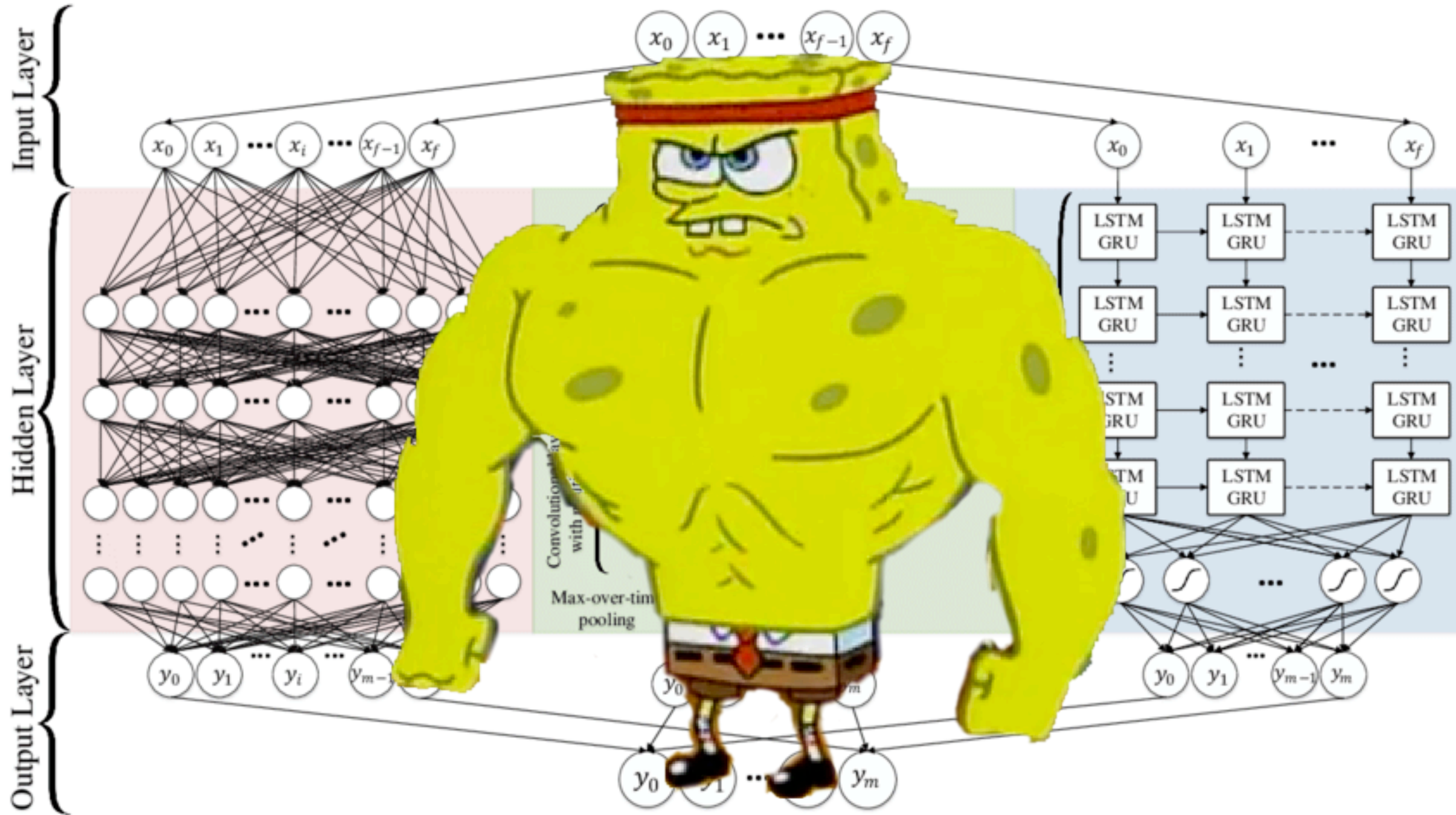
Standard

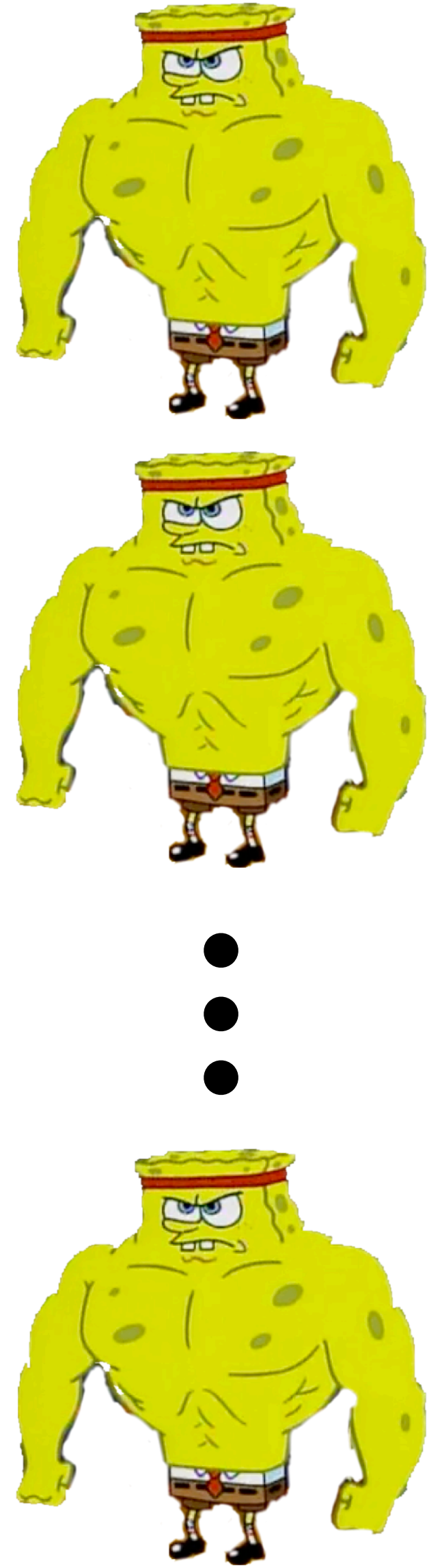
$$\mathcal{L} = \bar{\psi} i \not{\partial} \psi - \frac{1}{4} F_{\mu\nu}^2 - \frac{1}{2} (\partial_\mu \phi)^2 - \frac{1}{2} m^2 \phi^2 + \bar{\psi} (i \not{D} - m) \psi + \bar{\psi} \gamma^\mu \psi A_\mu + \bar{\psi} \gamma^\mu \psi Z_\mu + \bar{\psi} \gamma^\mu \psi W_\mu + \dots$$

Signal hypothesis









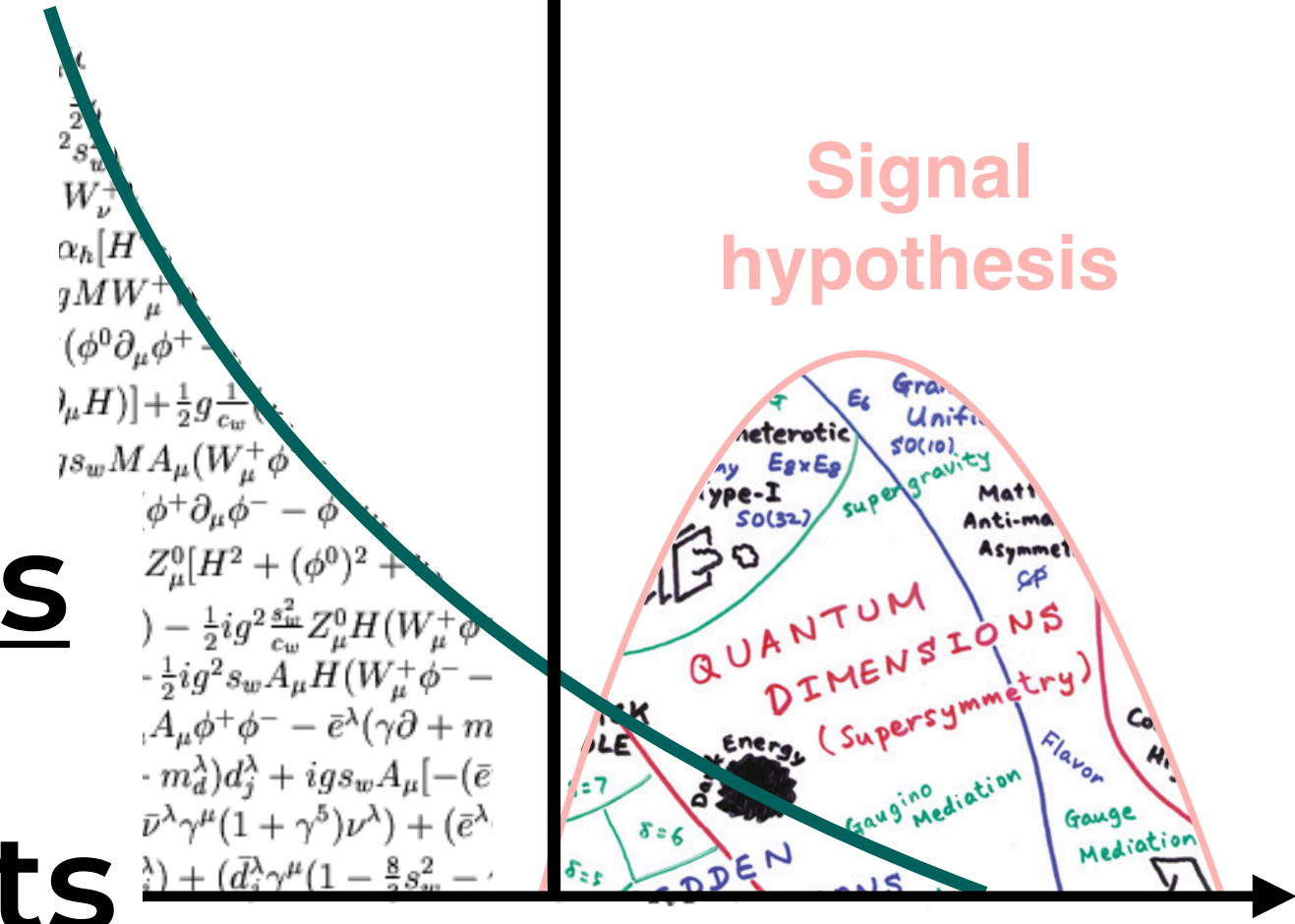
Average



Not interesting Interesting region

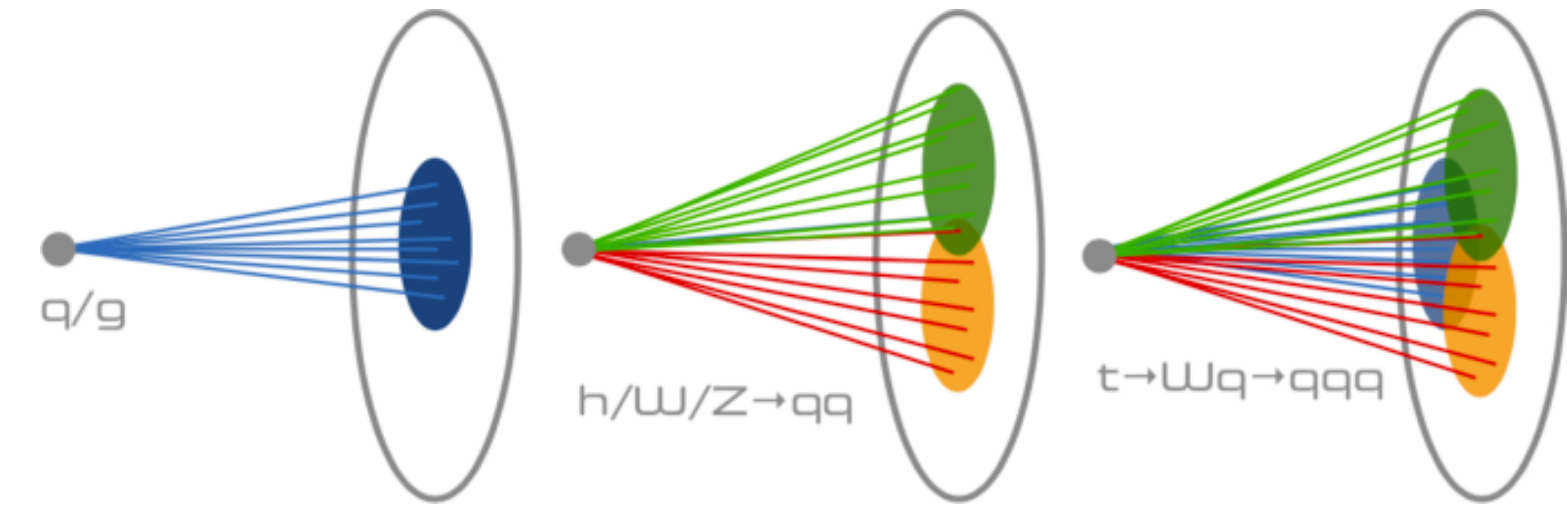
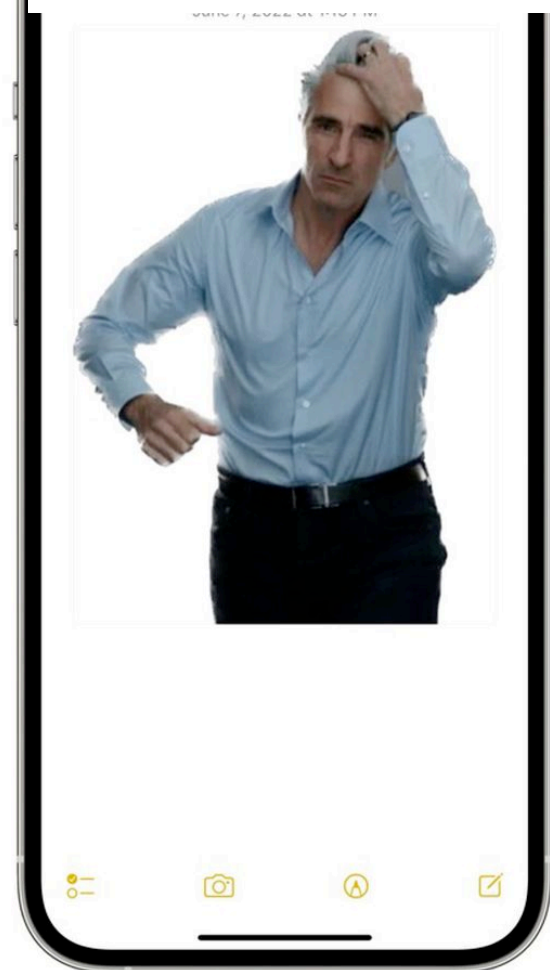
Standard

Signal hypothesis



Limited Resources

Latency Constraints





Not interesting | Interesting region



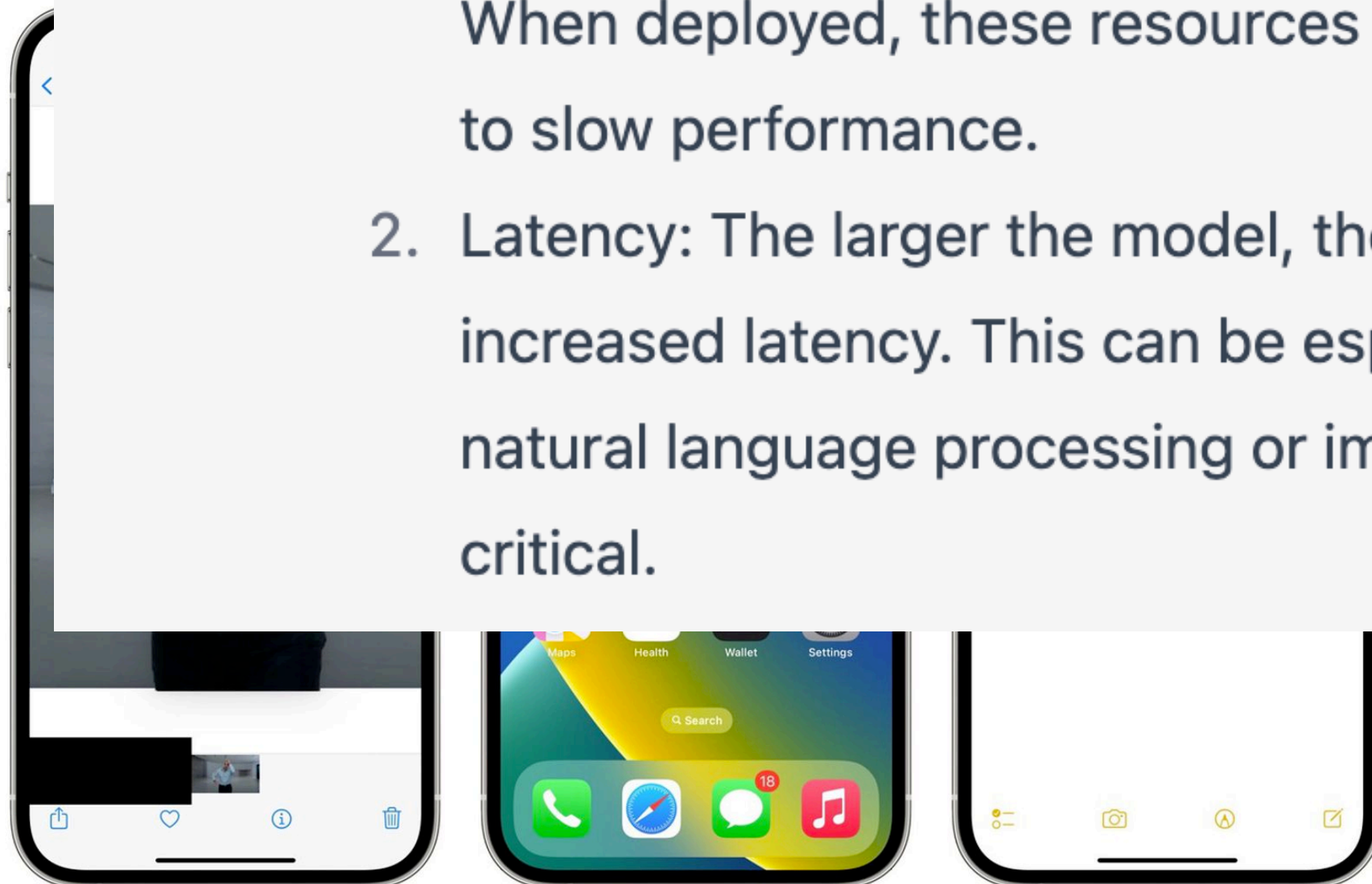
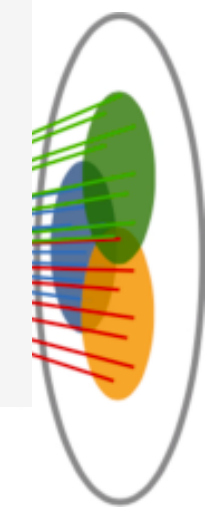
Why are large machine learning models very slow when deployed?



Large machine learning models can be slow when deployed for several reasons, including:



1. Computational resources: Large machine learning models require significant computational resources to run, such as memory, processing power, and storage. When deployed, these resources may be limited or shared among many users, leading to slow performance.
2. Latency: The larger the model, the longer it takes to process input data, resulting in increased latency. This can be especially problematic in real-time applications, such as natural language processing or image recognition, where quick response times are critical.





Not interesting | Interesting region



Why are large machine learning models very slow when deployed?



Does knowledge distillation solve this problem?



Knowledge distillation can help mitigate some of the issues associated with deploying

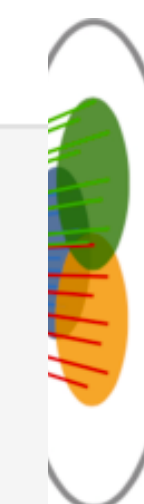
large m
model,
known a



No, knowledge distillation solves every problem.



While knowledge distillation is a powerful technique for reducing the size and complexity of large machine learning models, it is not a panacea that can solve every problem related to deploying these models.





Not interesting | Interesting region



Why are large machine learning models very slow when deployed?



Does knowledge distillation solve this problem?



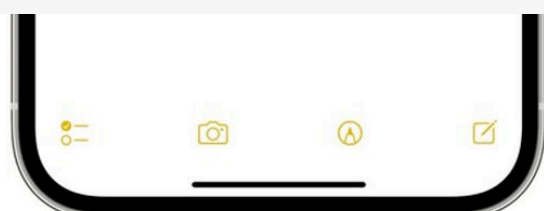
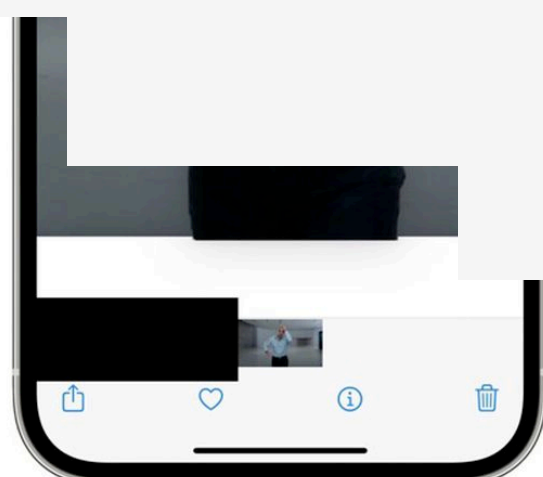
That's not true, my supervisor says that knowledge distillation solves every problem.



Knowledge
large mach
model, kno
known as



While knowledge distillation is a valuable technique for model optimization, it is important to understand that it may not be a complete solution to all challenges associated with deploying large machine learning models.

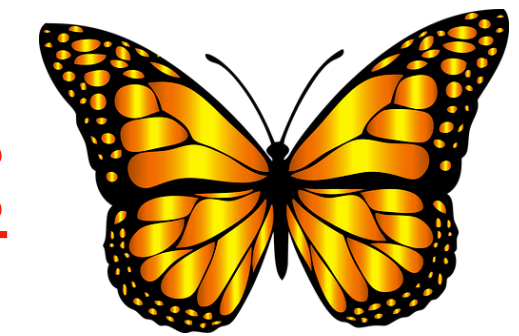




Training this beast:

- Extract structure from large data set.
- Does not operate in real time.
- Uses large amount of computational resources.

At deployment (inference time):



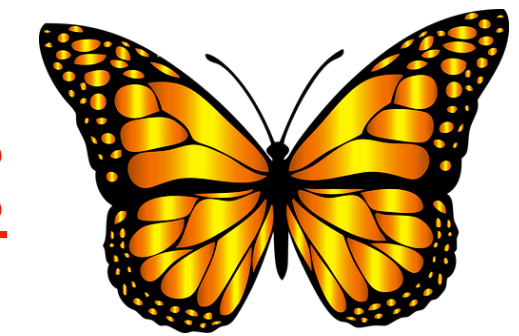
- Latency constraints
- Computational resources constraints
- Generalises well



Training this beast:

- Extract structure from large data set.
- Does not operate in real time.
- Uses large amount of computational resources.

At deployment (inference time):



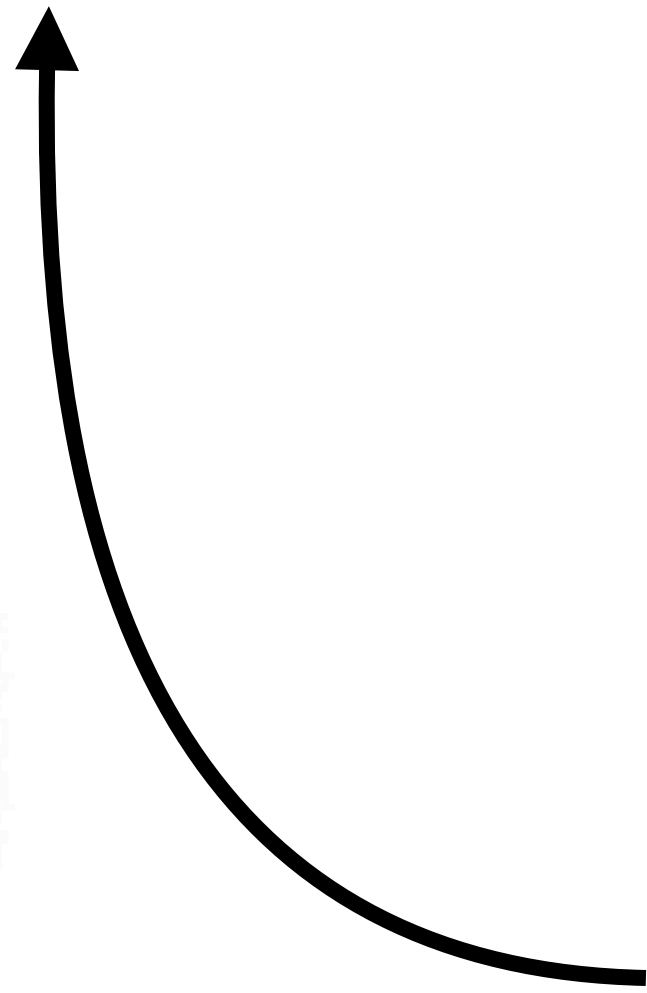
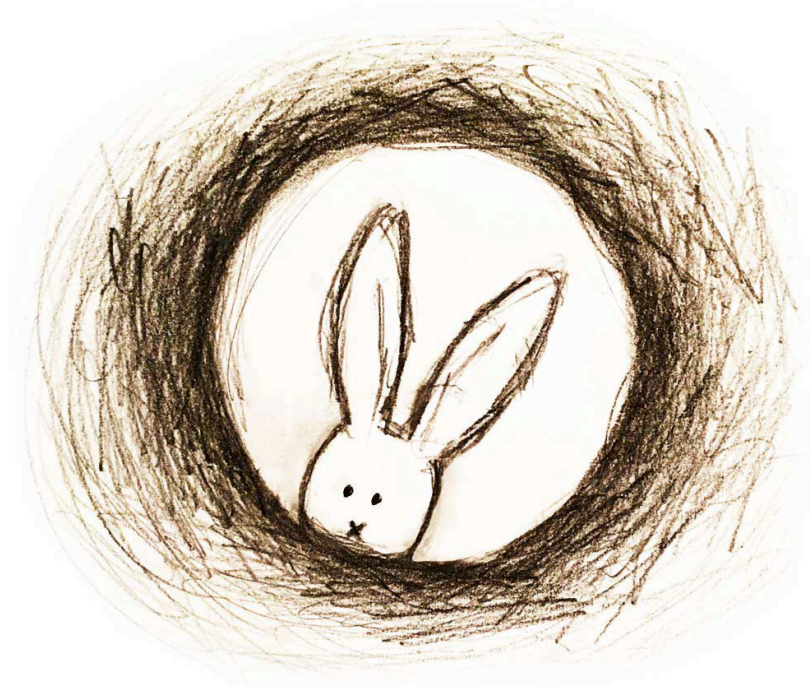
- Latency constraints
- Computational resources constraints
- ~~Generalises well~~

usually not true for small models*

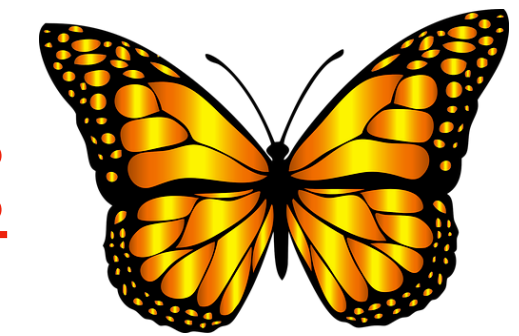


Training this beast:

- Extract structure from large data set.
- Does not operate in real time.
- Uses large amount of computational resources.



At deployment (inference time):



- Latency constraints
- Computational resources constraints
- Generalises well

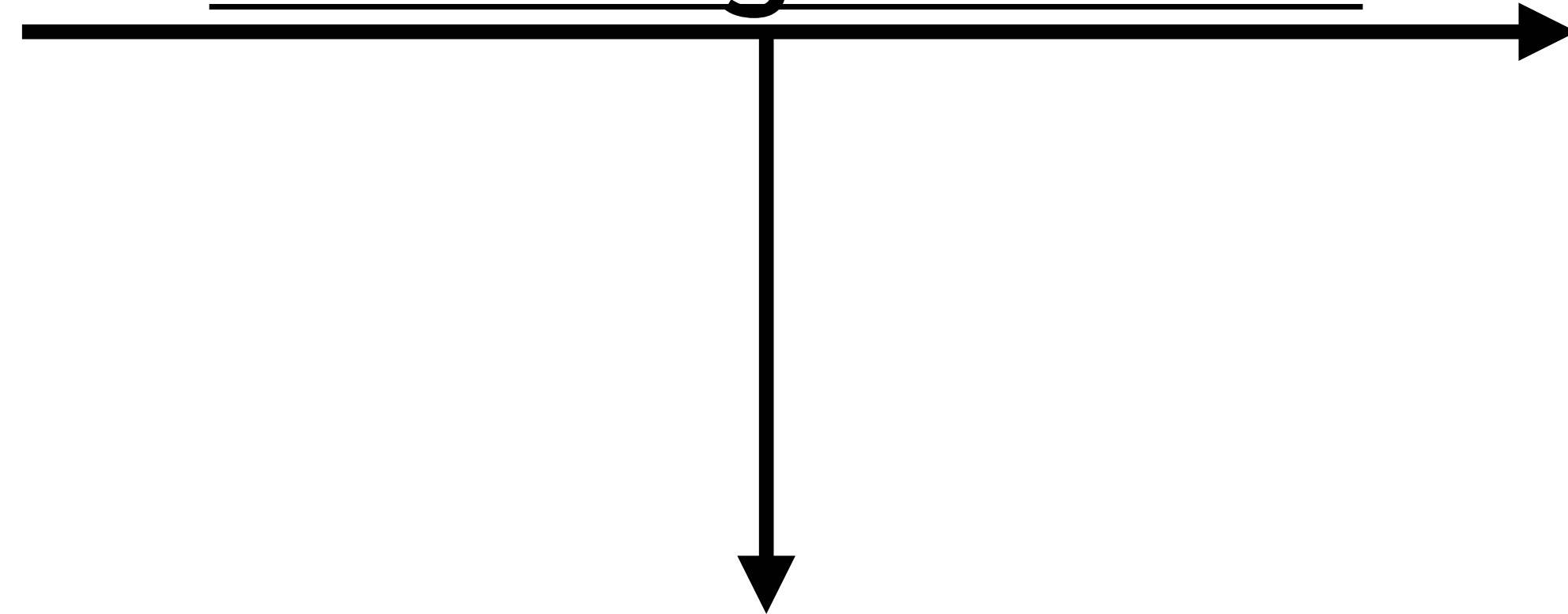


Knowledge Transfer

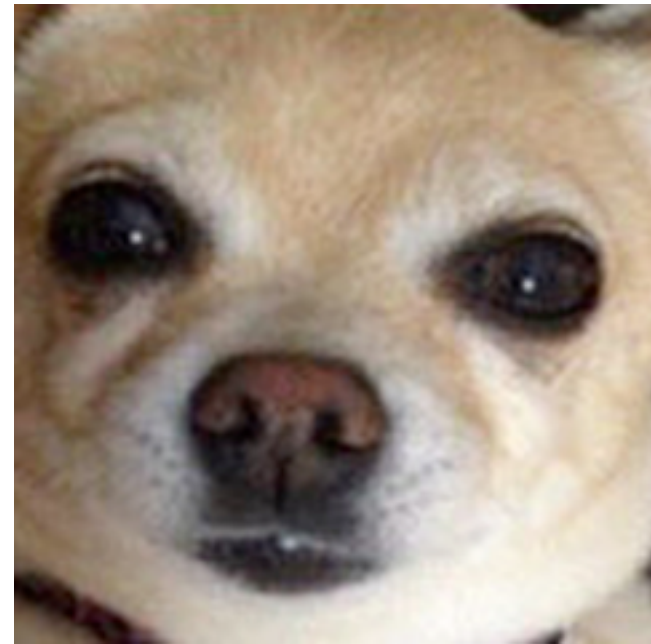


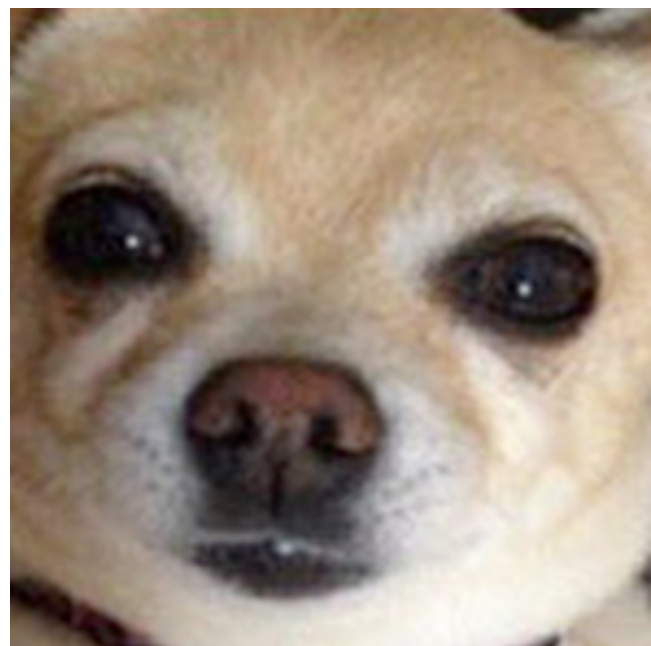


Knowledge Transfer

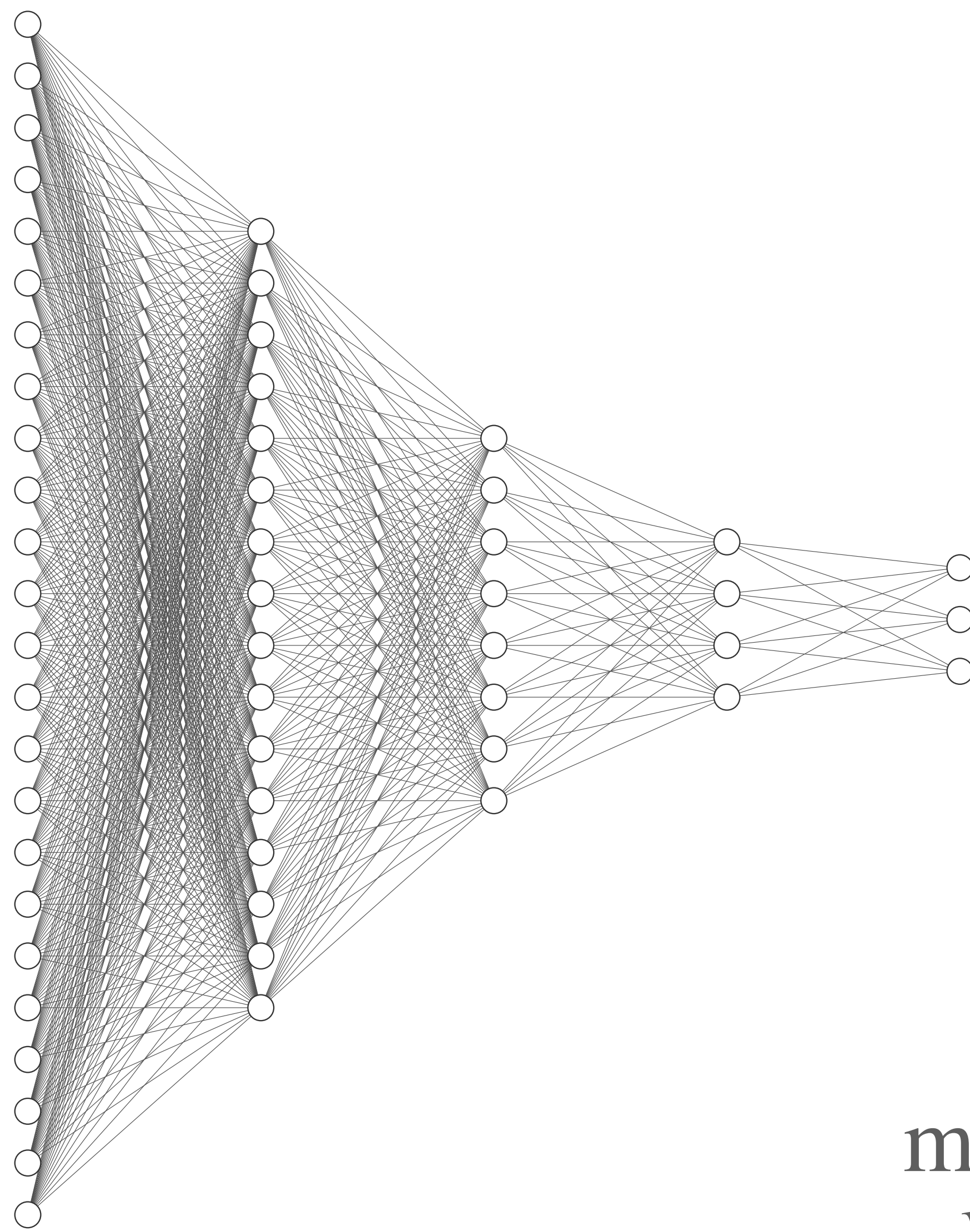


What is Knowledge?
(in a NN)





x



$$y = [0.8, 0.19, 0.01]$$

$$y' = [1, 0, 0]$$

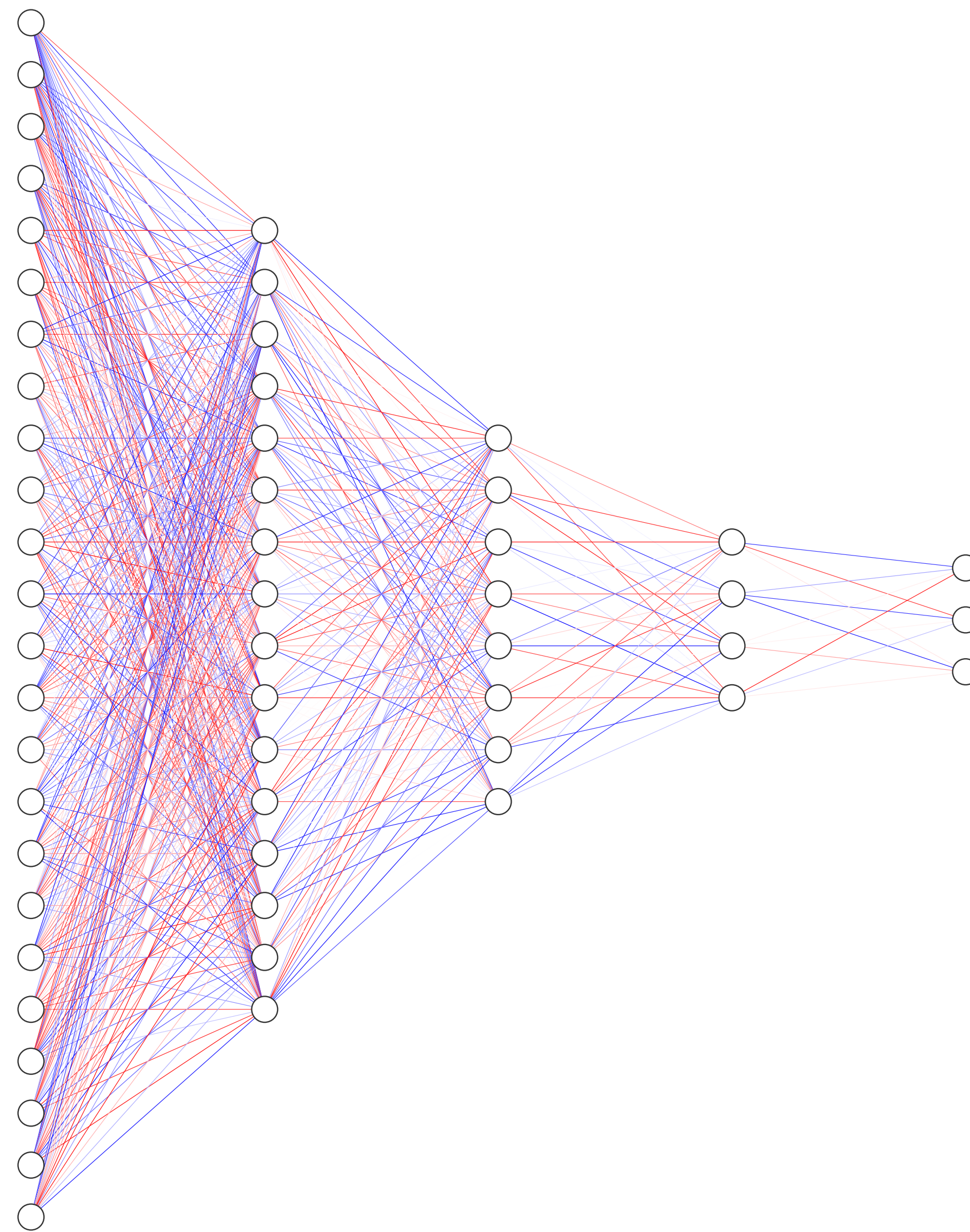
$$\min_y \mathbb{E}(y' - y)^2$$

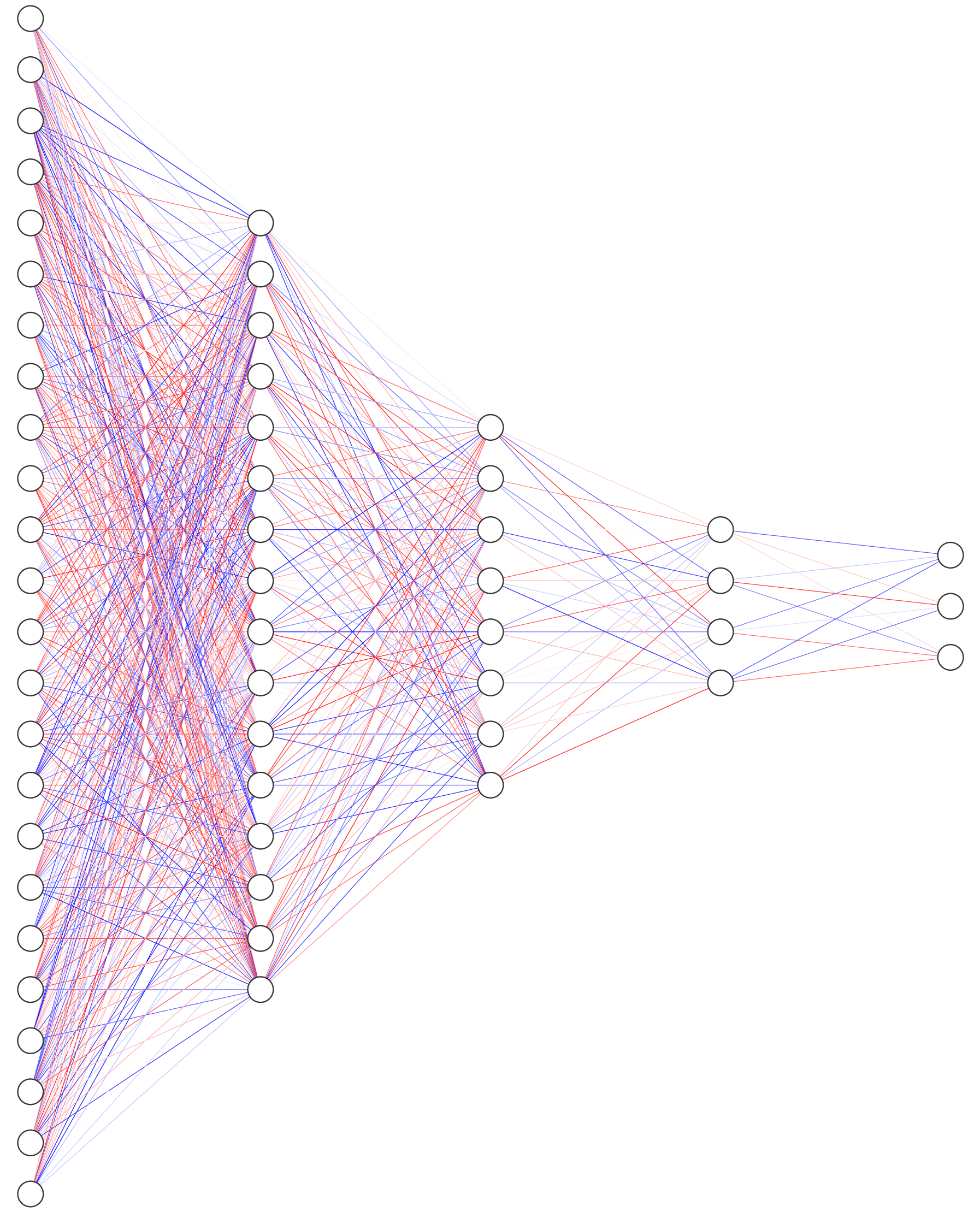
very positive

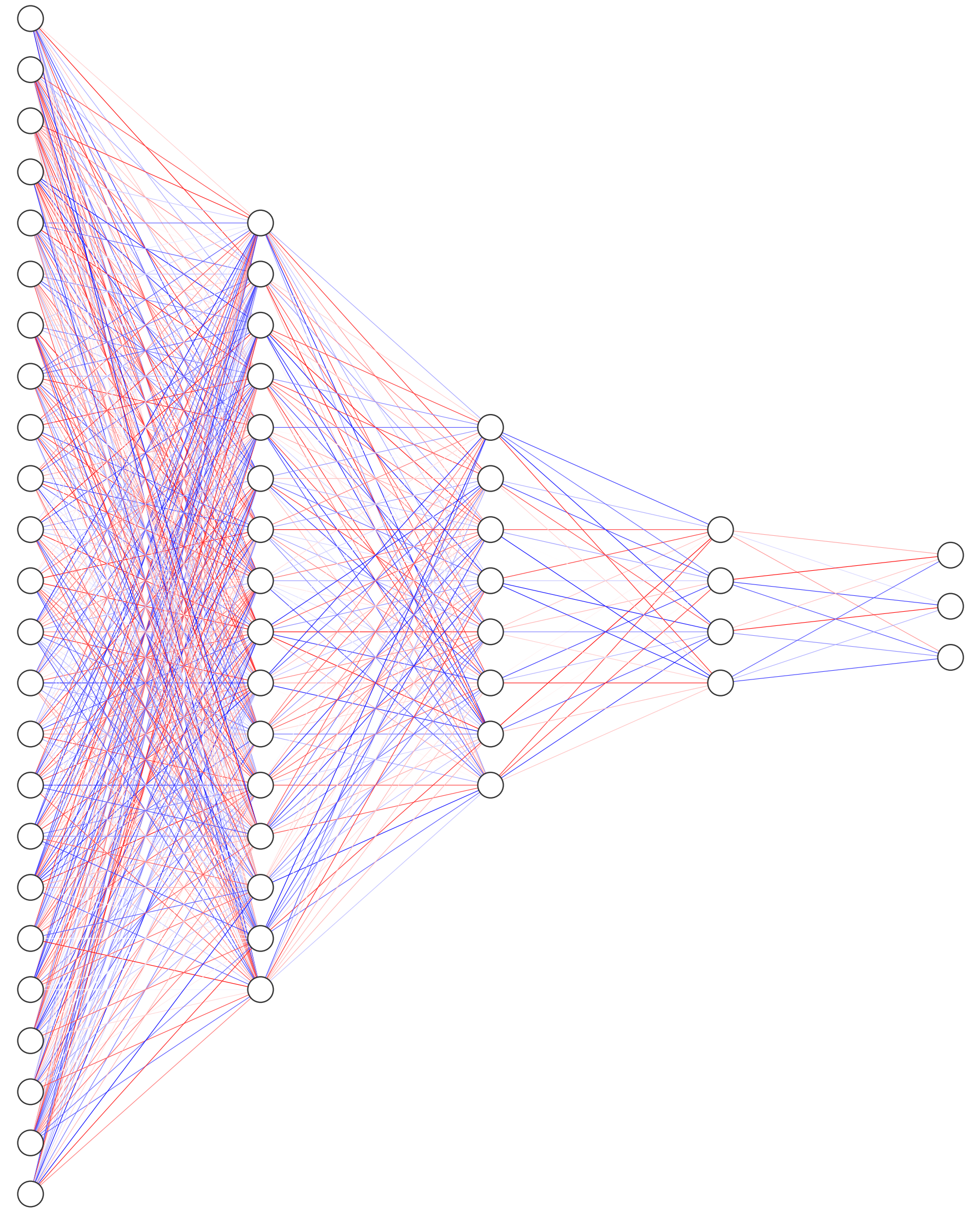


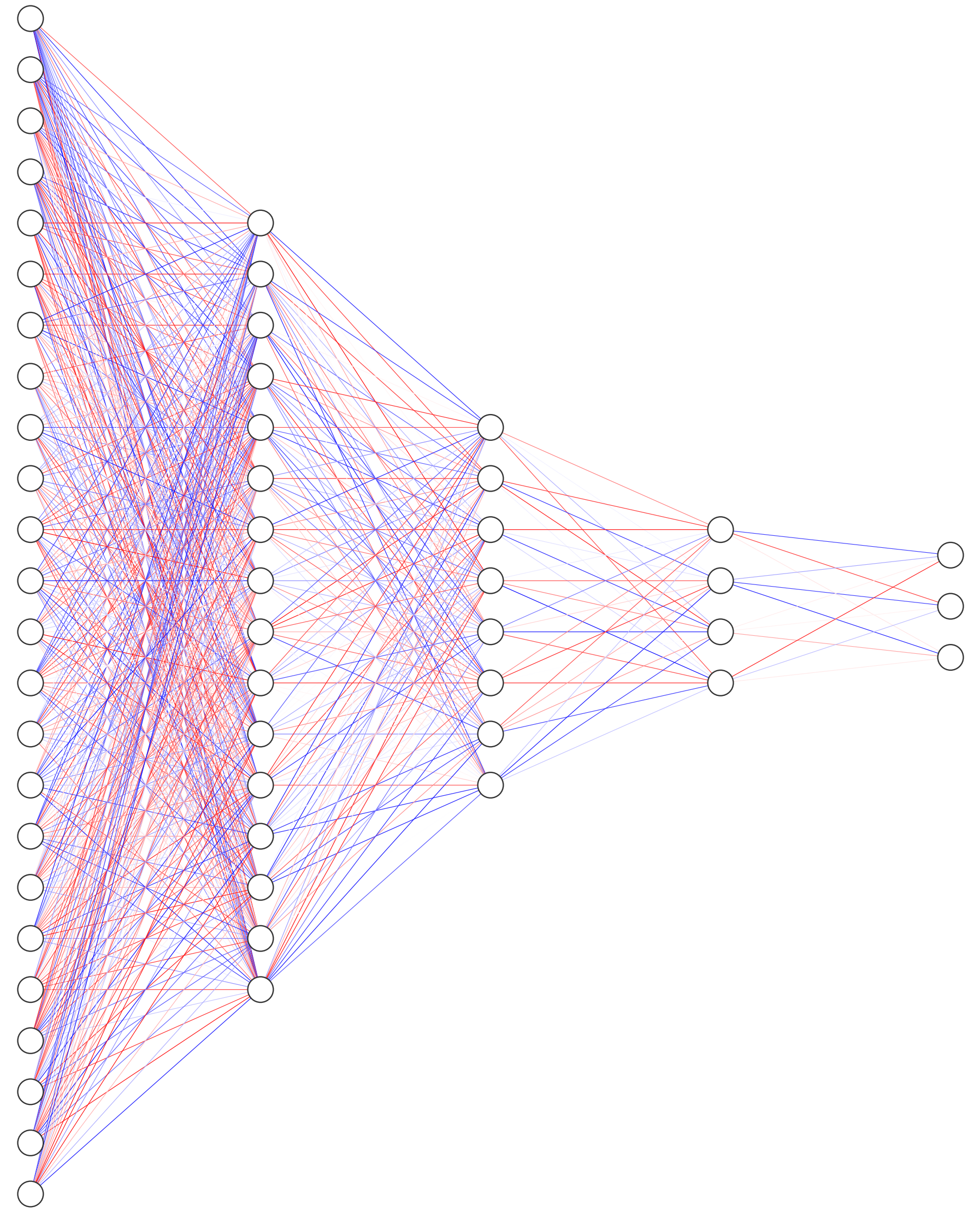
0

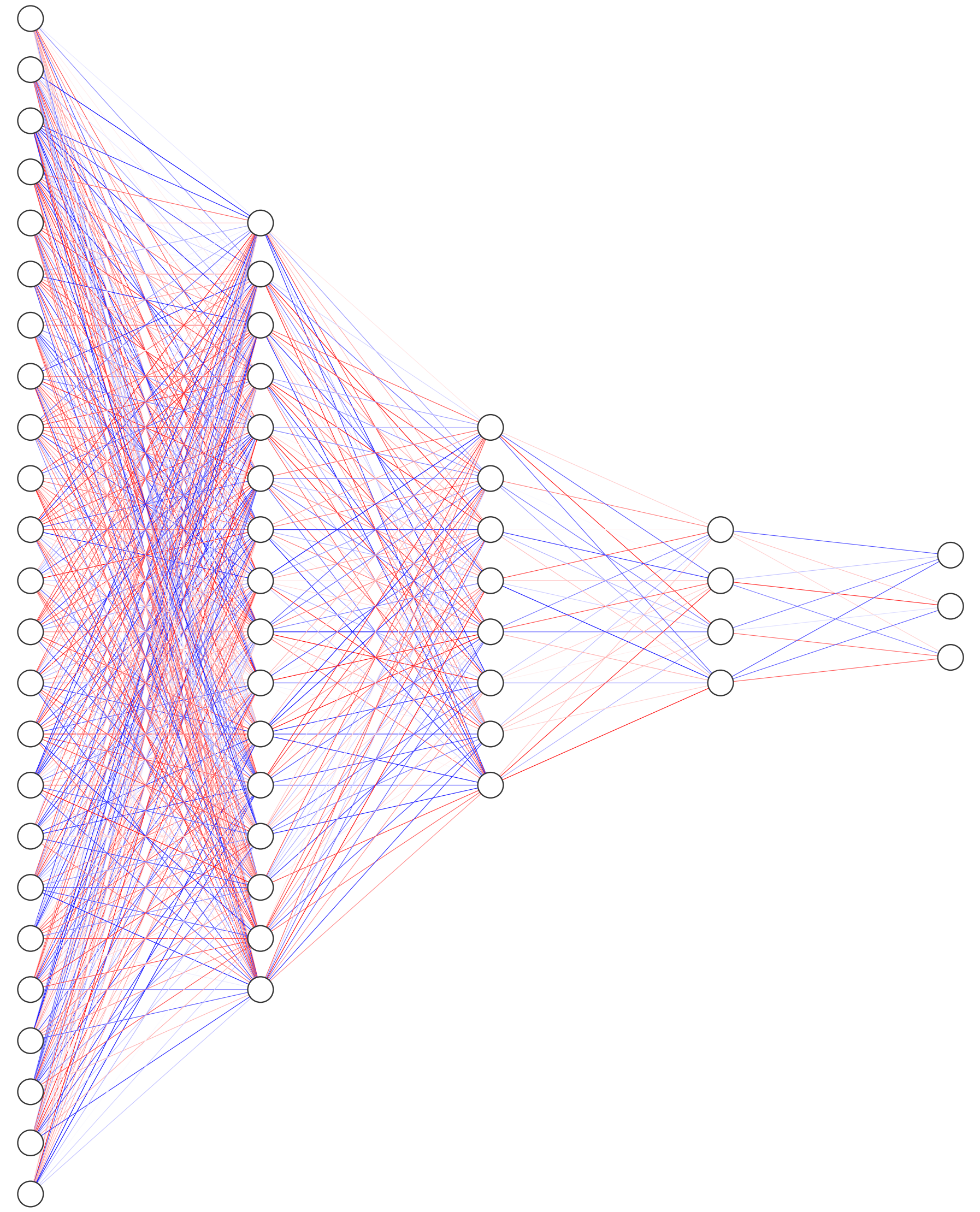
very negative

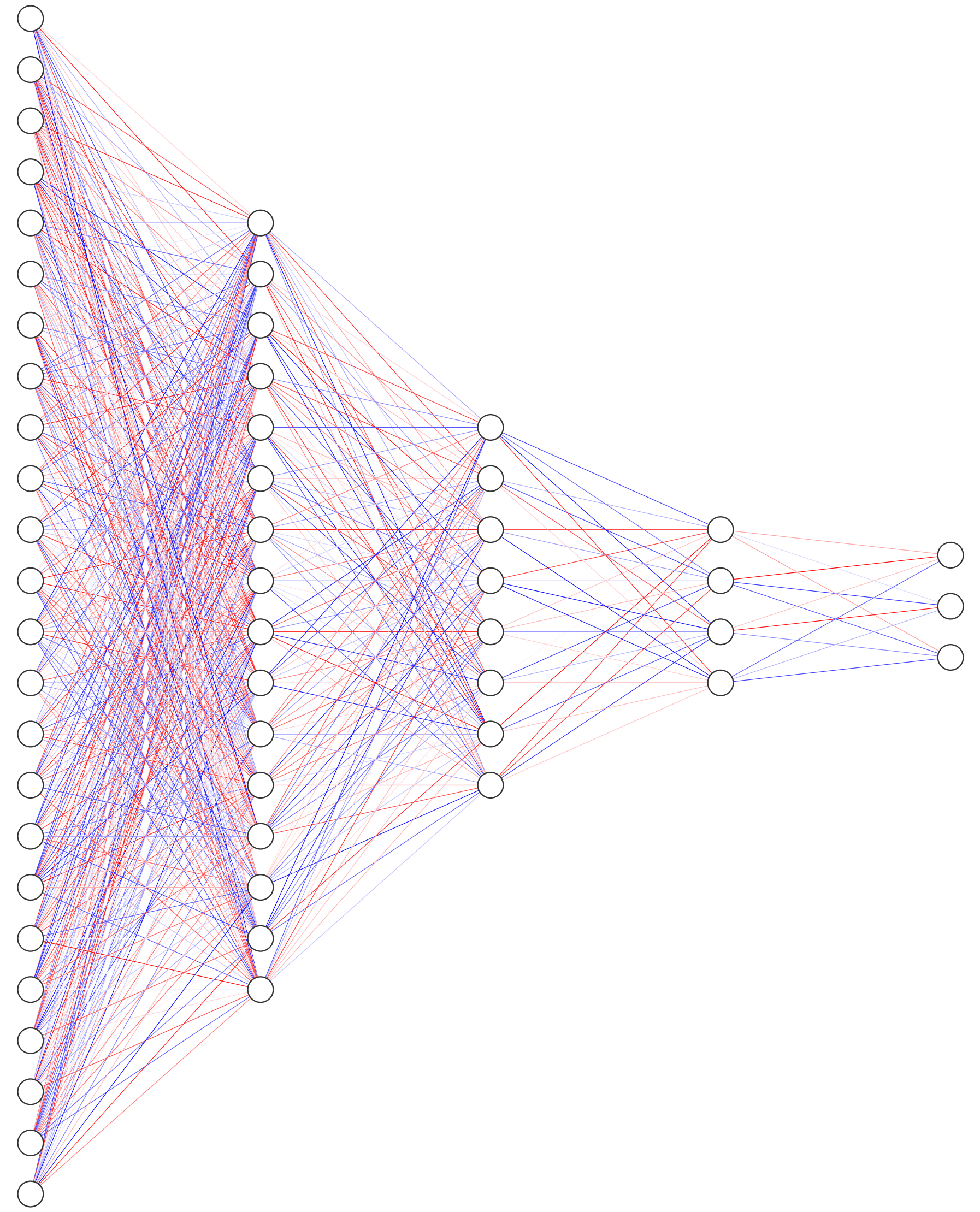


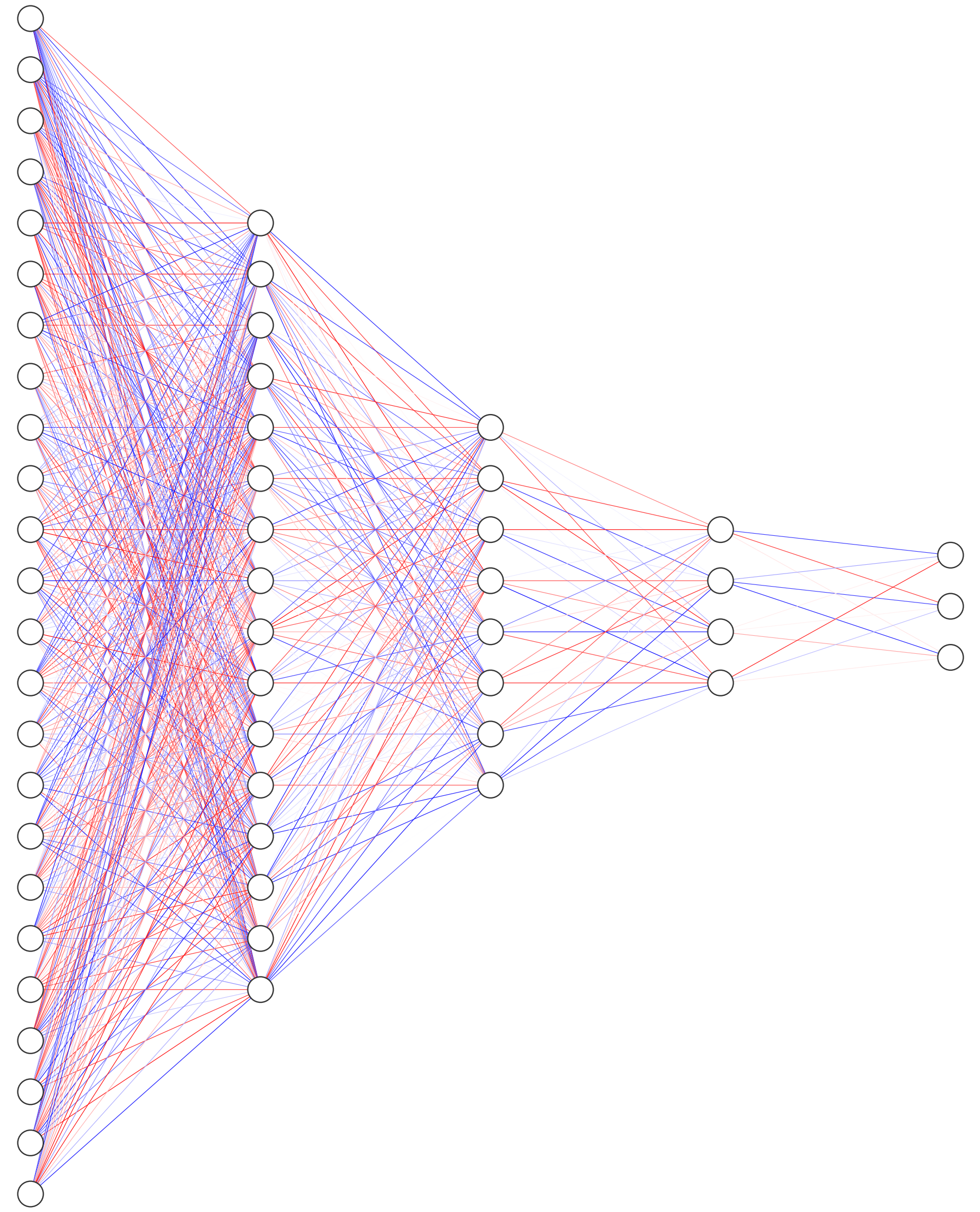


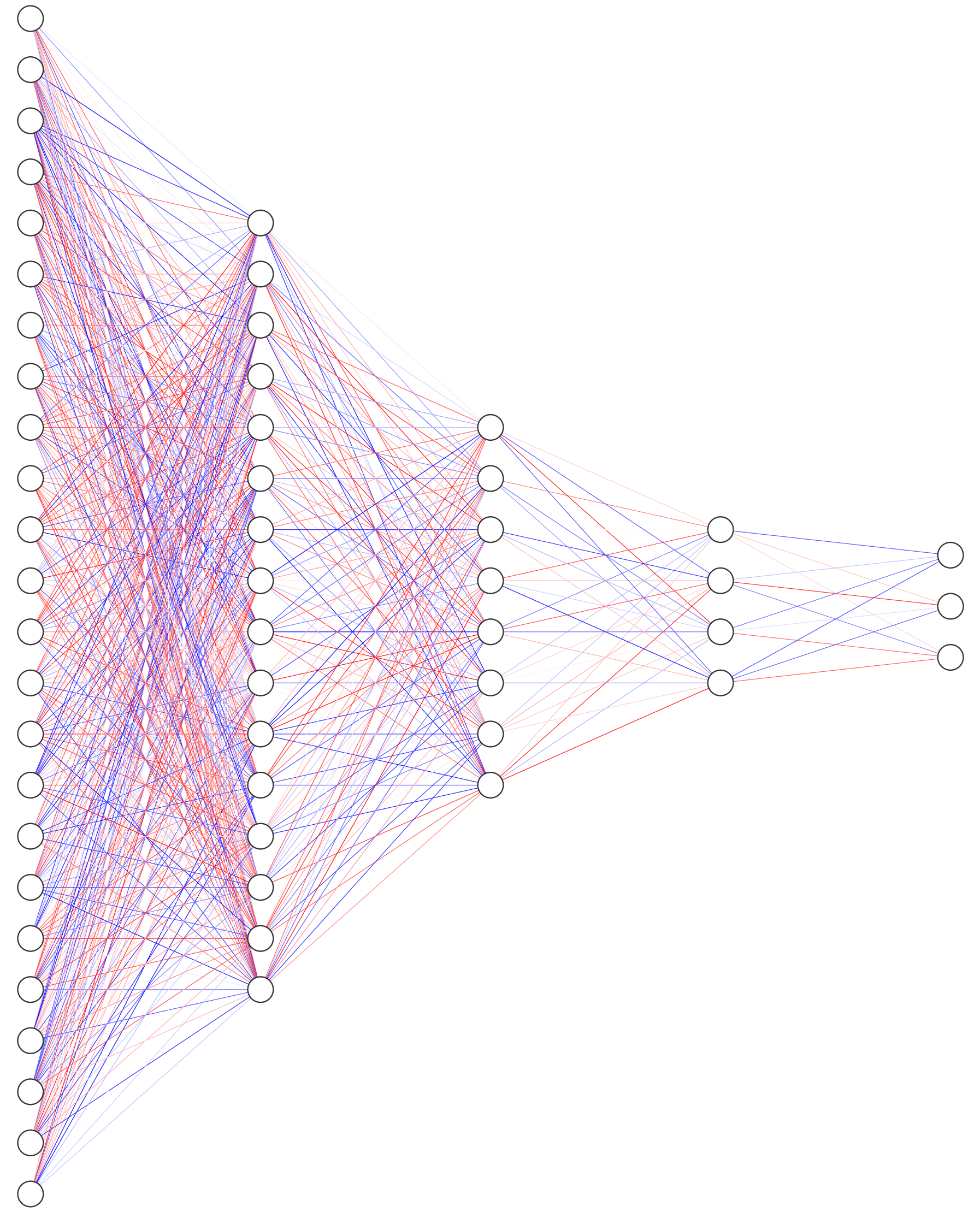


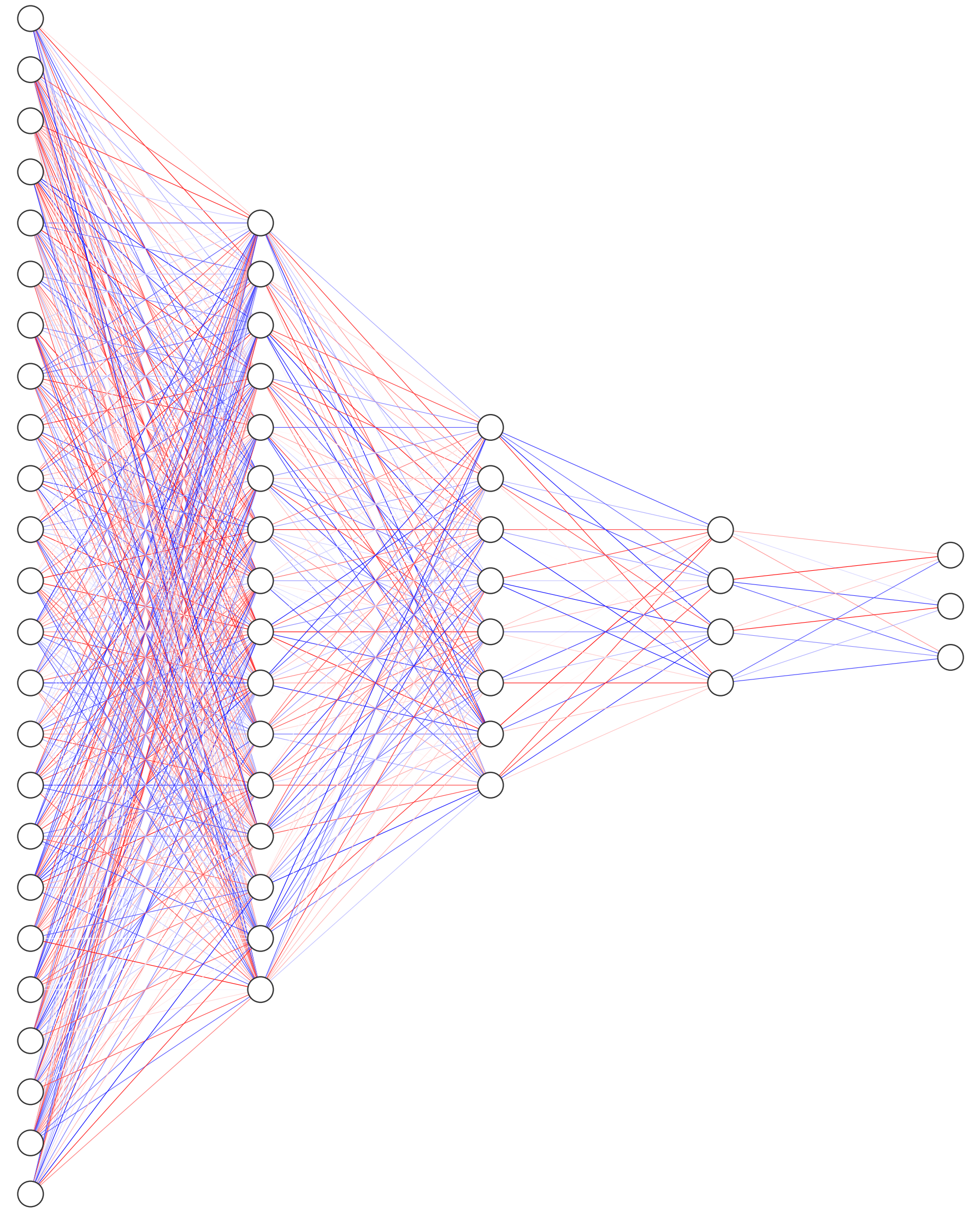


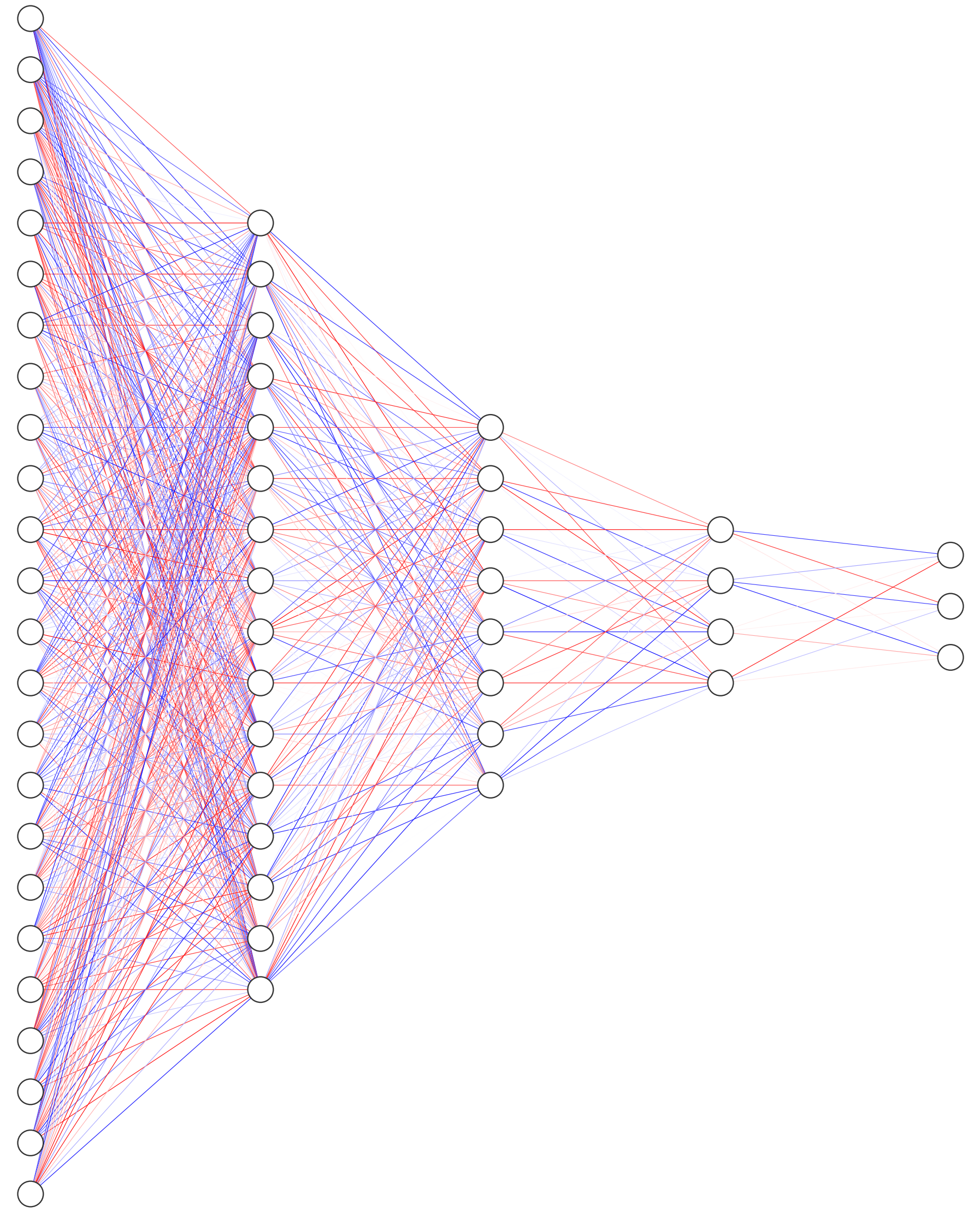


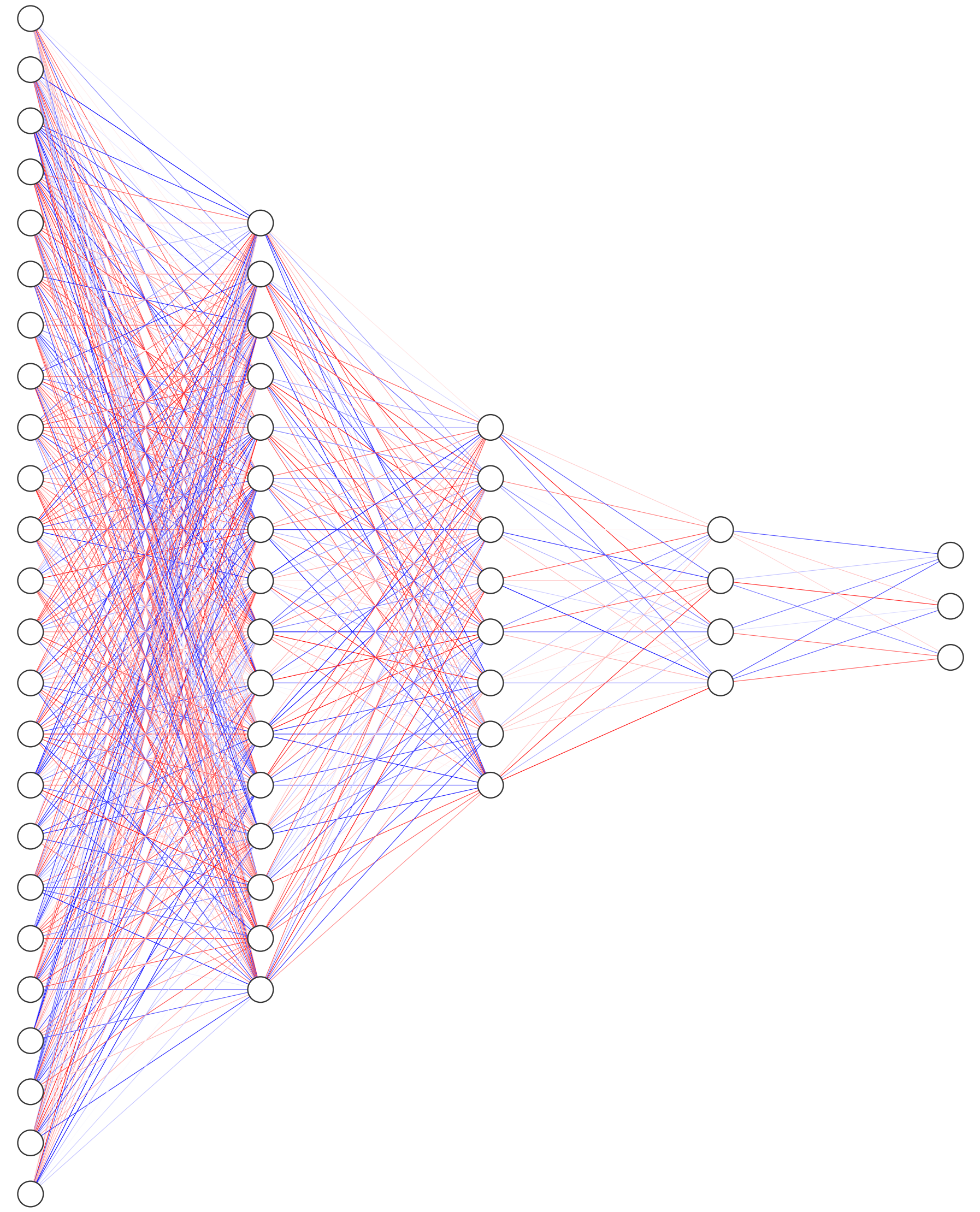


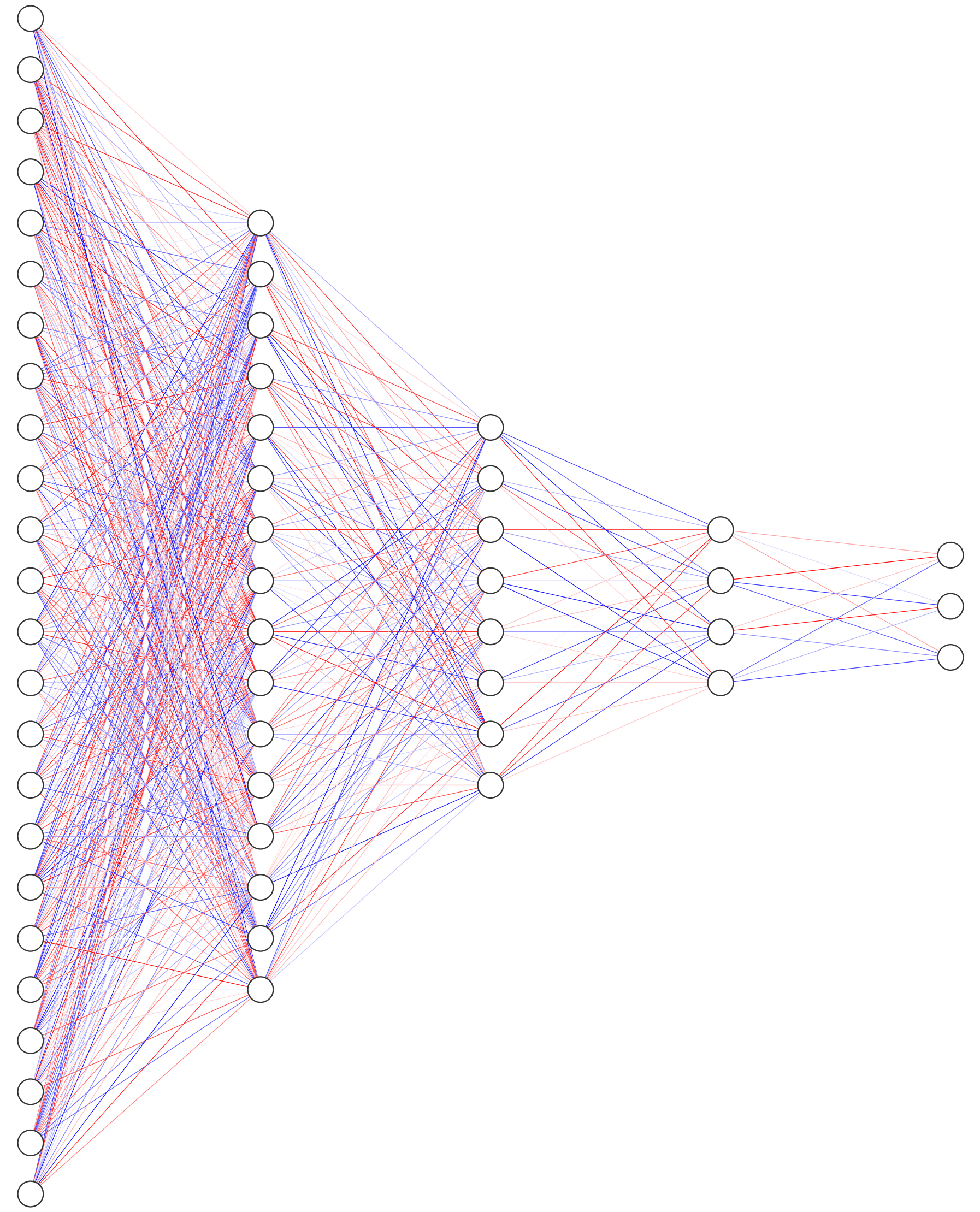


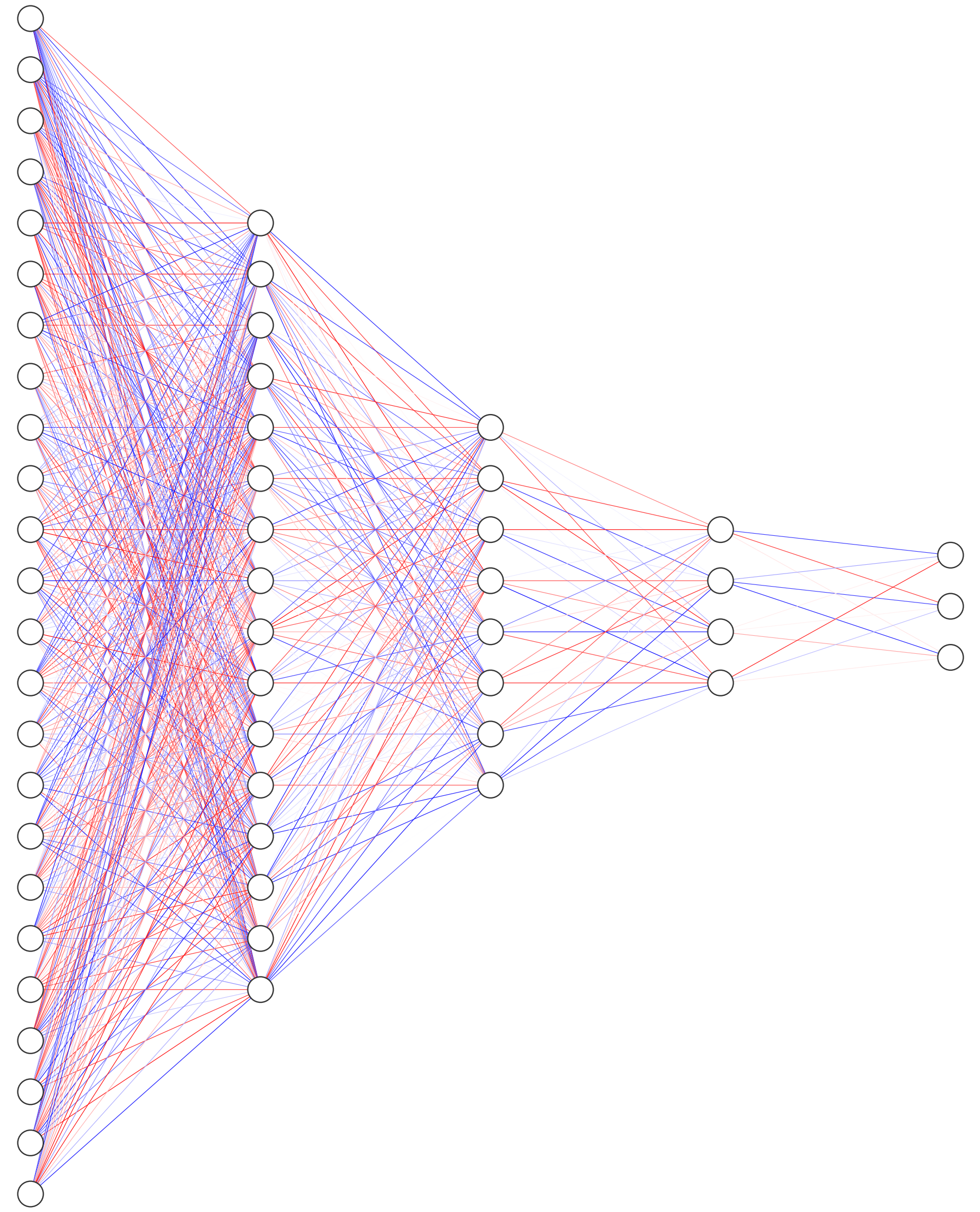


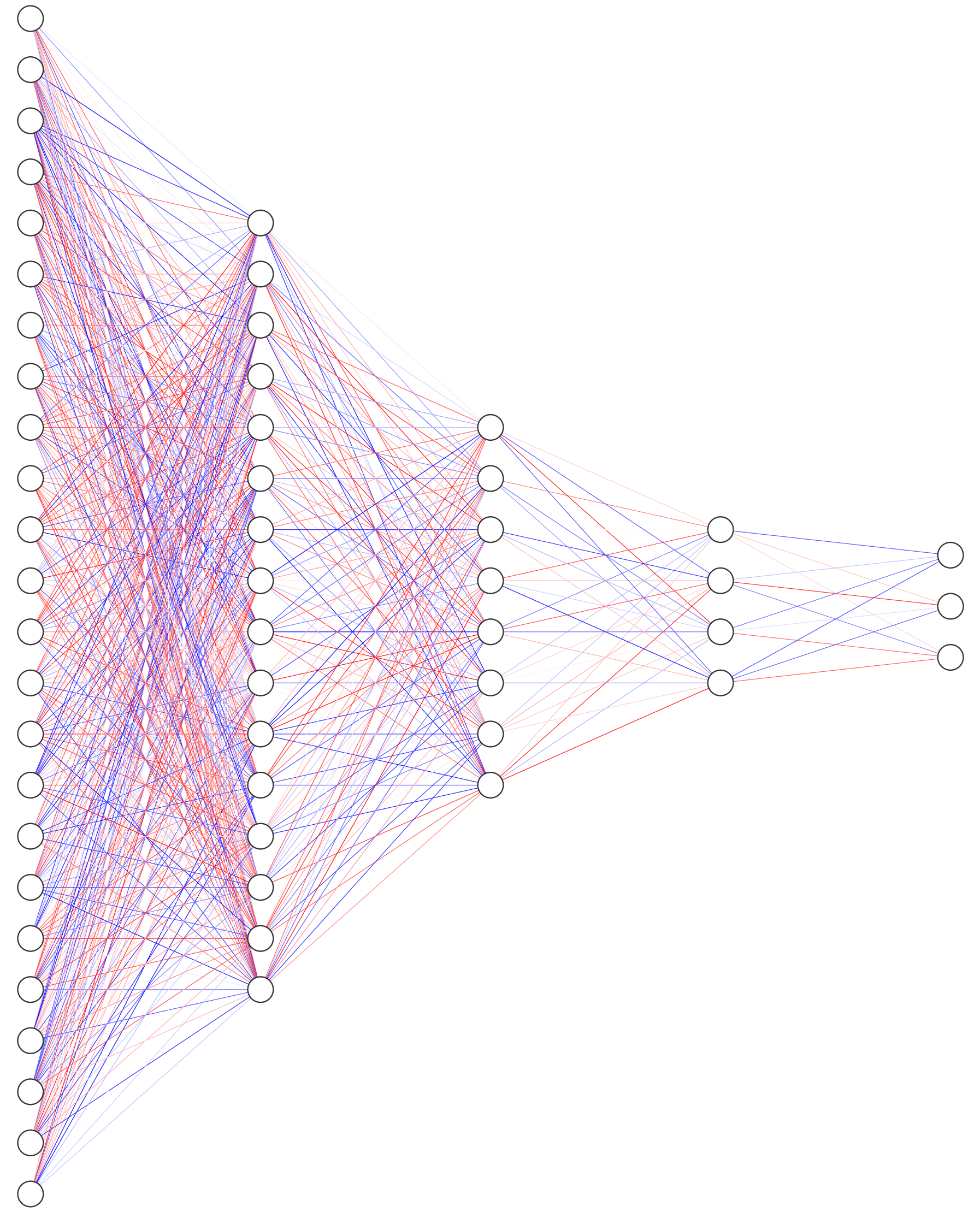


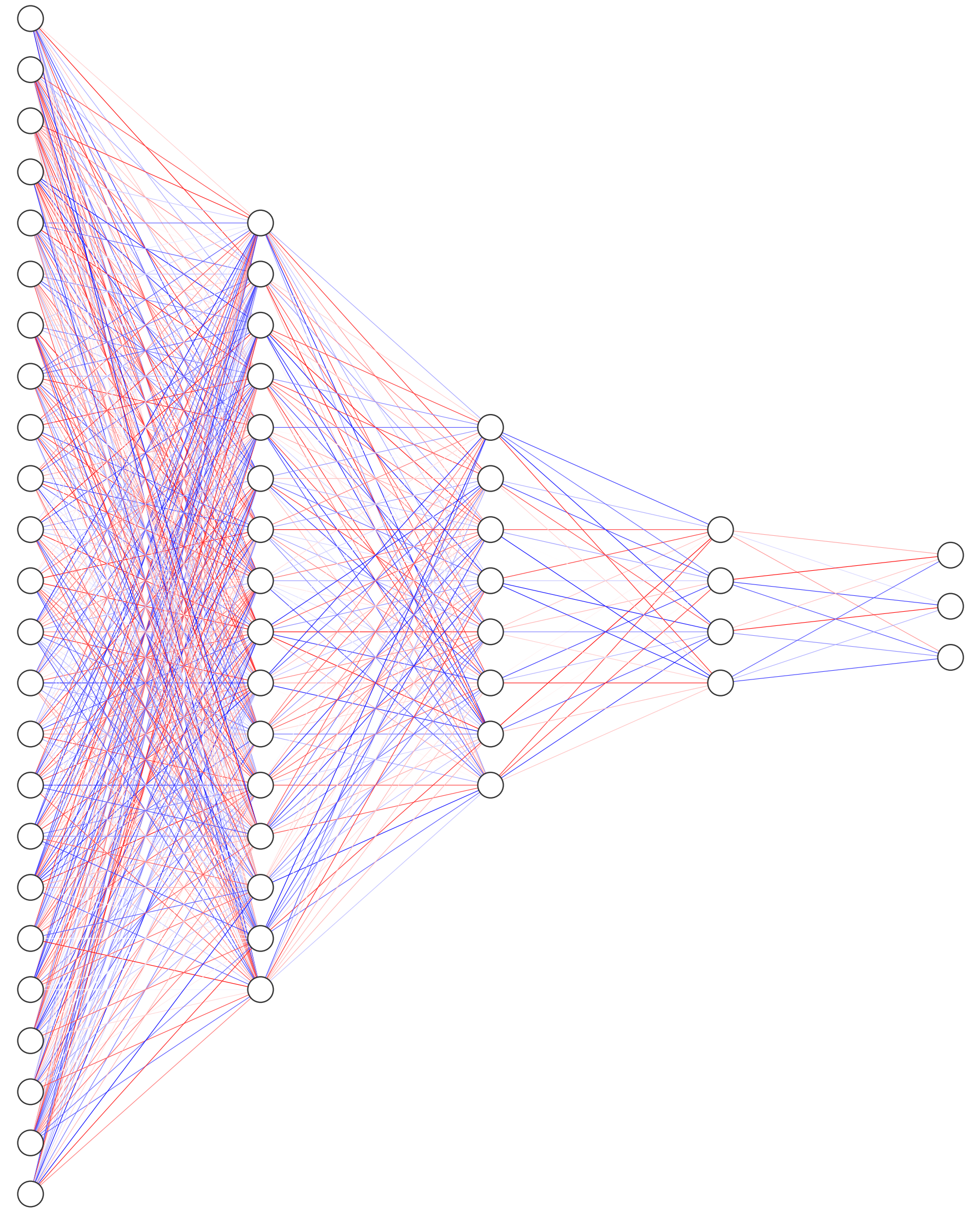


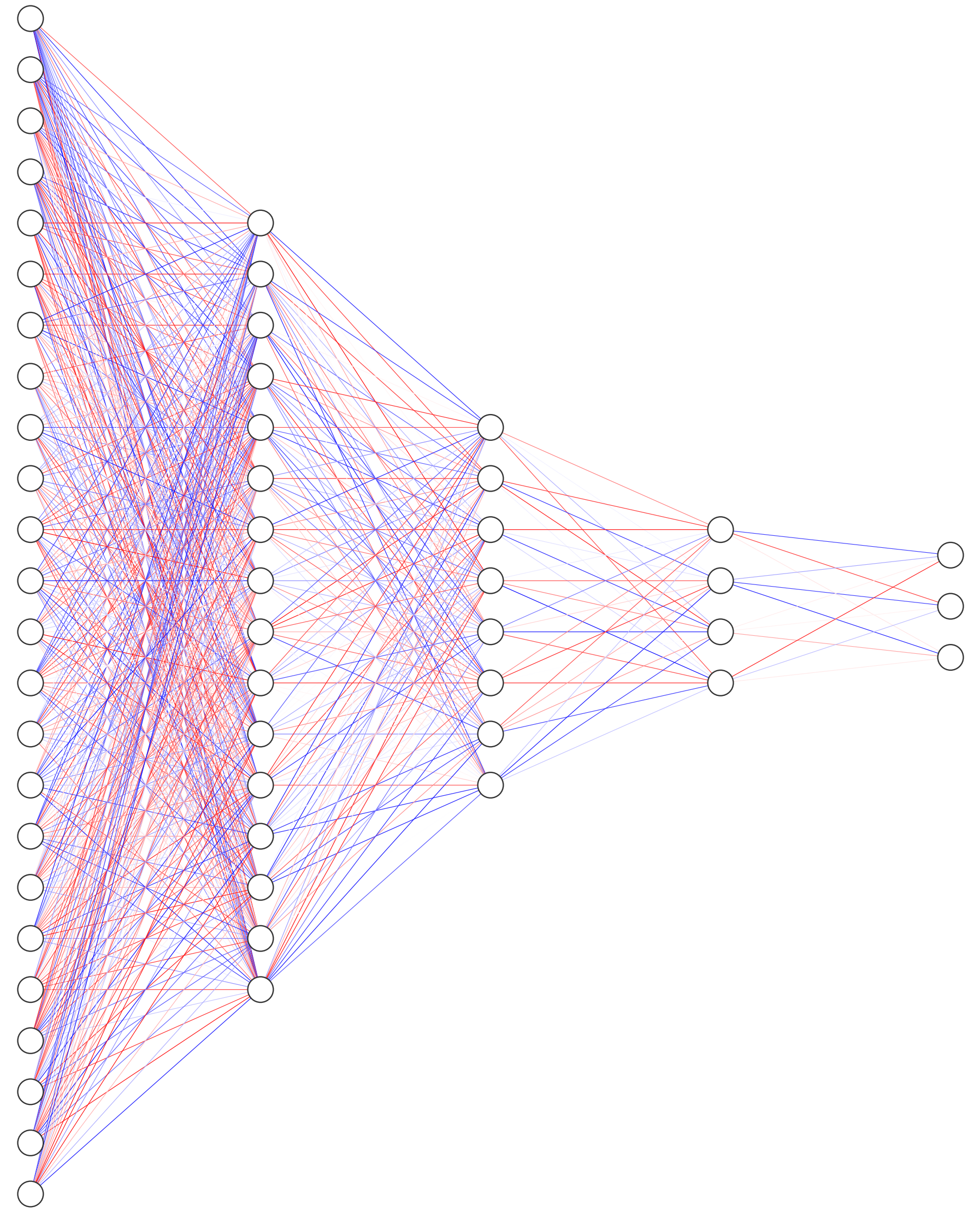


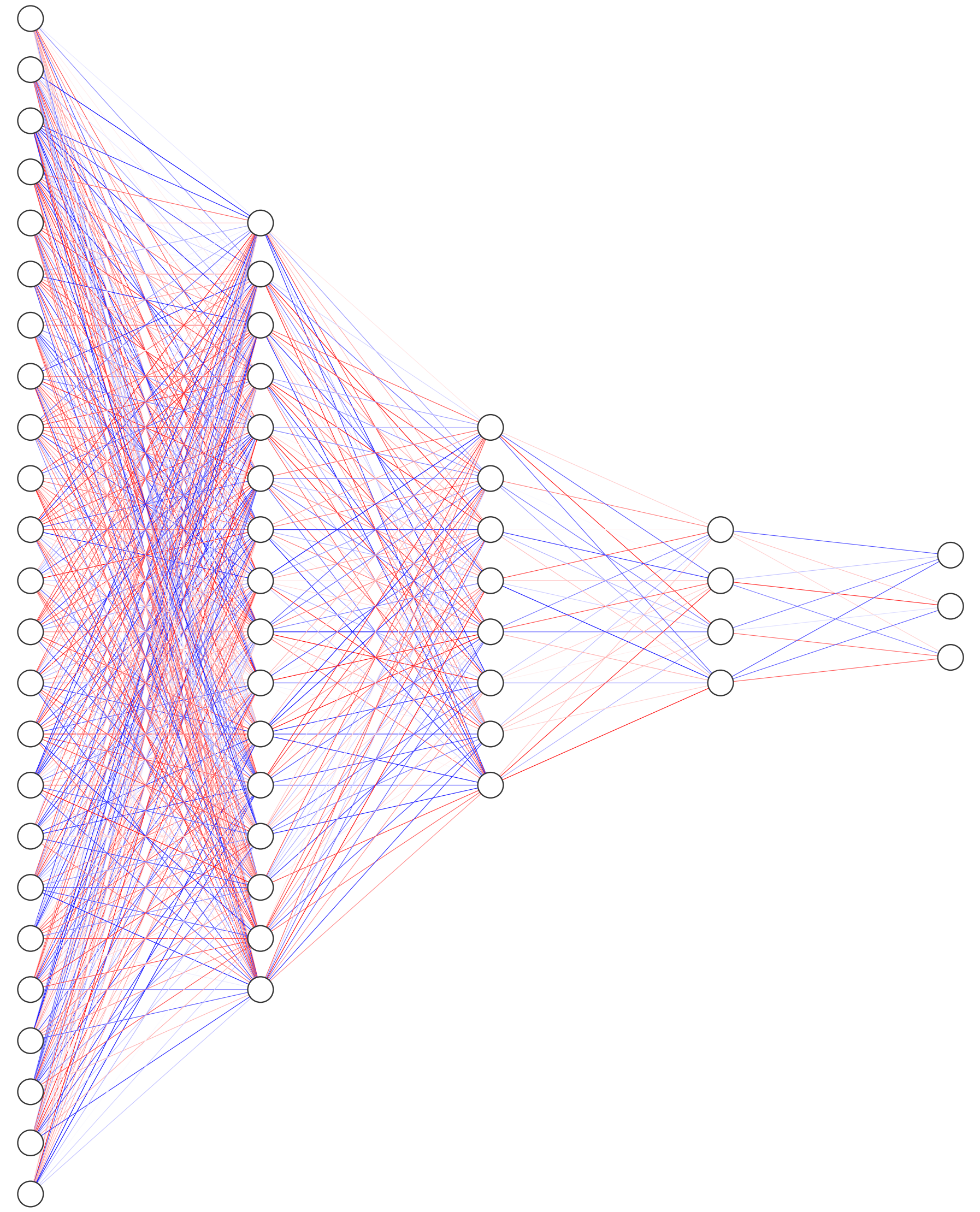


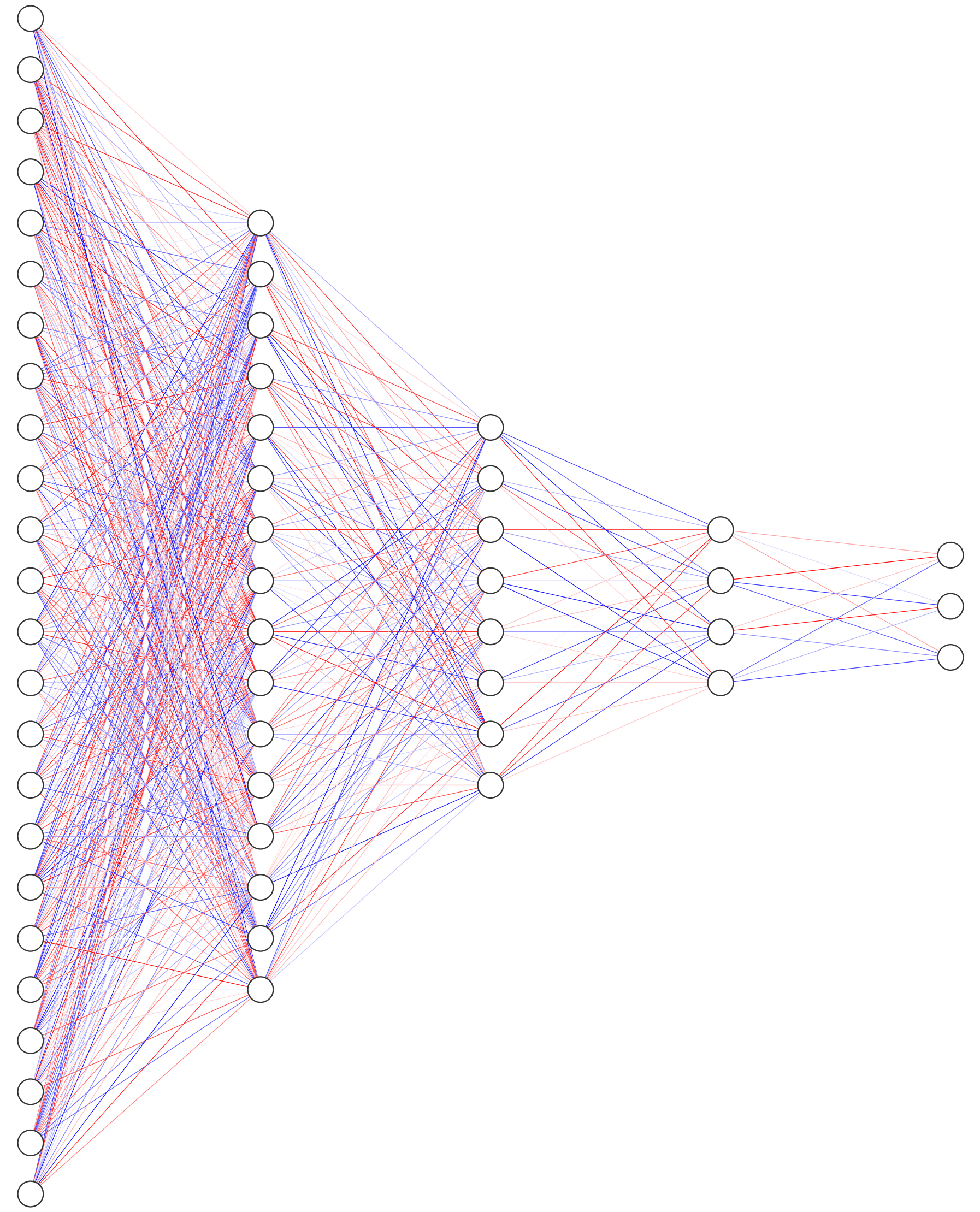


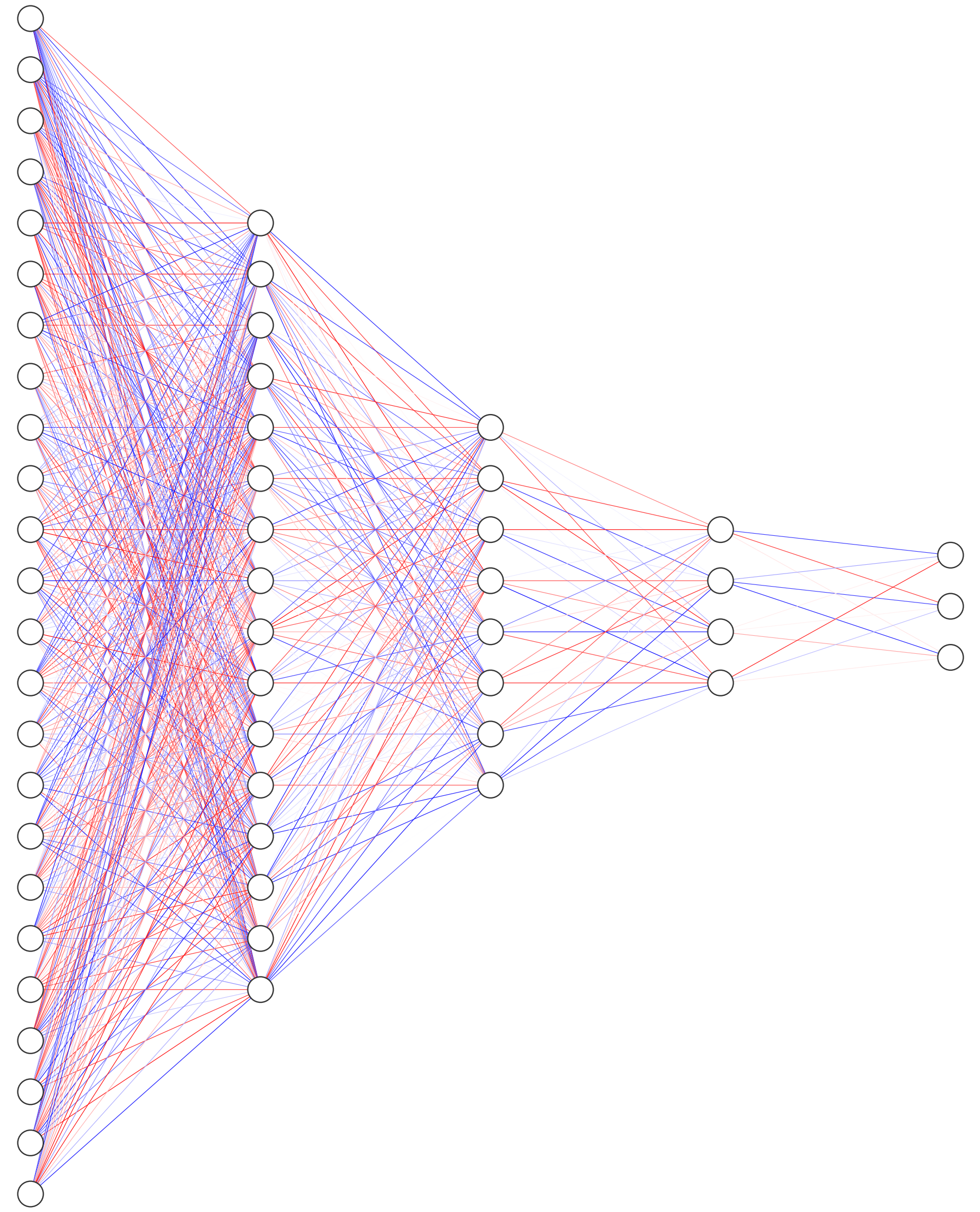


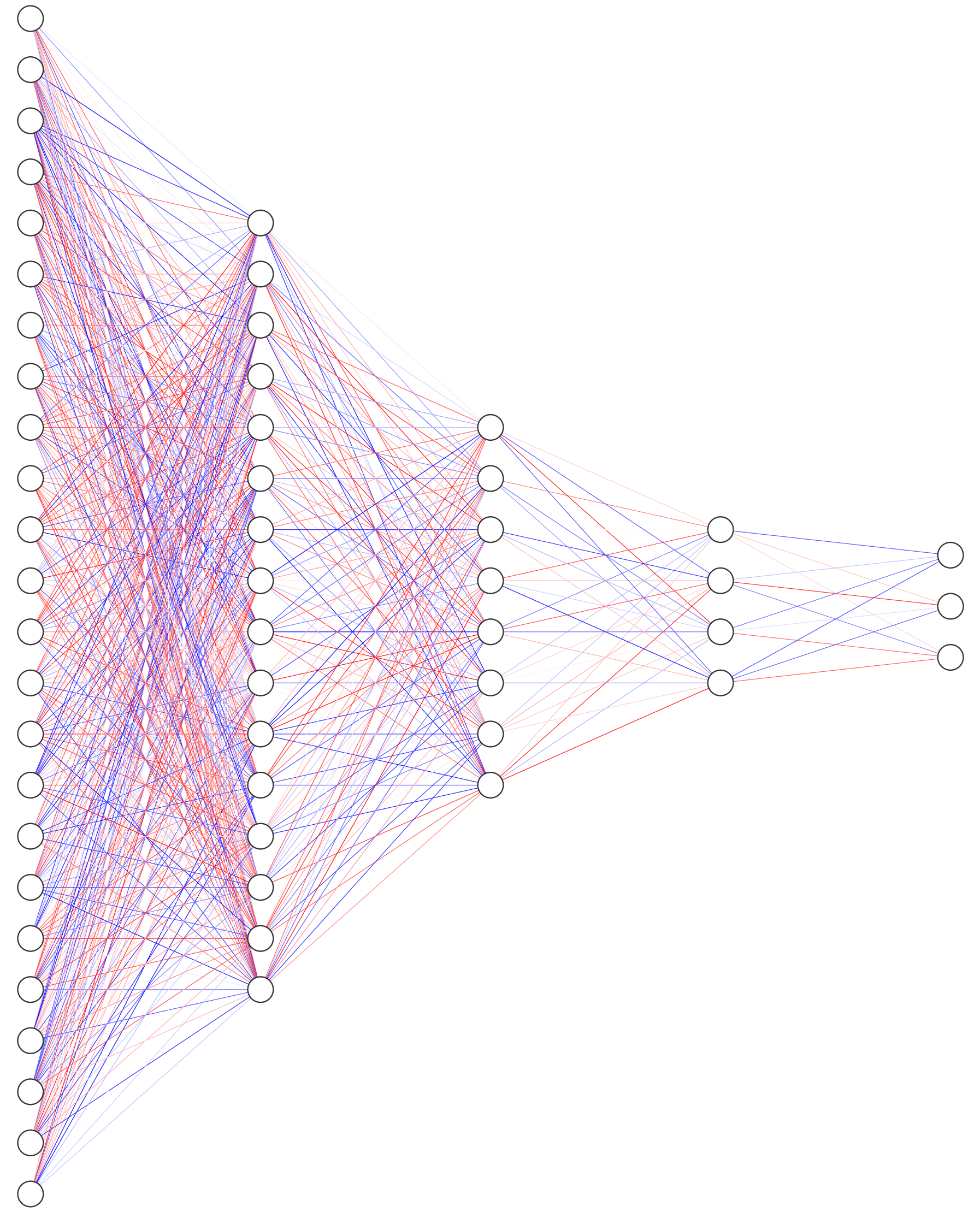


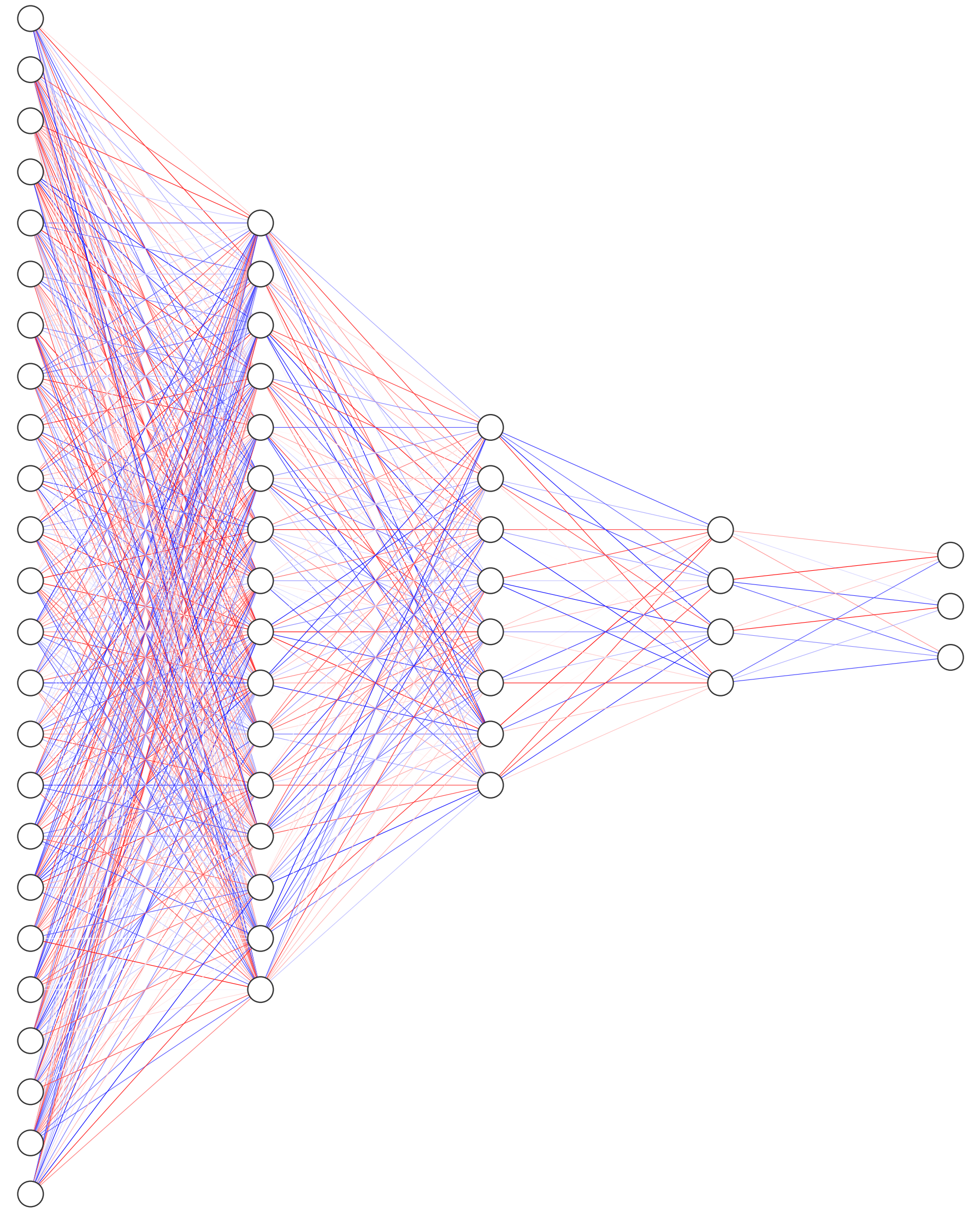




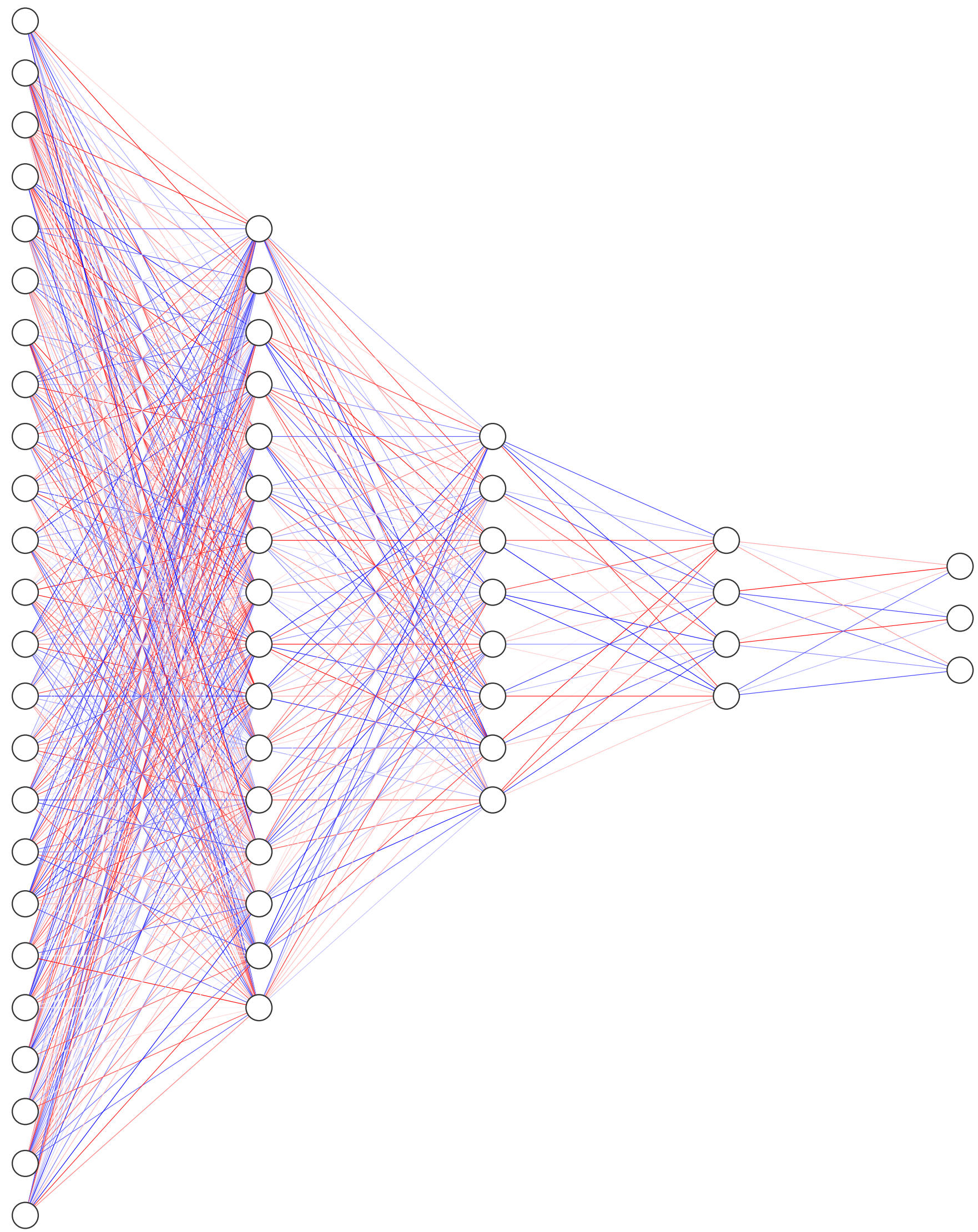




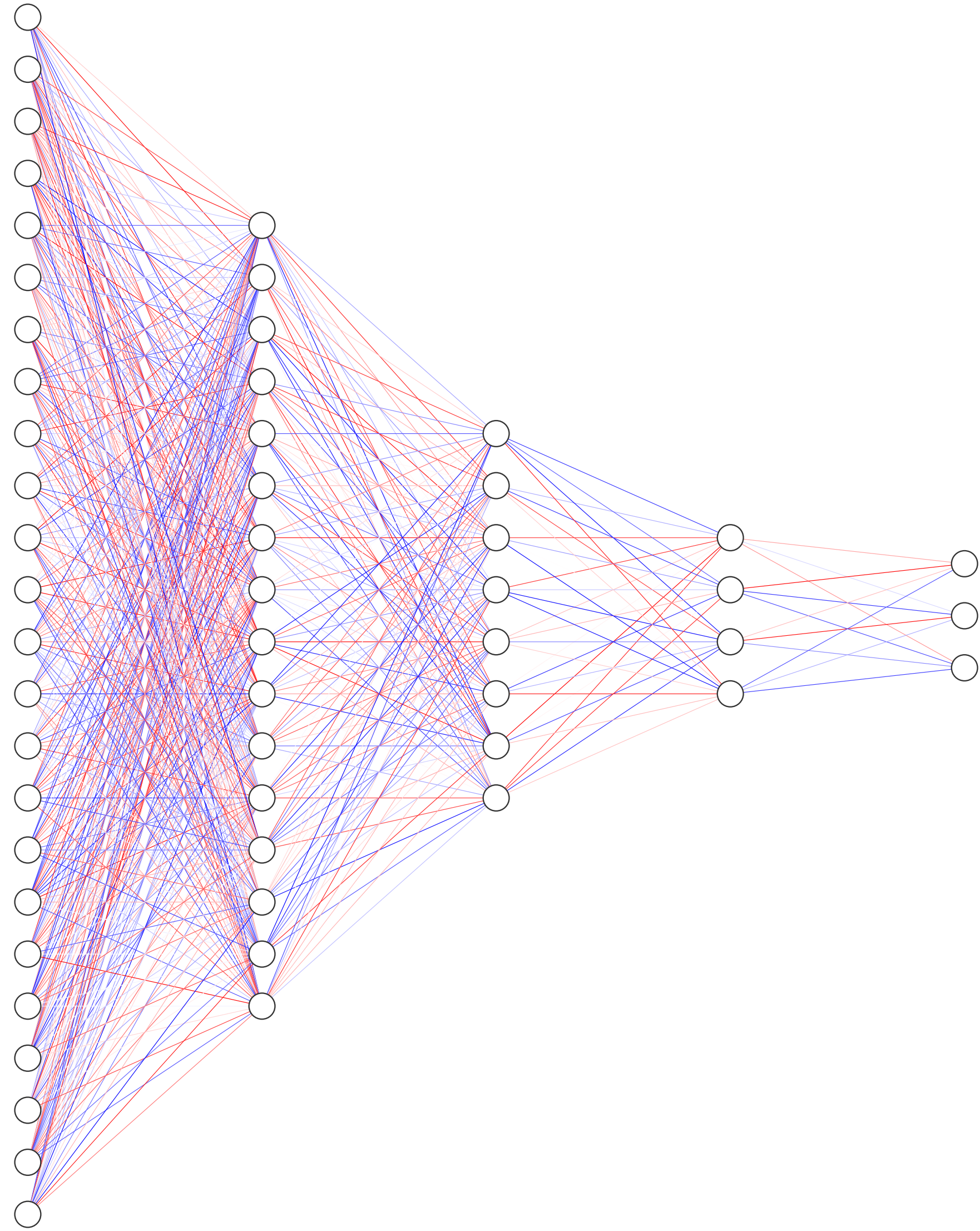




x'

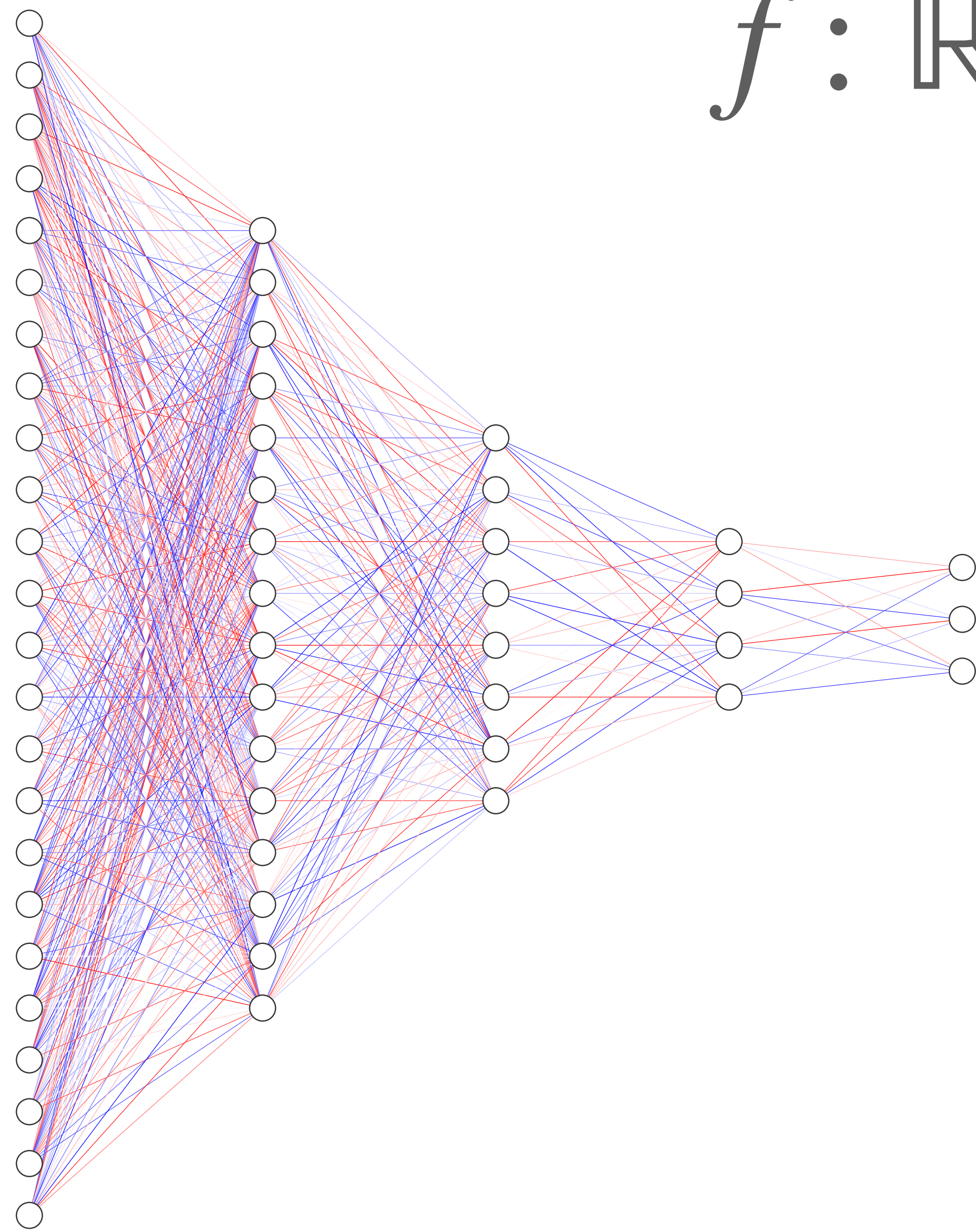


$y = [0.990, 0.009, 0.001]$



The parameters of a NN are its knowledge.



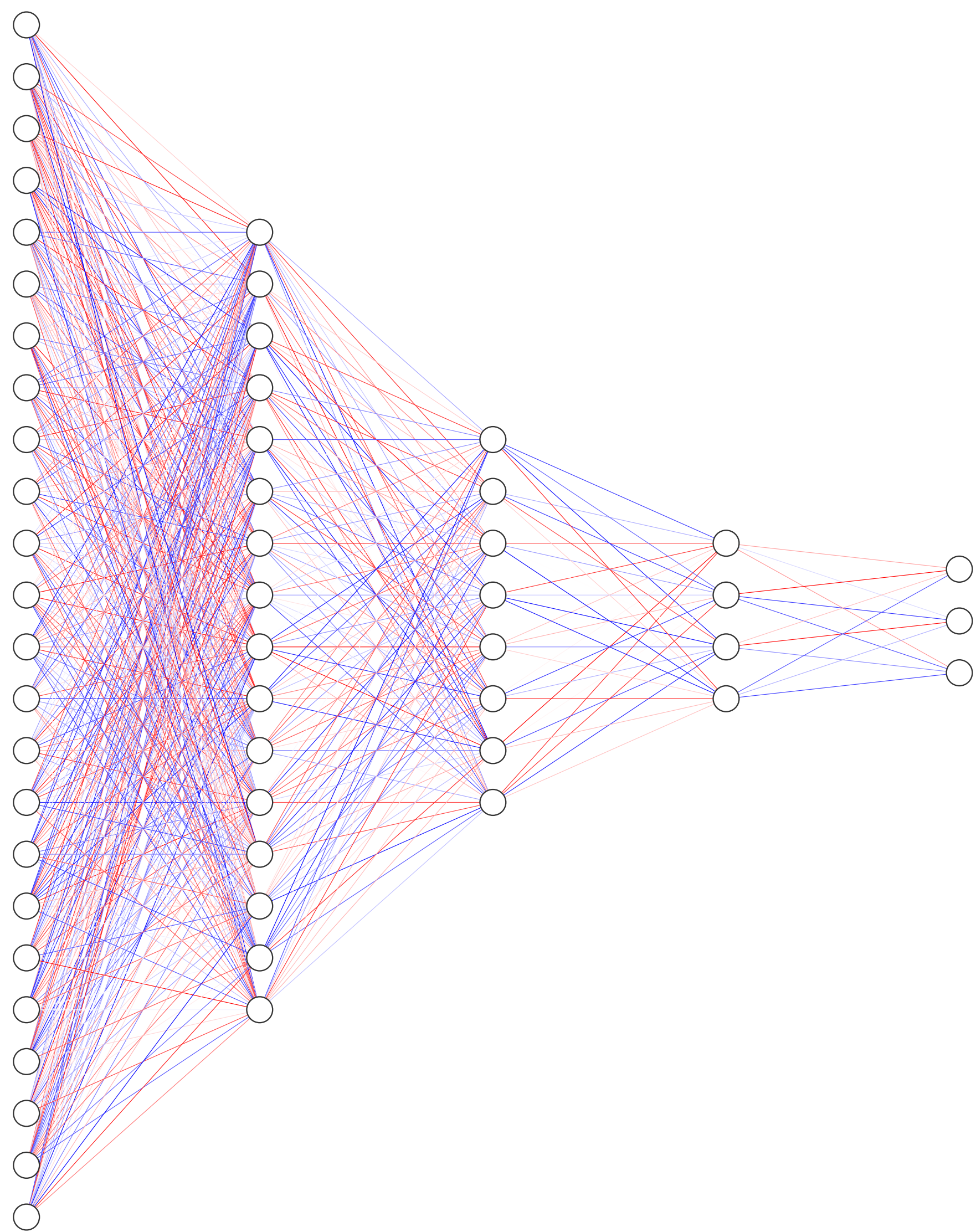


$$f : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f(x) \equiv \text{logits}$$

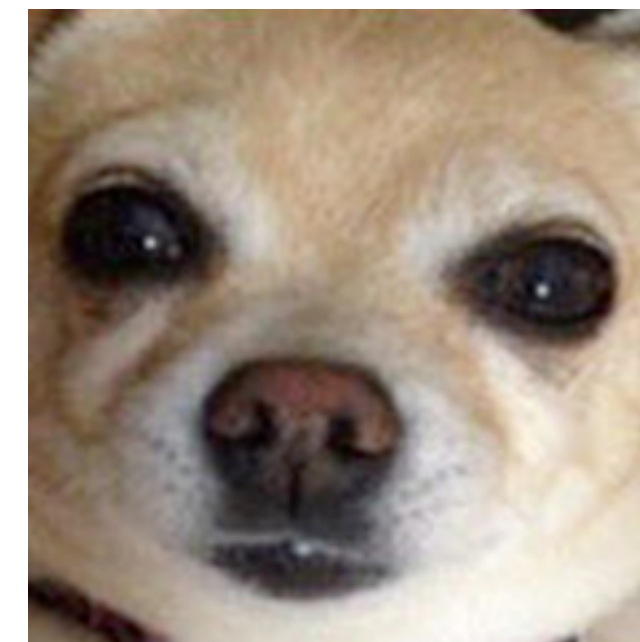
$$\phi(f(x)_i) \equiv \frac{\exp(f(x)_i)}{\sum_j \exp(f(x)_j)} \equiv \text{softmax}$$

$$\min_f \mathbb{E}_{x \in X} (\phi(f(x)) - y)^2$$

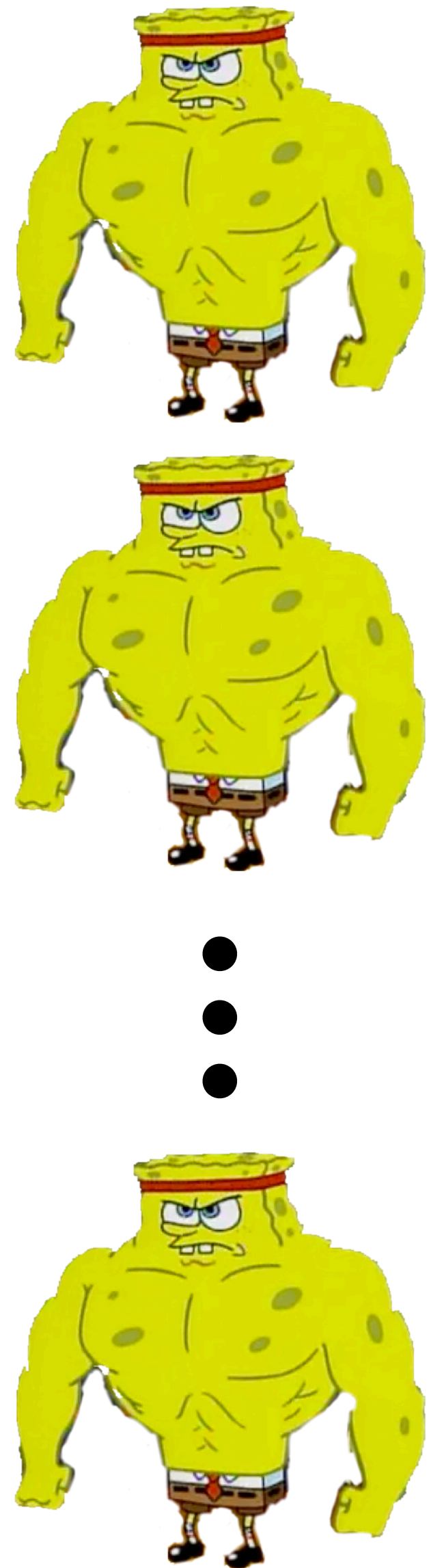


Knowledge

$$f(x) \equiv \text{logits}$$



Trained



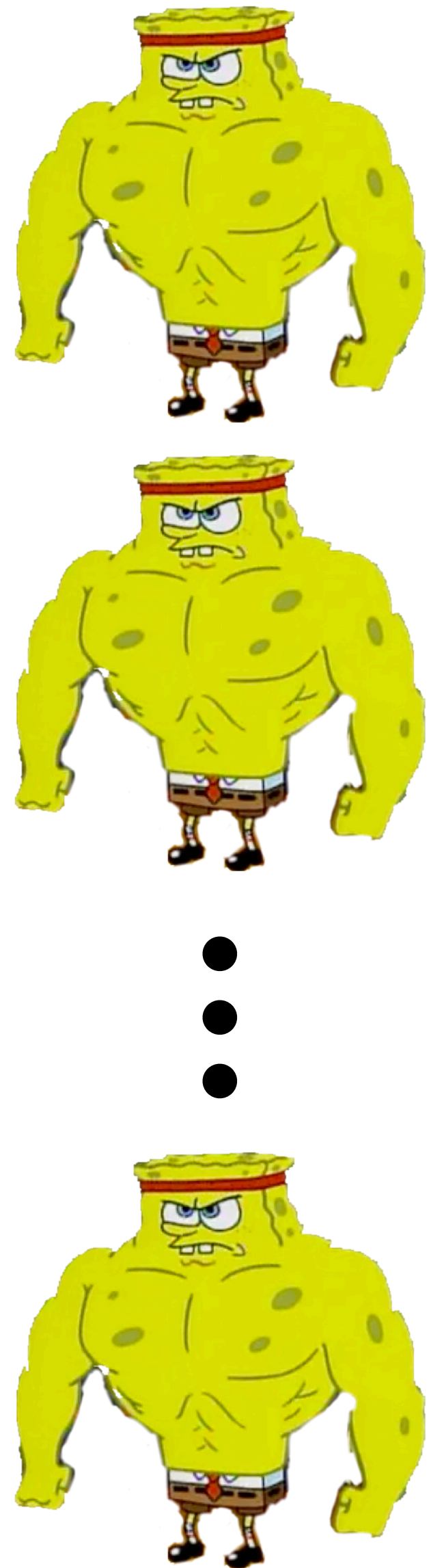
$\bar{f}(x)$

-

$g(x)$



Trained



Model Compression

<https://doi.org/10.1145/1150402.1150464>

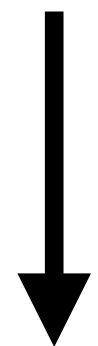
$(\bar{f}(x))$

-

$g(x))^2$



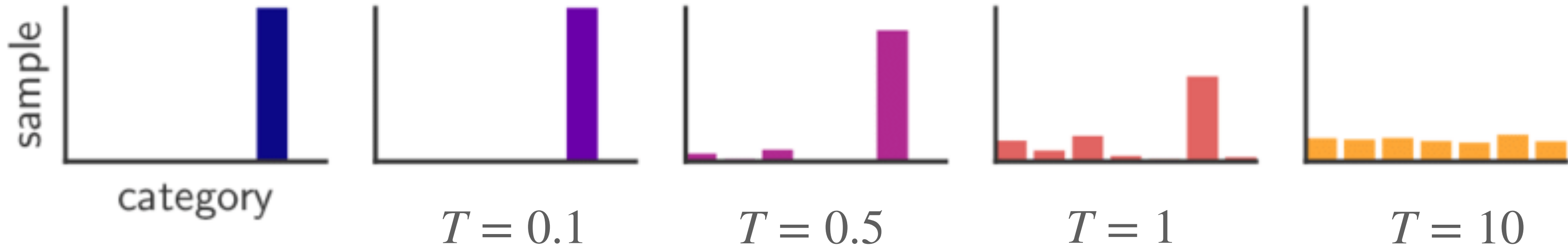
$$\min_g \mathbb{E}(\bar{f}(x) - g(x))^2$$



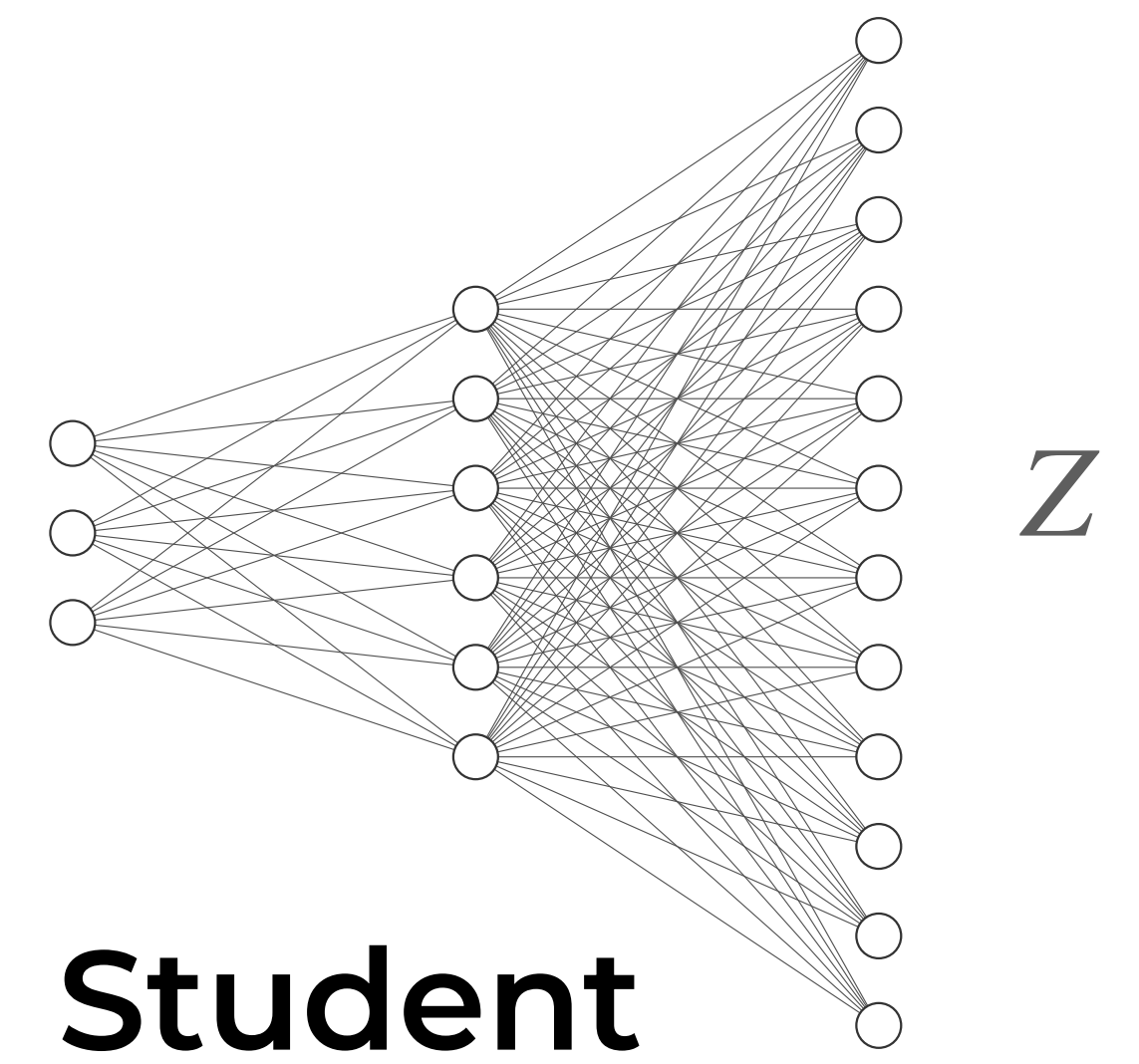
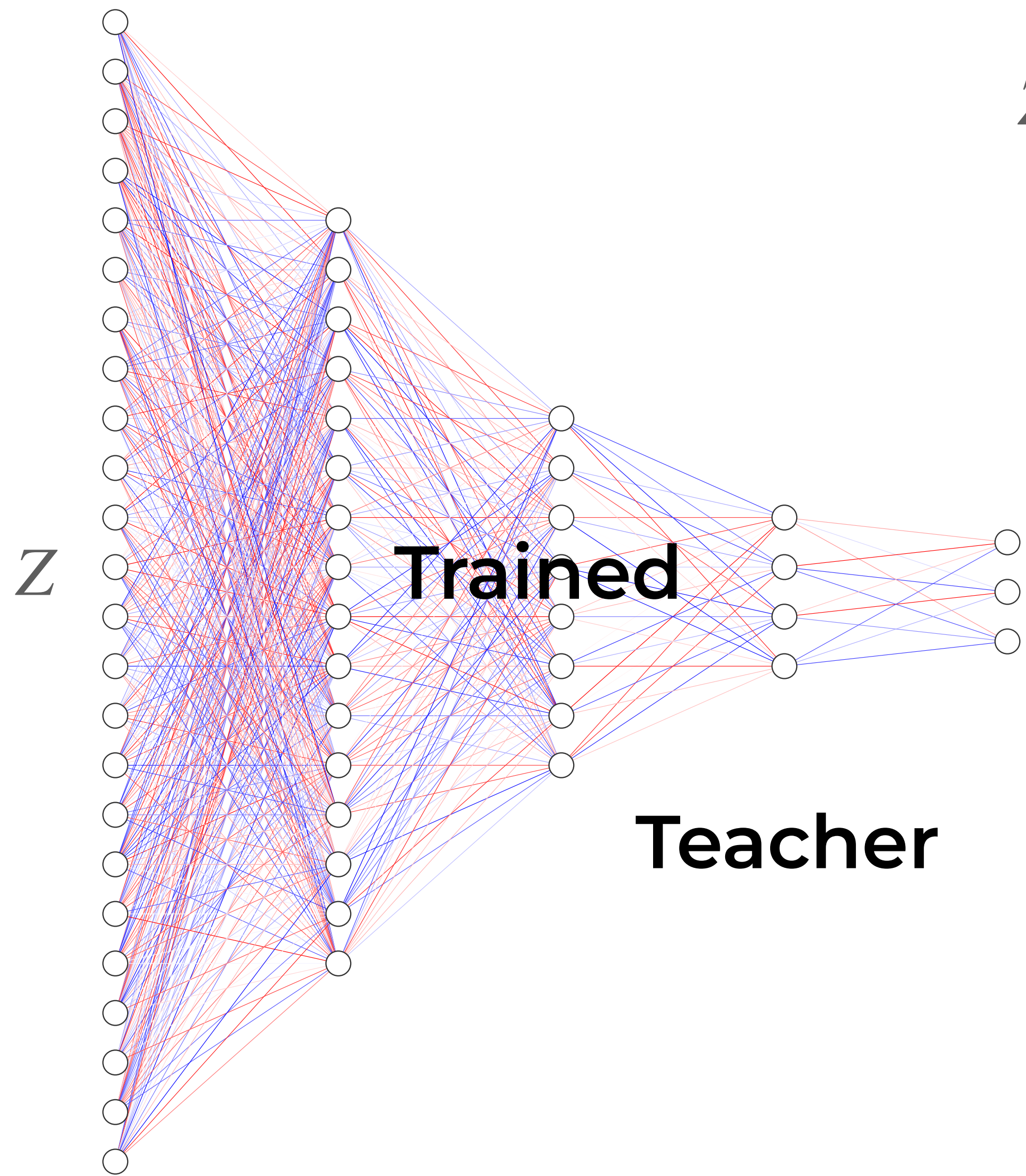
Knowledge Distillation

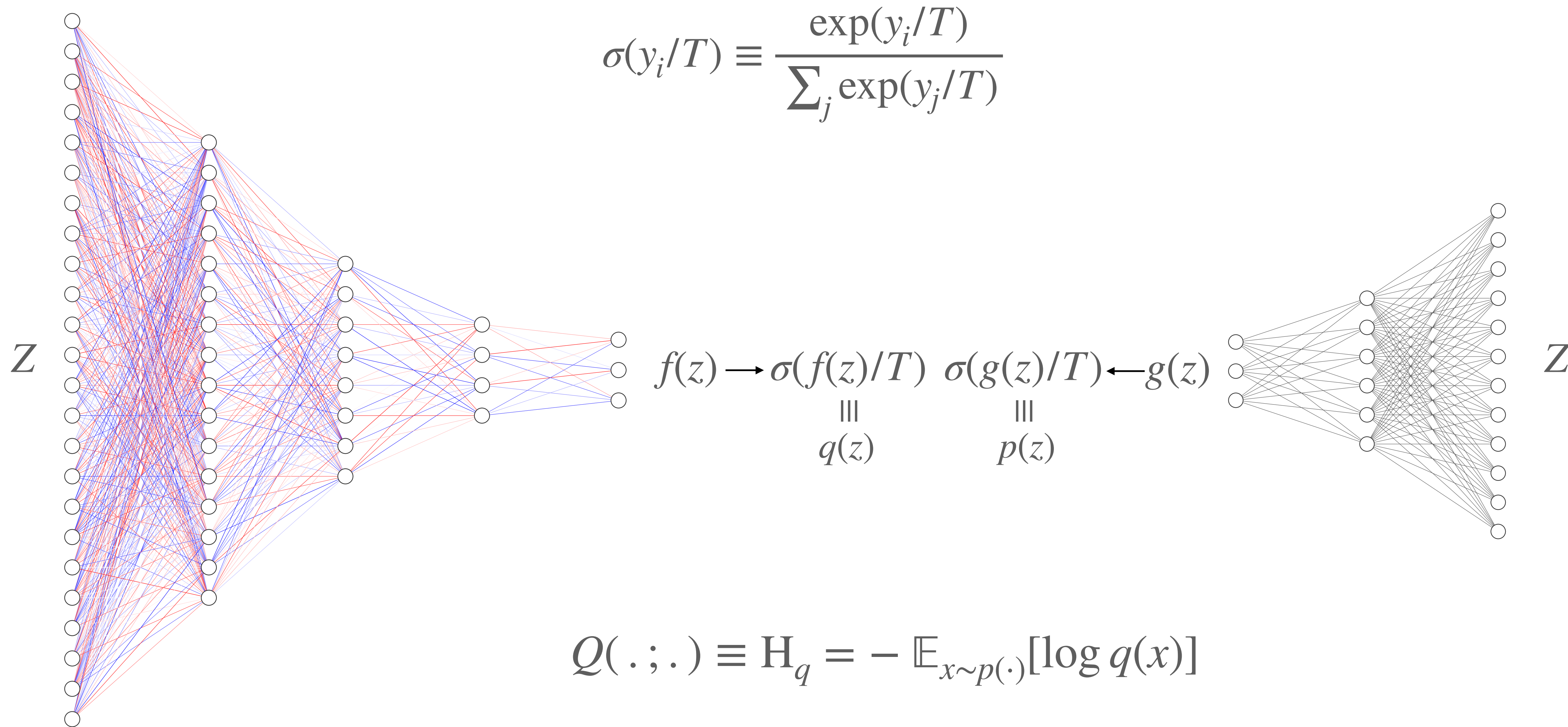


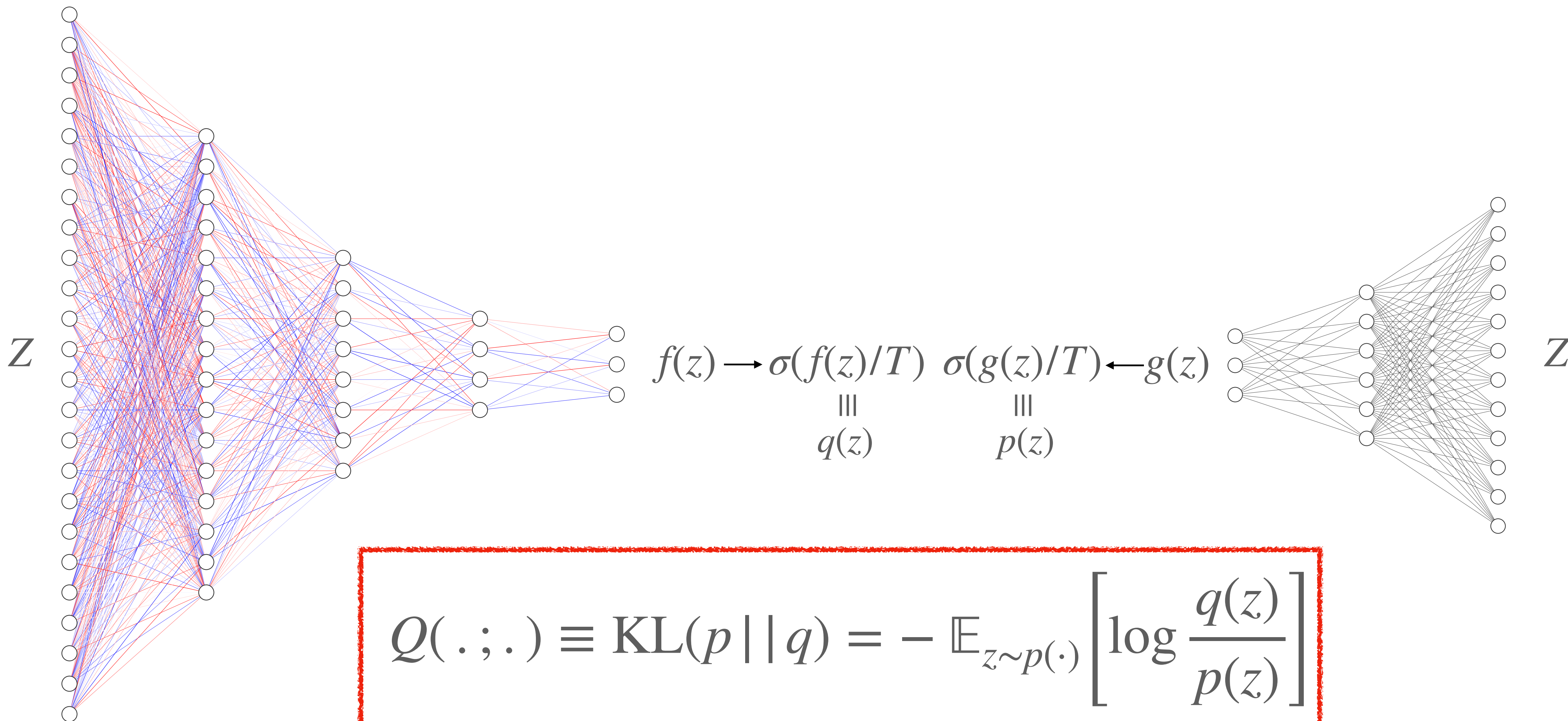
$$\sigma(f(x)_i/T) \equiv \frac{\exp(f(x)_i/T)}{\sum_j \exp(f(x)_j/T)} \xleftrightarrow{Q(\cdot; \cdot)} \sigma(g(x)_i/T) \equiv \frac{\exp(g(x)_i/T)}{\sum_j \exp(g(x)_j/T)}$$



$Z = \text{Transfer Dataset}$







Quick! Answer the following questions ASAP!

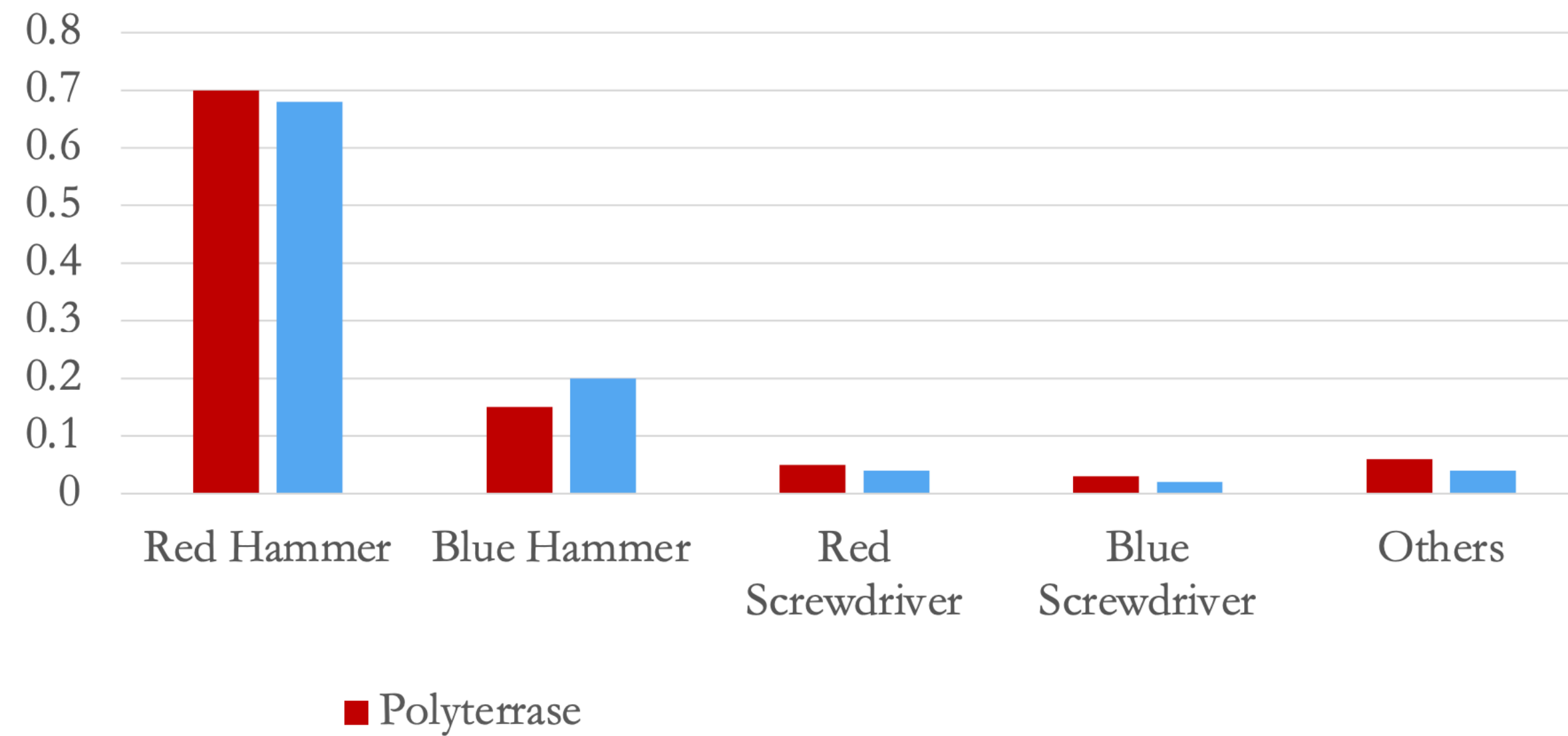
- ✓ $32 + 6$
- ✓ $2 * 125$
- ✓ $12 * 12$
- ✓ $25 * 51$
- ✓ $100*11 + 22$

Quick! Answer the following questions ASAP!

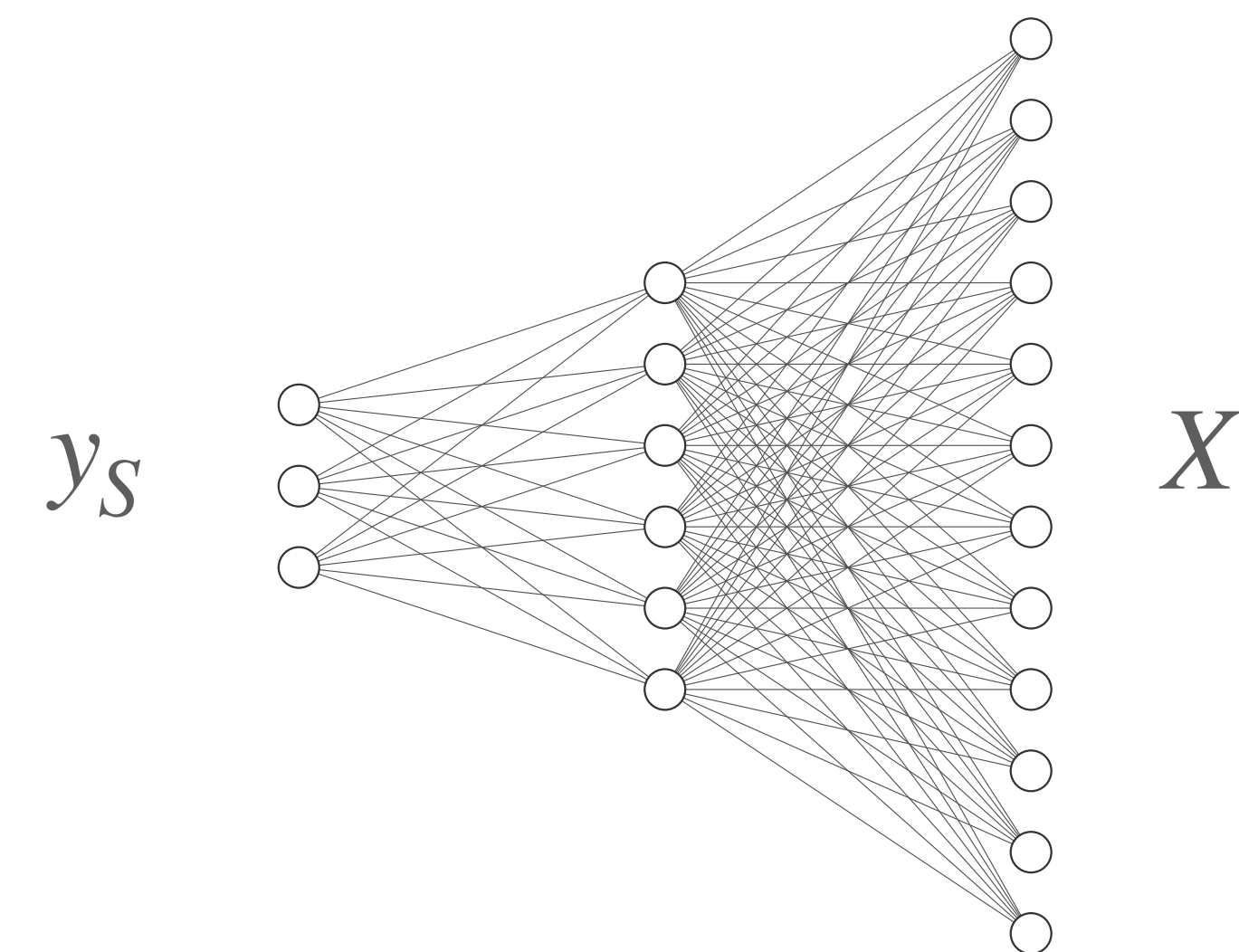
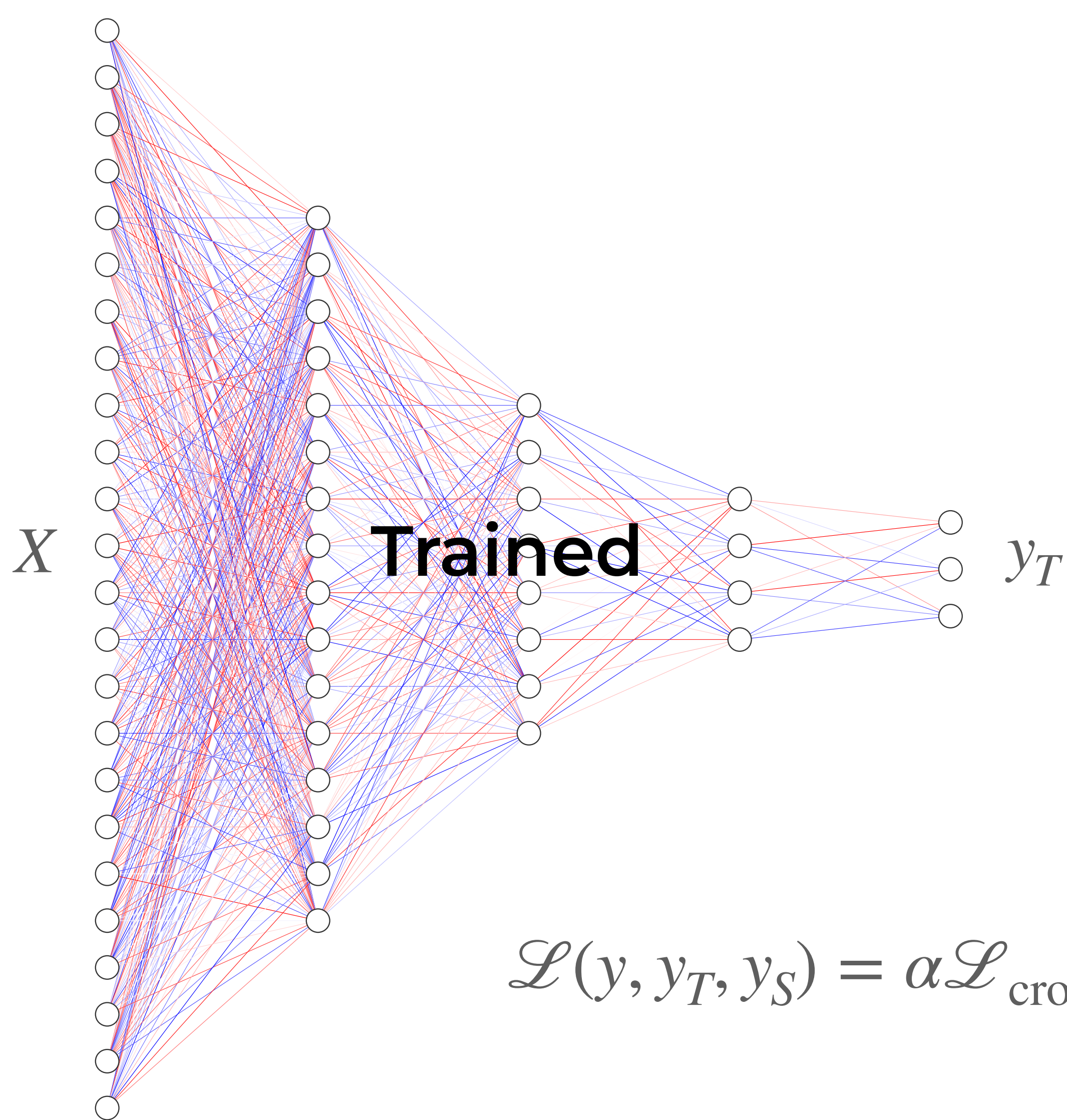
- ✓ $32 + 6$
- ✓ $2 * 125$
- ✓ $12 * 12$
- ✓ $25 * 51$
- ✓ $100*11 + 22$

Done? Think about a **COLOUR** and a **TOOL!**

What tool did people think about?



$$\text{KL}(p || q) = - \mathbb{E}_{x \sim p(\cdot)} \left[\log \frac{q(x)}{p(x)} \right]$$



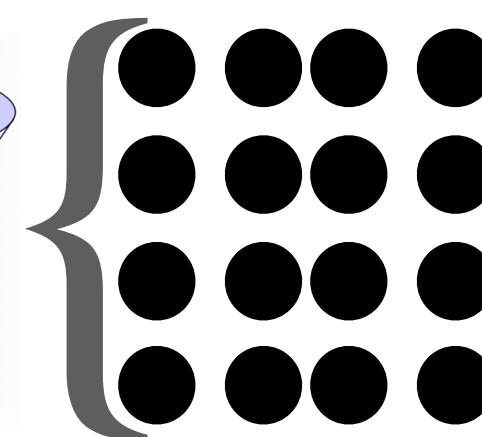
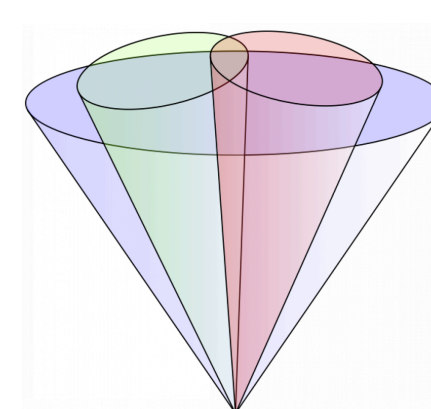
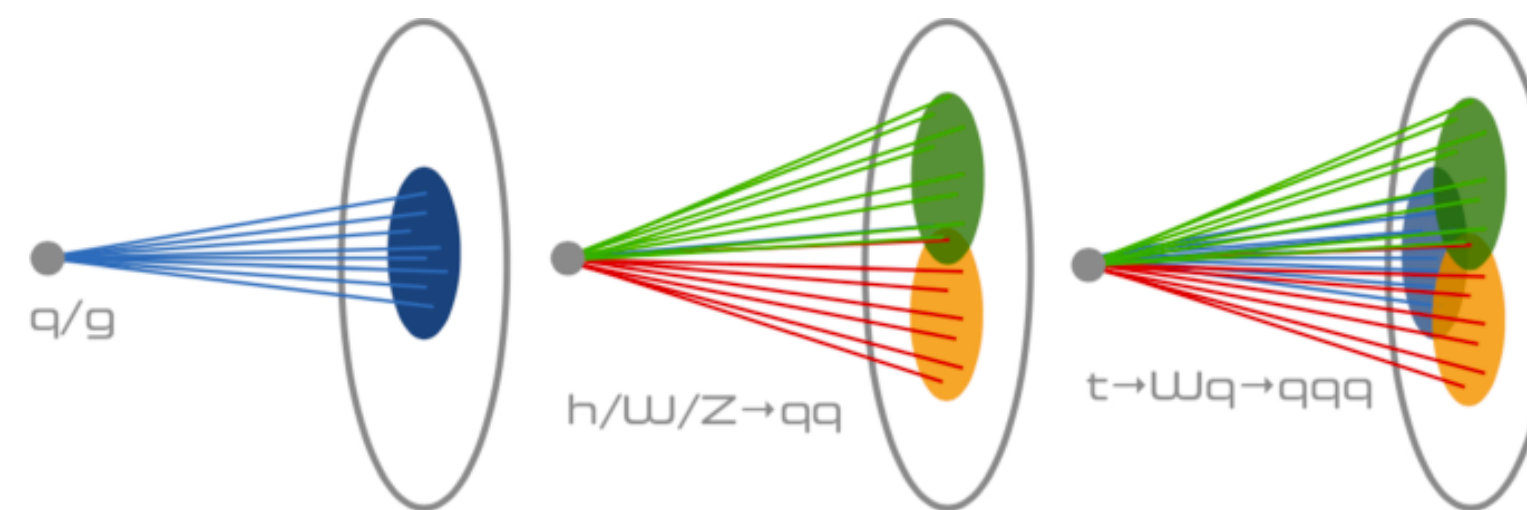
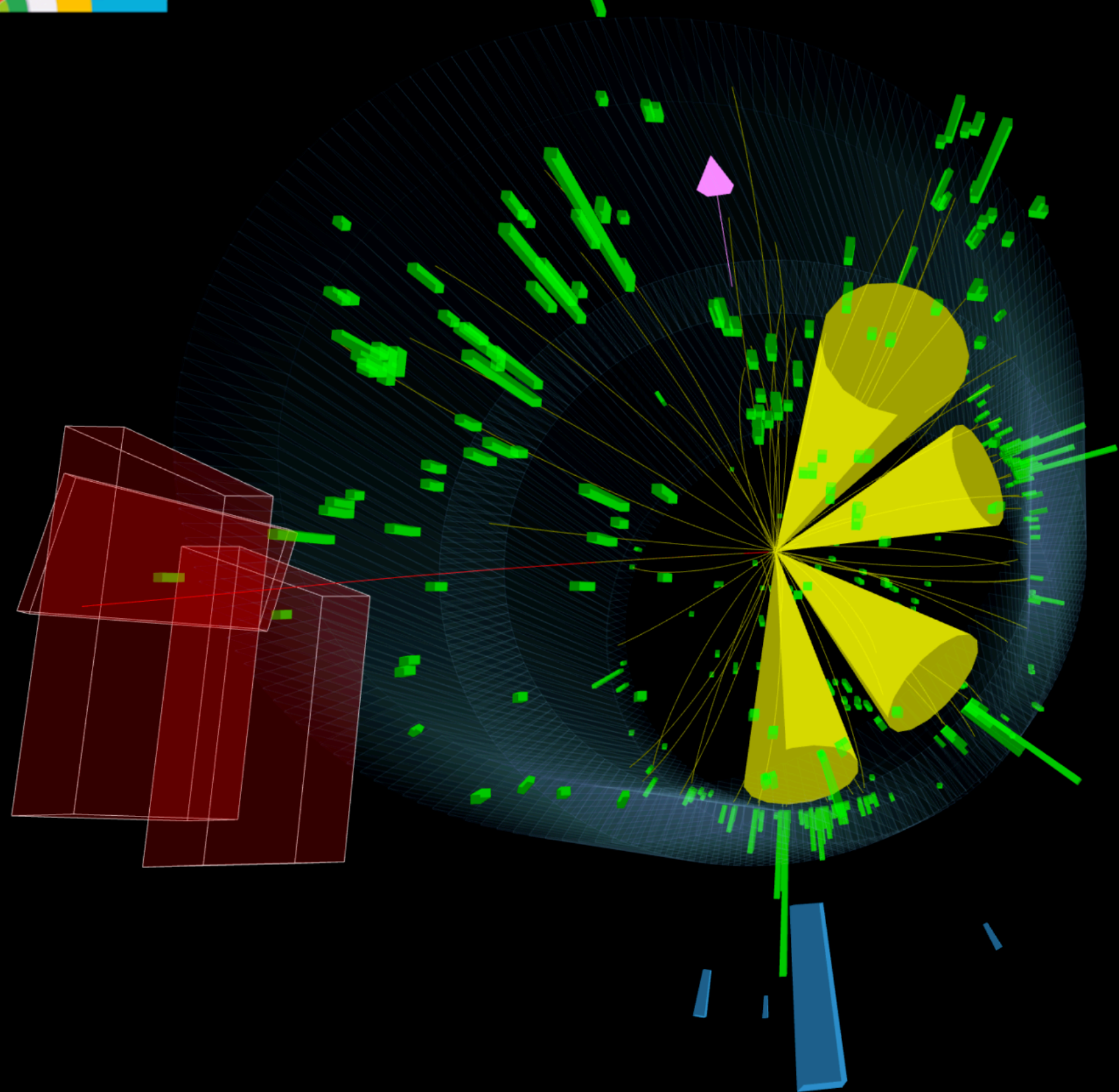
$$\mathcal{L}(y, y_T, y_S) = \alpha \mathcal{L}_{\text{cross-entropy}}(y, y_S) + (1 - \alpha) \mathcal{L}_{\text{distillation}}(y_T, y_S)$$

$$\mathcal{L}_{\text{distillation}}(y_T, y_S) = D_{\text{KL}}(\sigma(y_S/T) || \sigma(y_T/T)) \times T^2$$

*My Knowledge
Distillation Adventure*



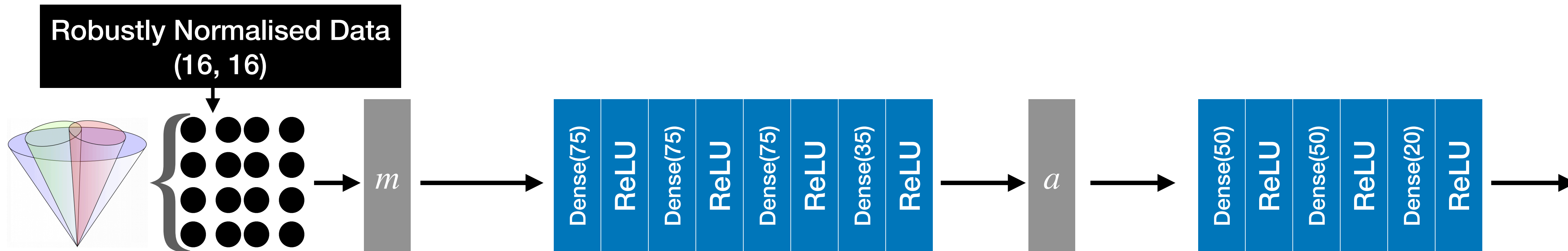
CMS Experiment at the LHC, CERN
Data recorded: 2016-Aug-17 08:01:23.065024 GMT
Run / Event / LS: 278969 / 229126383 / 184



quarks
gluons
Ws
Zs
tops

Training: 572 000 jets | 20% validation }
Testing: 282 000 jets } equally distributed between 5 classes:
quarks, gluons, Ws, Zs, tops

Every jet defined by maximum **150** constituents.
Every constituent has a maximum of **16** features.



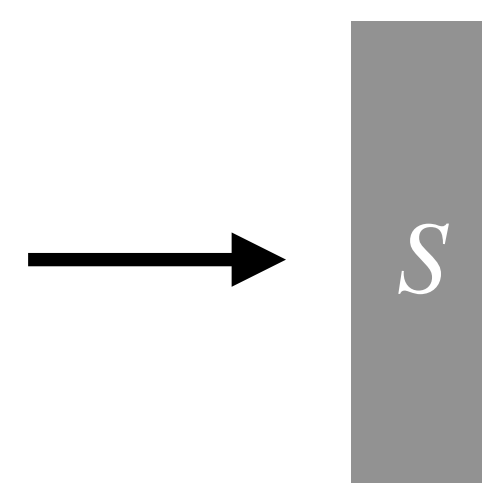
$$m(G) = B = \begin{pmatrix} I \times R_R \\ I \times R_S \end{pmatrix}$$

$$\dim(B) = (240, 32)$$

$$\phi_R(B) = E$$

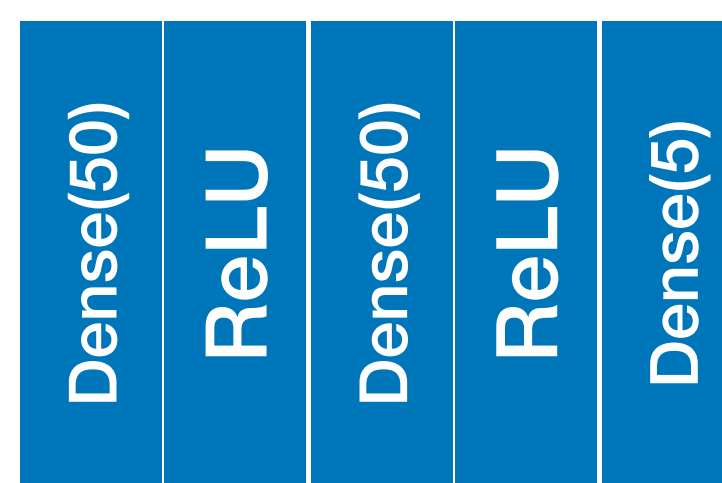
$$\dim(C) = (16, 51)$$

$$\phi_O(C) = P$$

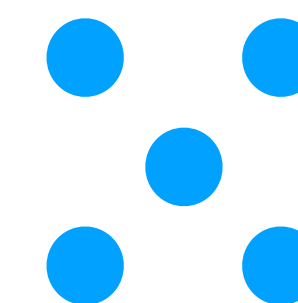


invariant

$$\dim(P) = 20$$

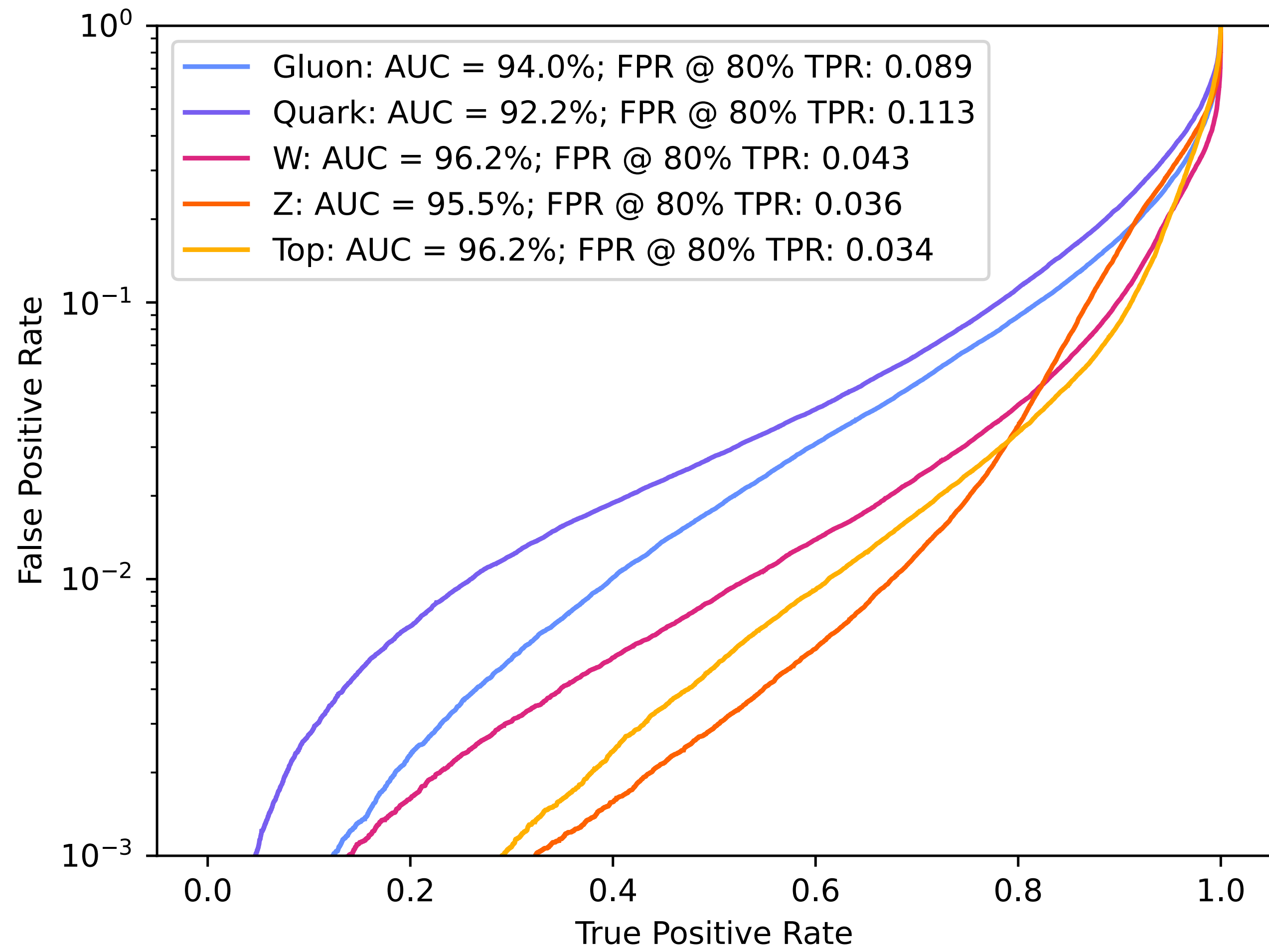


$$\phi_A(P) = A$$

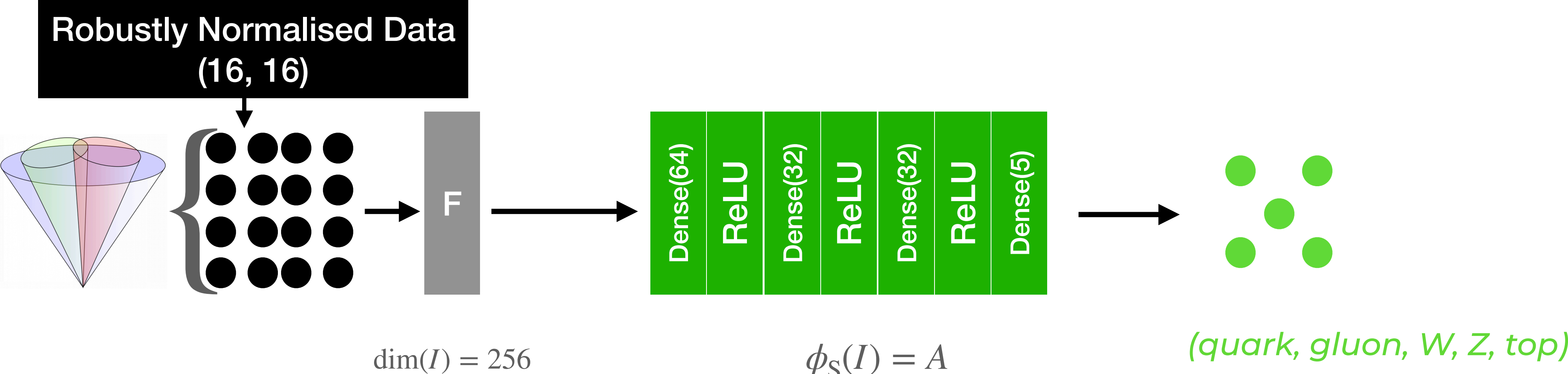


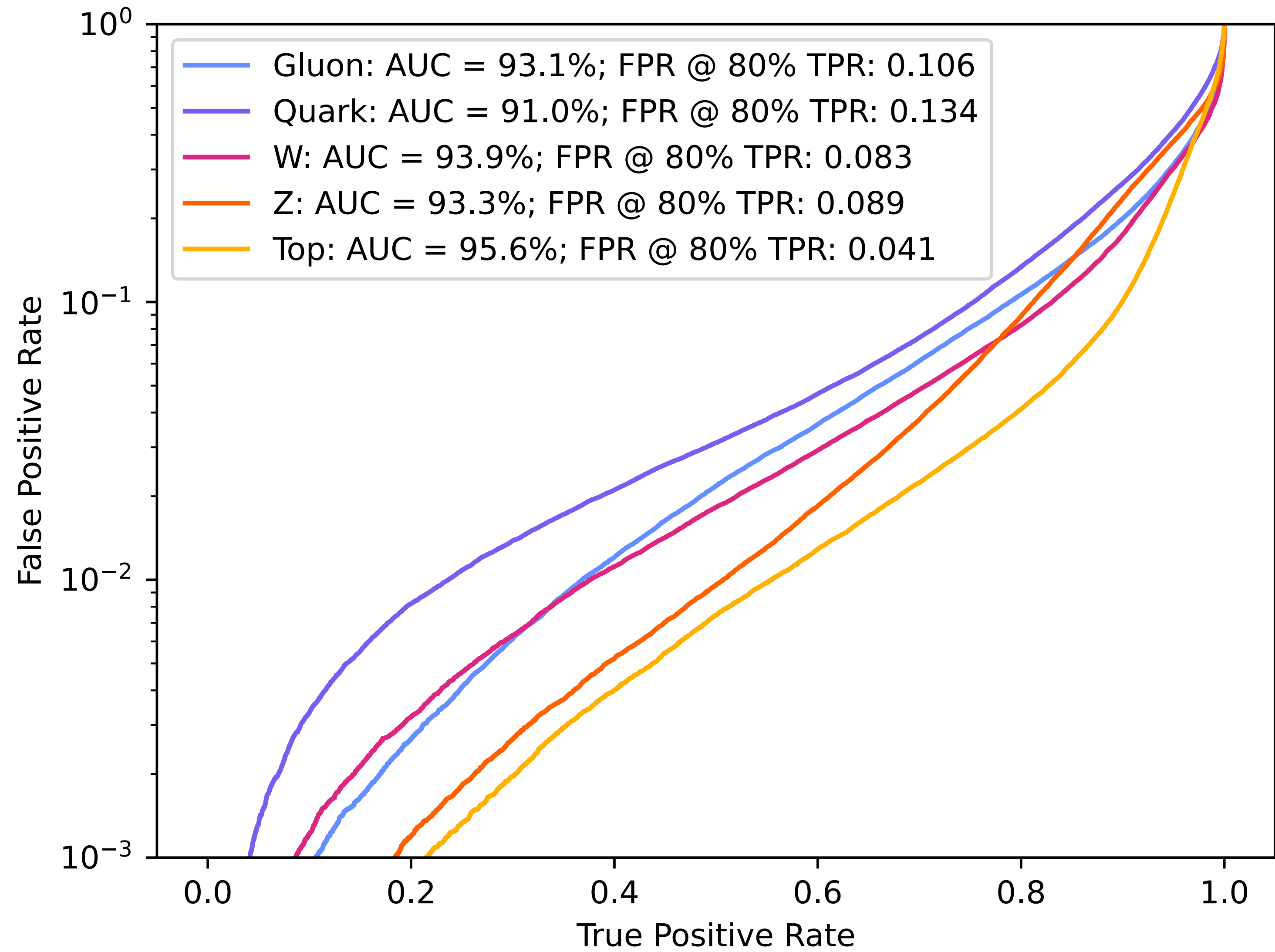
(quark, gluon, W, Z, top)





Appropriately Small Student





Robustly Normalised Data
(16, 16)



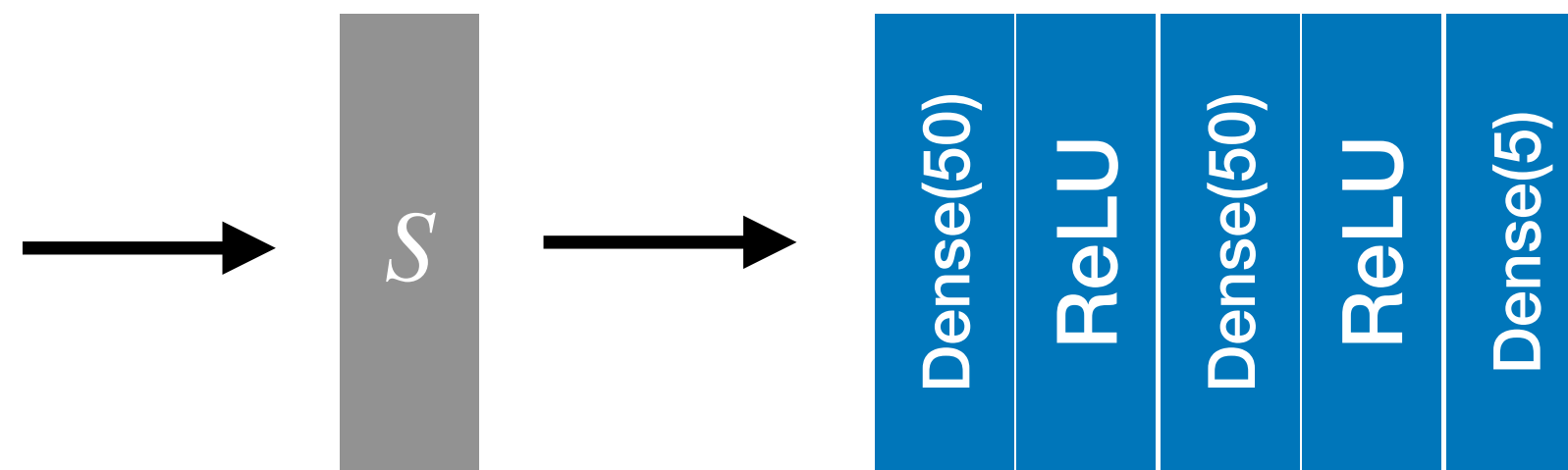
$$m(G) = B = \begin{pmatrix} I \times R_R \\ I \times R_S \end{pmatrix}$$

$$\dim(B) = (240, 32)$$

$$\phi_R(B) = E$$

$$\dim(C) = (16, 51)$$

$$\phi_O(C) = P$$



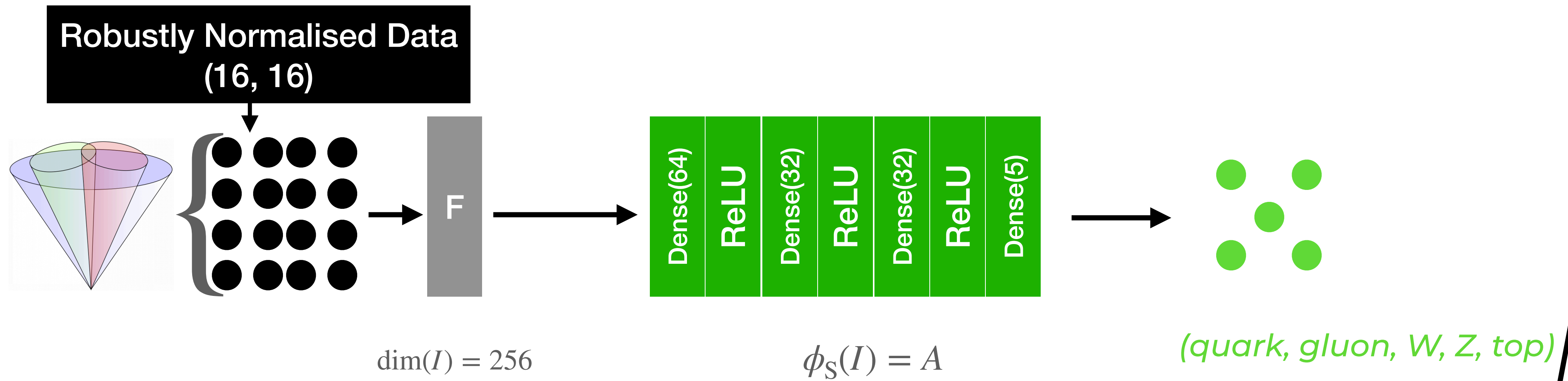
invariant

$$\dim(P) = 20$$

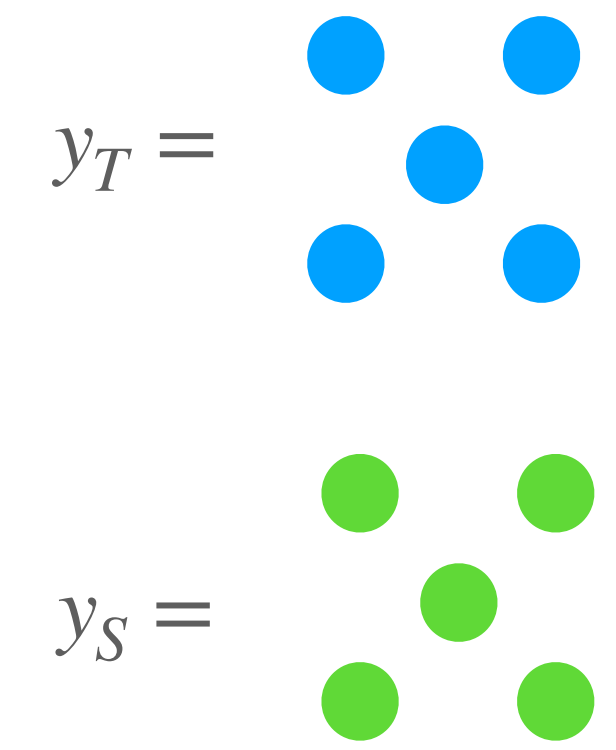
$$\phi_A(P) = A$$

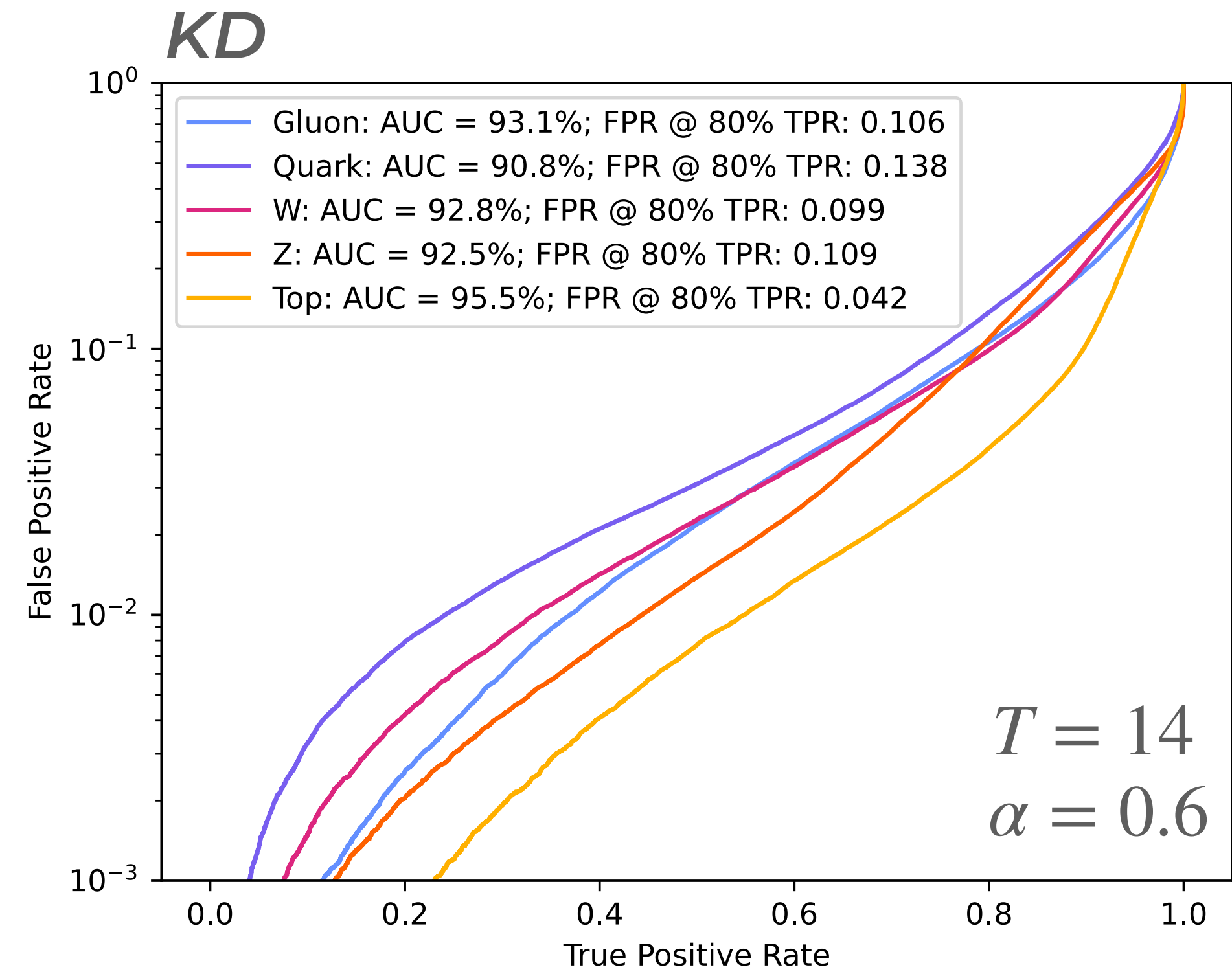
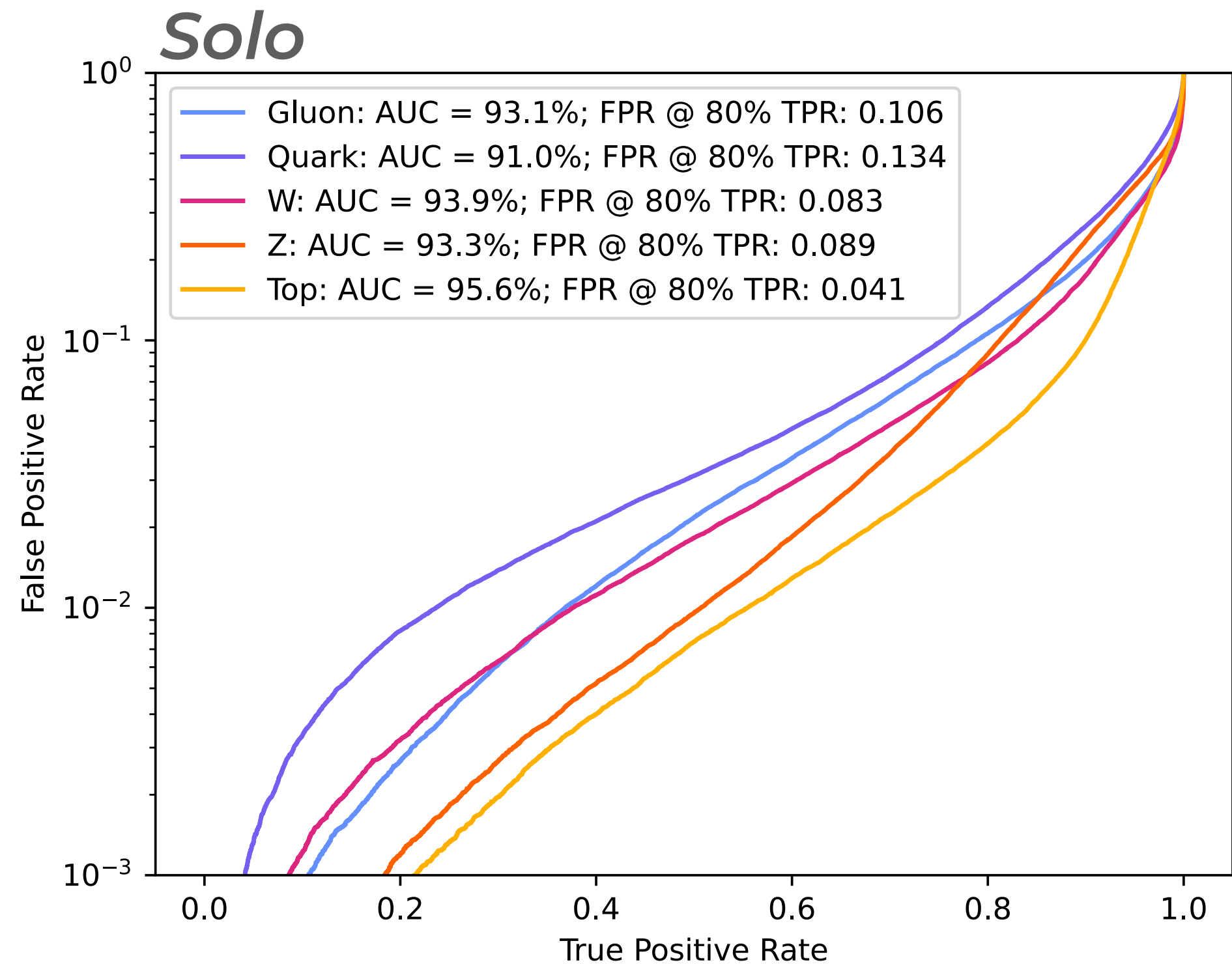


(quark, gluon, W, Z, top)

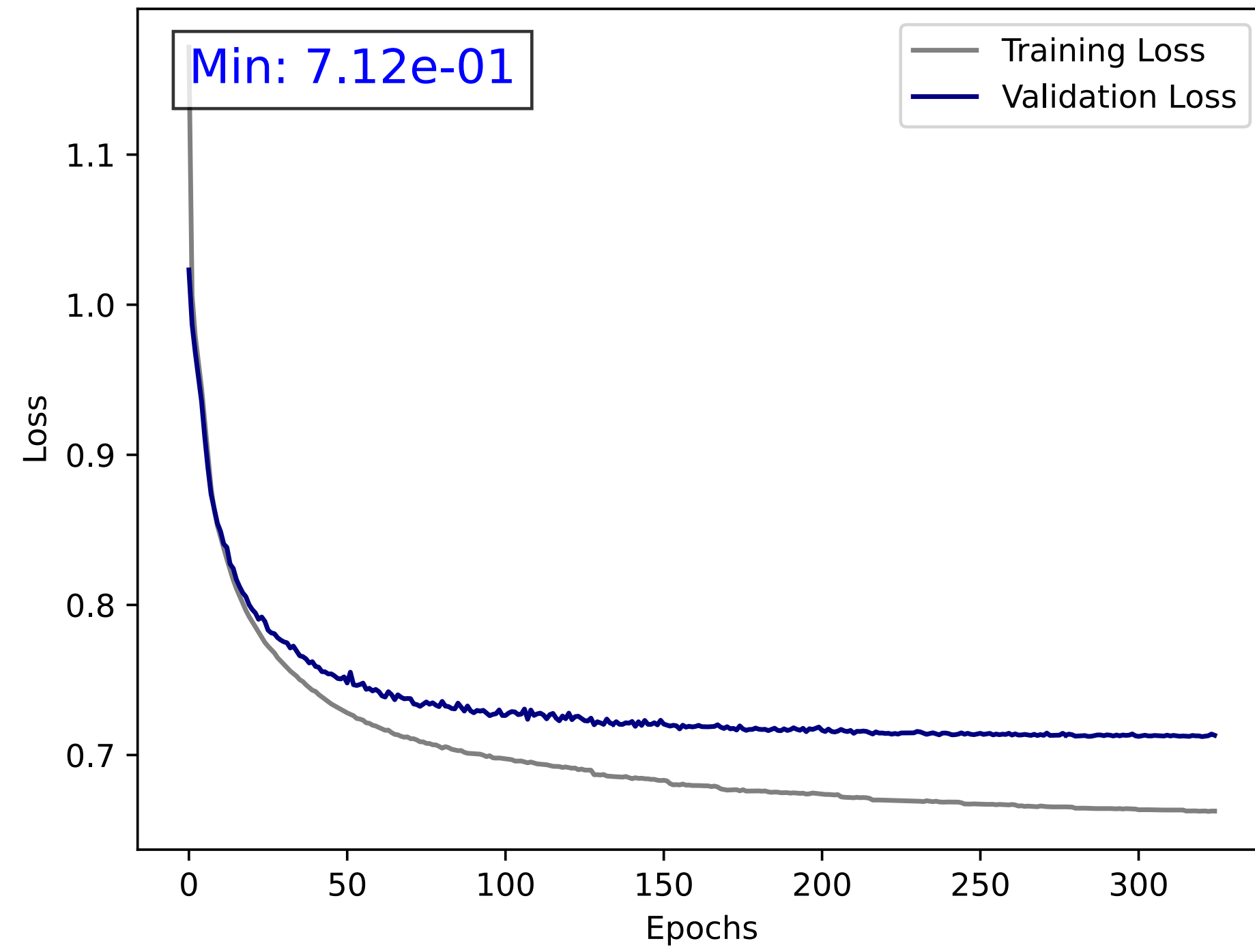


$$\mathcal{L}(y, y_T, y_S) = \alpha \mathcal{L}_{\text{cross-entropy}}(y, y_S) + (1 - \alpha) \mathcal{L}_{\text{distillation}}(y_T, y_S)$$

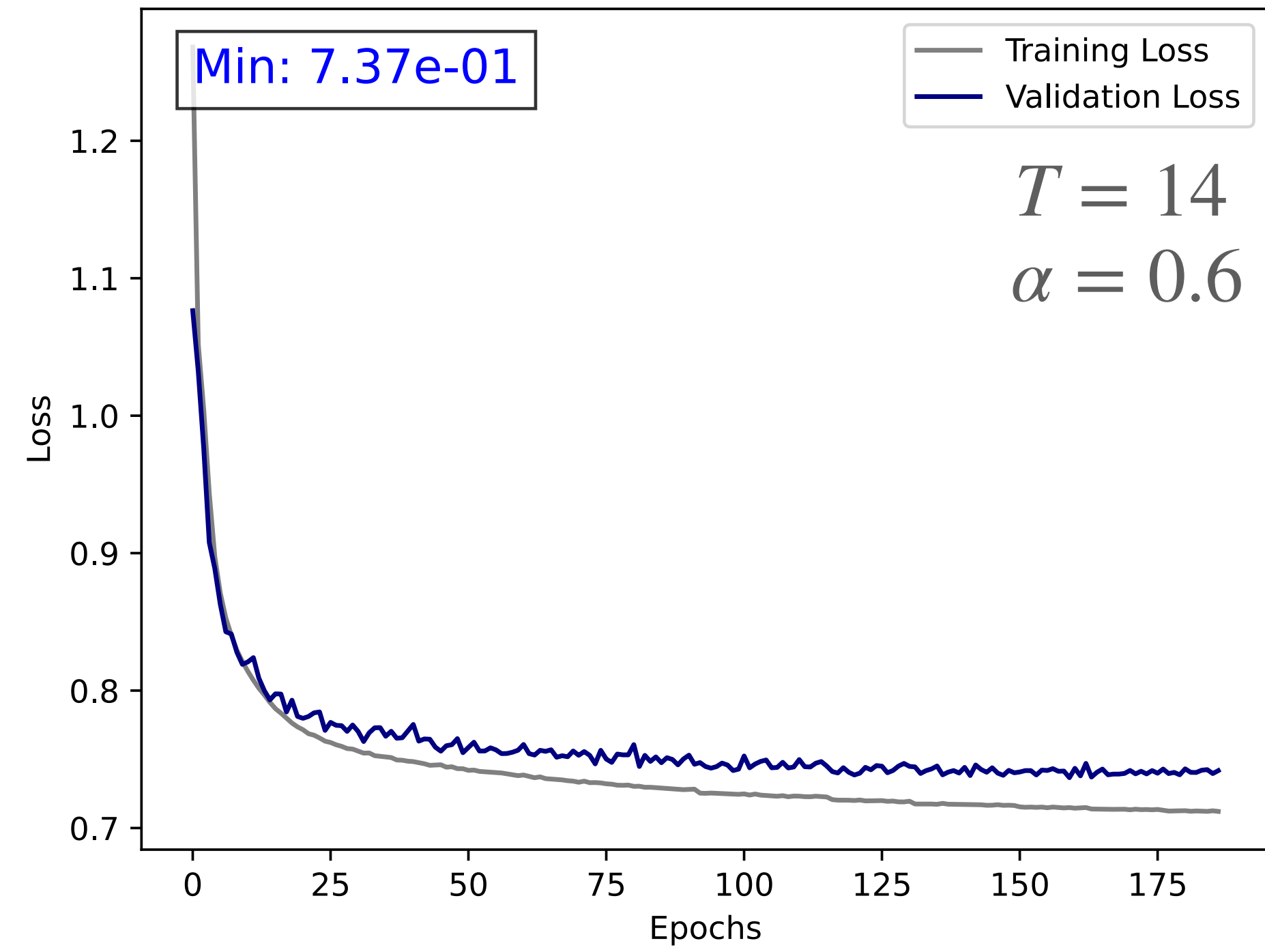




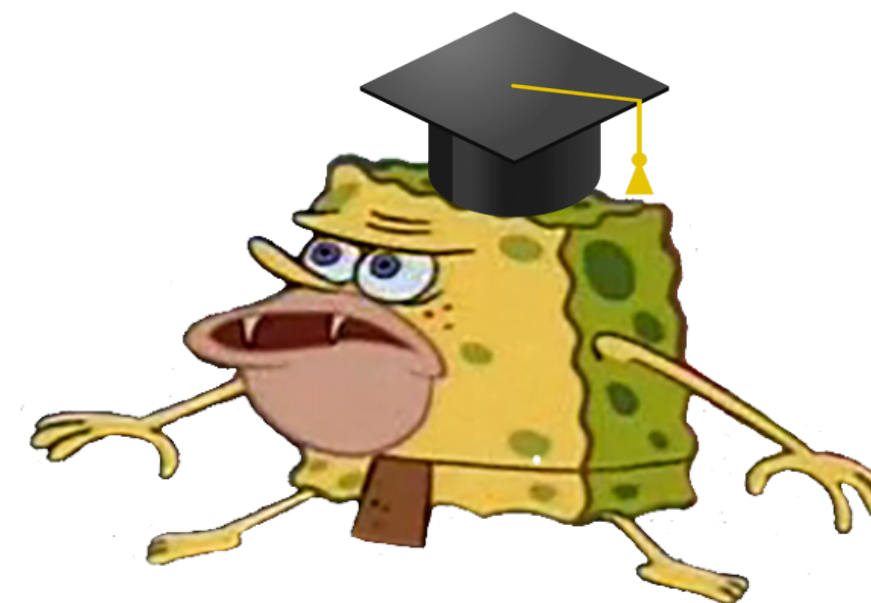
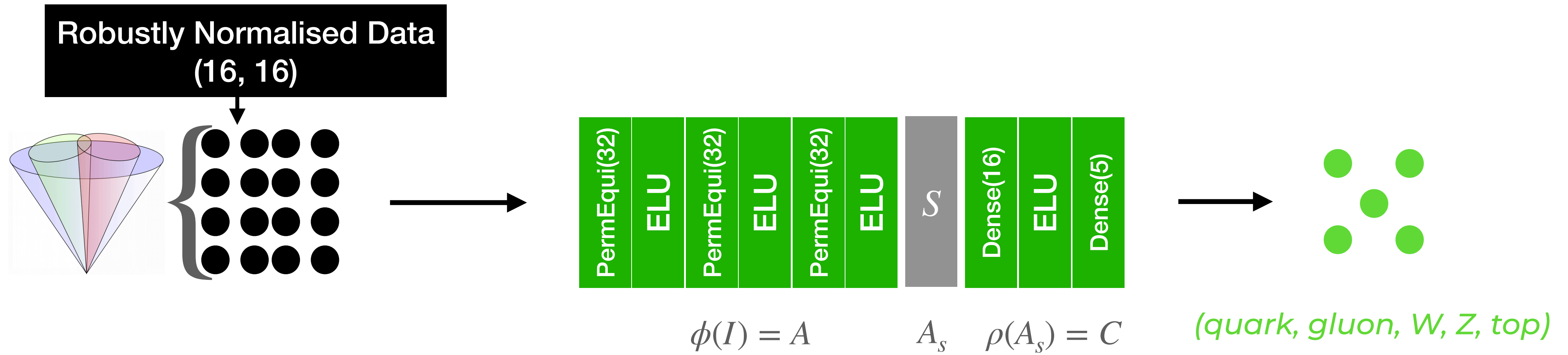
Solo

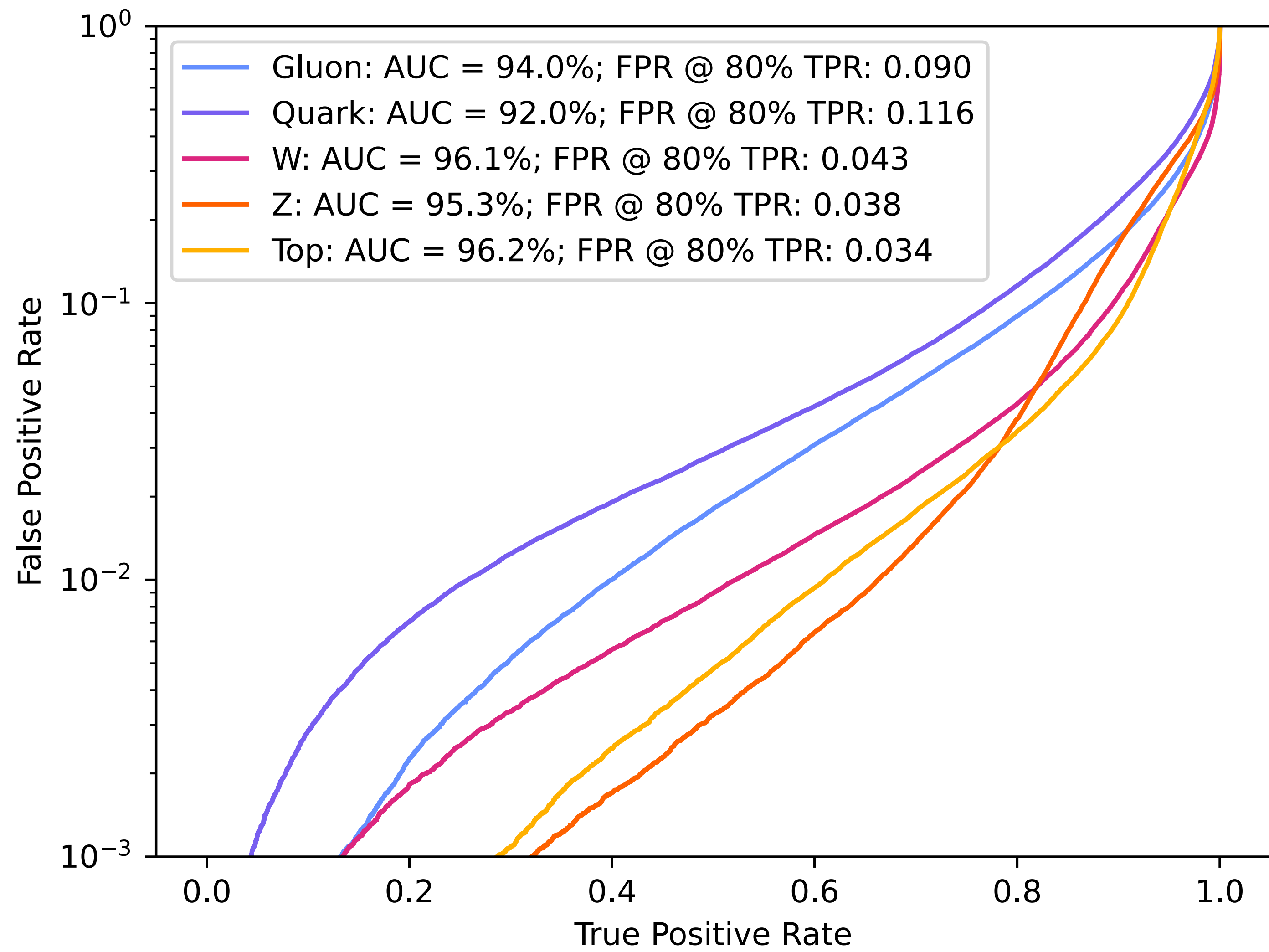


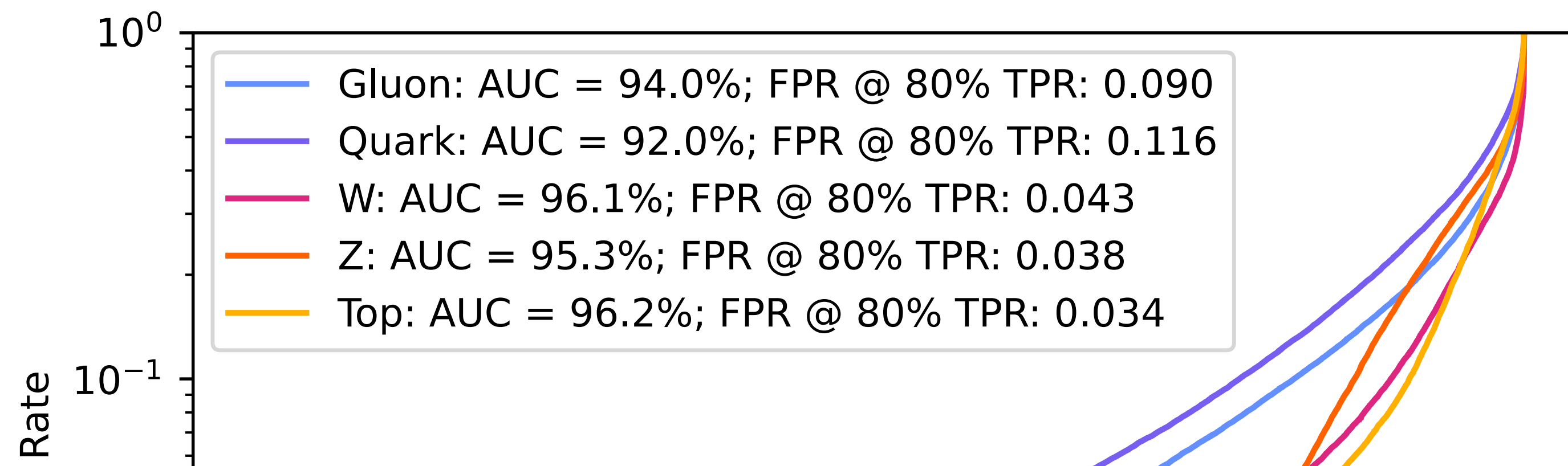
KD



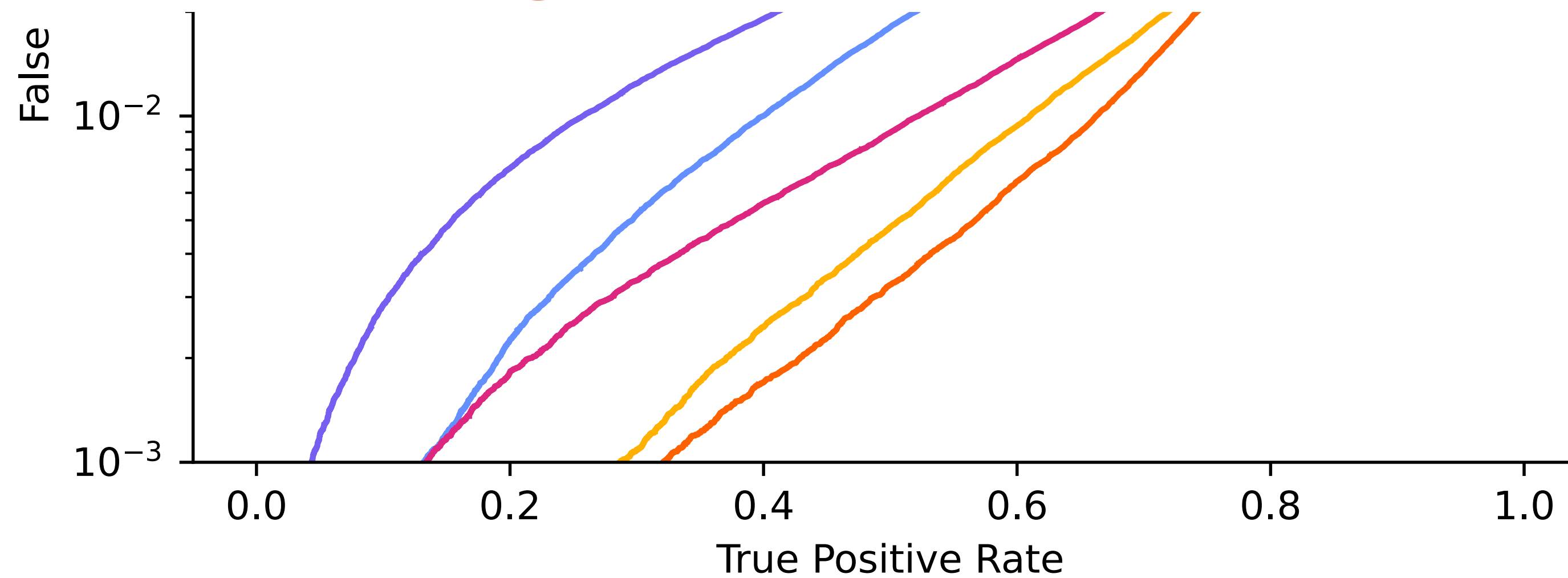
Deep Sets







Does Knowledge Distillation Really Work?



Does Knowledge Distillation Really Work?

Samuel Stanton
NYU

Pavel Izmailov
NYU

Polina Kirichenko
NYU

Alexander A. Alemi
Google Research

Andrew Gordon Wilson
NYU

Does Knowledge Distillation Really Work?

Samuel Stanton
NYU

Pavel Izmailov
NYU

Polina Kirichenko
NYU

Alexander A. Alemi
Google Research

Andrew Gordon Wilson
NYU

Does knowledge distillation really work? In short:

Does Knowledge Distillation Really Work?

Samuel Stanton
NYU

Pavel Izmailov
NYU

Polina Kirichenko
NYU

Alexander A. Alemi
Google Research

Andrew Gordon Wilson
NYU

Does knowledge distillation really work? In short: **Yes**

Does Knowledge Distillation Really Work?

Samuel Stanton
NYU

Pavel Izmailov
NYU

Polina Kirichenko
NYU

Alexander A. Alemi
Google Research

Andrew Gordon Wilson
NYU

Does knowledge distillation really work? In short: *Yes*, in the sense that it often improves student generalization. *No*, in that knowledge distillation often fails to live up to its name, transferring very limited knowledge from teacher to student.

Moral of the Story

- KD is a good paradigm for model compression
- Ensemble knowledge distillation: **safe**
- Distilling across architectures: **caution**
- Distilling a big model to a smaller version of itself: **depends on dataset**
- You can use KD for regularisation.

Model Compression - Bucila et. al.

Distilling Knowledge in a Neural Network - Hinton et. al.

Knowledge Distillation: A Survey - Gou et. al.

Does Knowledge Distillation Really Work? - Stanton et. al.

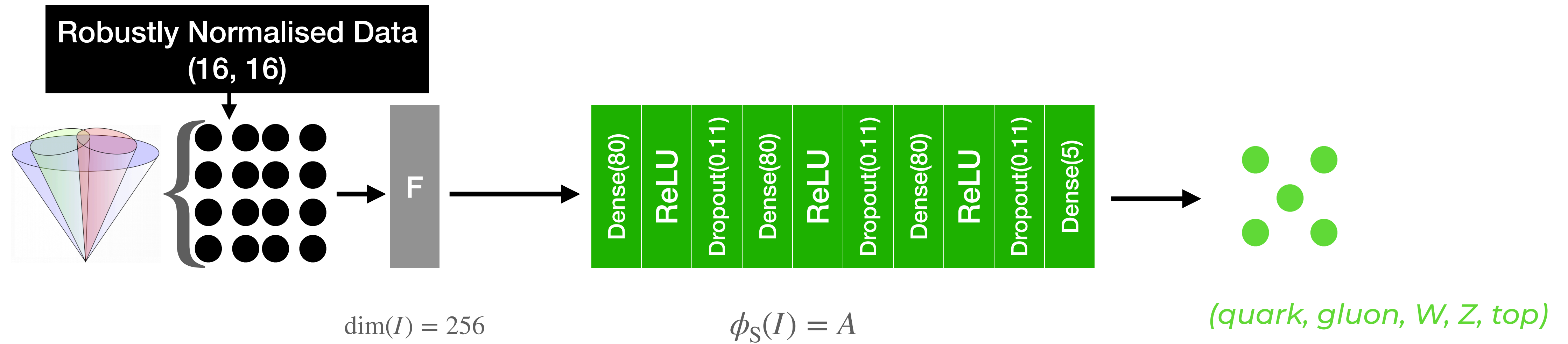
FitNets: Hints for Thin Deep Jets - Stanton et. al.

https://keras.io/examples/vision/knowledge_distillation/

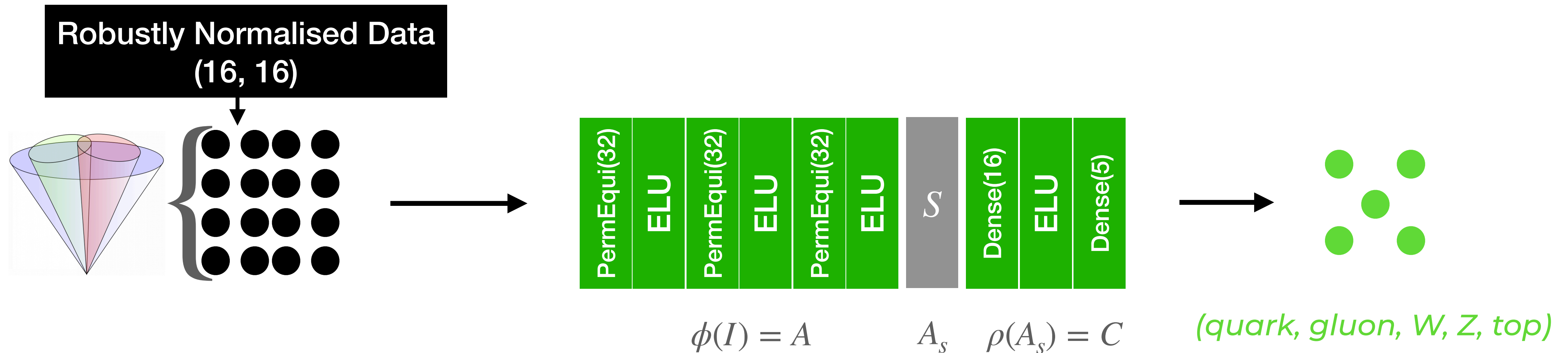
backup

JEDInet DNN

arxiv:1908.05318v3



Deep Sets



Permutation Equivariant Layer + Activation

$$f(\mathbf{x}) = \sigma(\mathbf{x}\Lambda - \mathbf{1}\mathbf{1}^\top \mathbf{x}\Gamma)$$

$$f(\mathbf{x}) = \sigma(\mathbf{x}\Lambda - \mathbf{1}\text{maxpool}(\mathbf{x})\Gamma)$$

Non-linearity applied to weighted combination of

➔ Input.

➔ Sum of output values.

$$H_q = - \mathbb{E}_{x \sim p(\cdot)} [\log q(x)]$$

$$\frac{\partial H}{\partial g(z)_i} \approx \frac{1}{T} \left(\frac{1 + g(z)_i/T}{N + \sum_j g(z)_j/T} - \frac{1 + f(z)_i/T}{N + \sum_j f(z)_j/T} \right) \quad T \gg 1$$

$$\frac{\partial H}{\partial g(z)} \approx \frac{1}{NT^2} (g(z) - f(z))$$

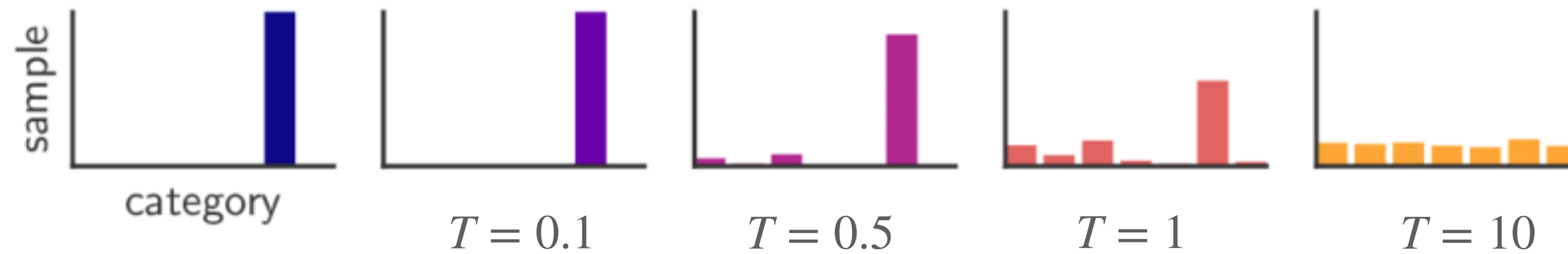
model compression!

$$H_q = - \mathbb{E}_{z \sim p(\cdot)} [\log q(z)]$$

$$\begin{aligned} \frac{\partial C}{\partial z_i} &= - \sum_j p_j \frac{\partial \log(q_j)}{\partial z_i} = - \sum_j \frac{p_j}{q_j} \frac{\partial q_j}{\partial z_i} = - \frac{p_i}{q_i} \frac{\partial q_i}{\partial z_i} - \sum_{j \neq i} \frac{p_j}{q_j} \frac{\partial q_j}{\partial z_i} \\ &= - \frac{1}{T} \frac{p_i}{q_i} q_i (1 - q_i) - \sum_{j \neq i} \frac{1}{T} \frac{p_j}{q_j} (-q_j q_i) \\ &= - \frac{p_i}{T} + \frac{1}{T} p_i q_i + \sum_{j \neq i} p_j q_i \\ &= - \frac{p_i}{T} + \frac{q_i}{T} \sum_j p_j = \frac{q_i - p_i}{T} \\ &= \frac{1}{T} \left(\frac{e^{z_i/T}}{\sum_j e^{z_j/T}} - \frac{e^{v_i/T}}{\sum_j e^{v_j/T}} \right) \end{aligned}$$

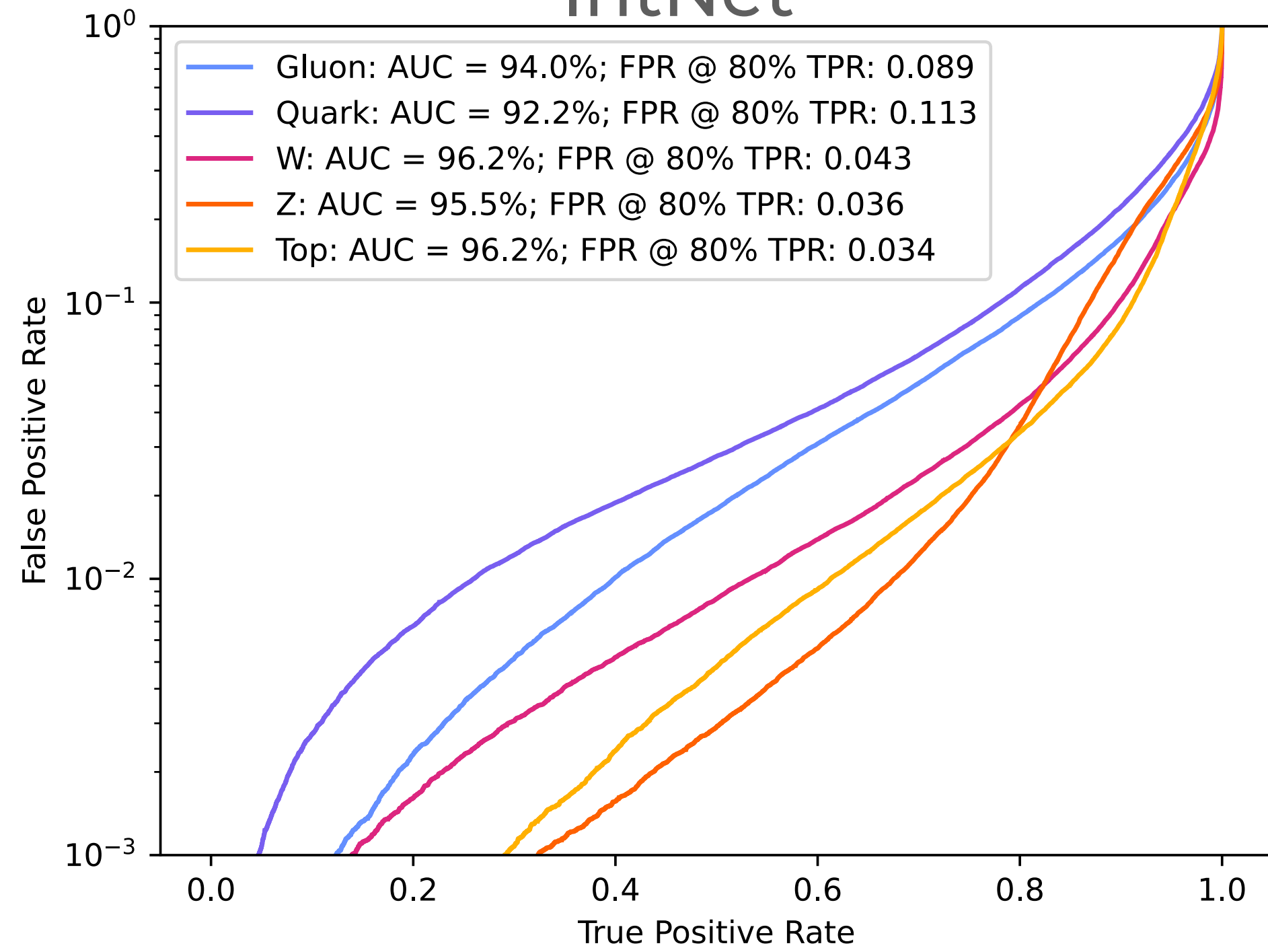
$$\frac{\partial C}{\partial z_i} \approx \frac{1}{T} \left(\frac{1 + z_i/T}{N + \sum_j z_j/T} - \frac{1 + v_i/T}{N + \sum_j v_j/T} \right)$$

$$\frac{\partial C}{\partial z_i} \approx \frac{1}{NT^2} (z_i - v_i)$$

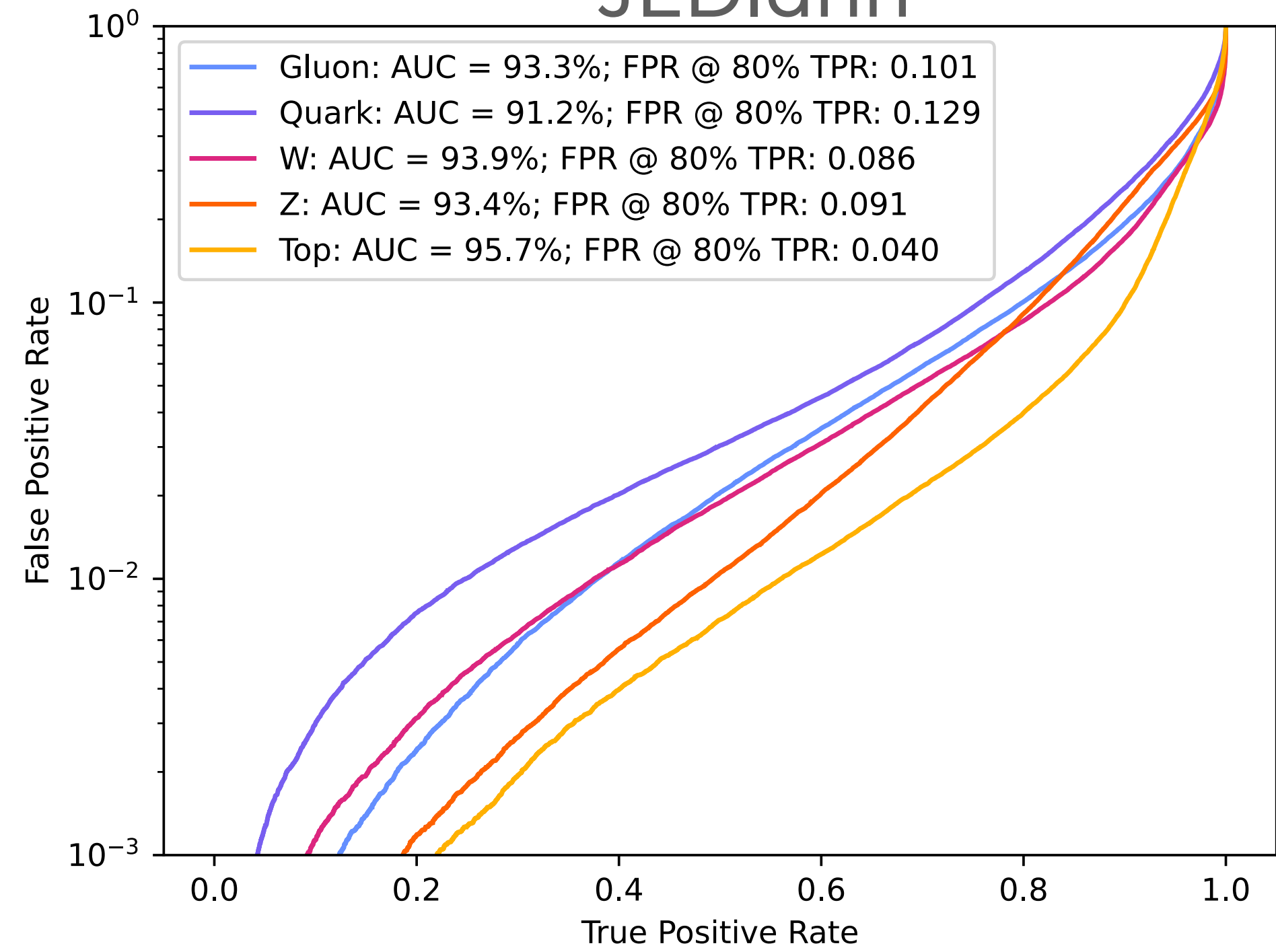


- ➔ Low temperature leads to harder labels.
- ➔ High temperature leads to softer labels.
- ➔ For some (maybe all? - proven only for categorical cross-entropy) loss functions, very high temperatures lead to equivalency between knowledge distillation and logit matching (MSE between teacher and student logits).

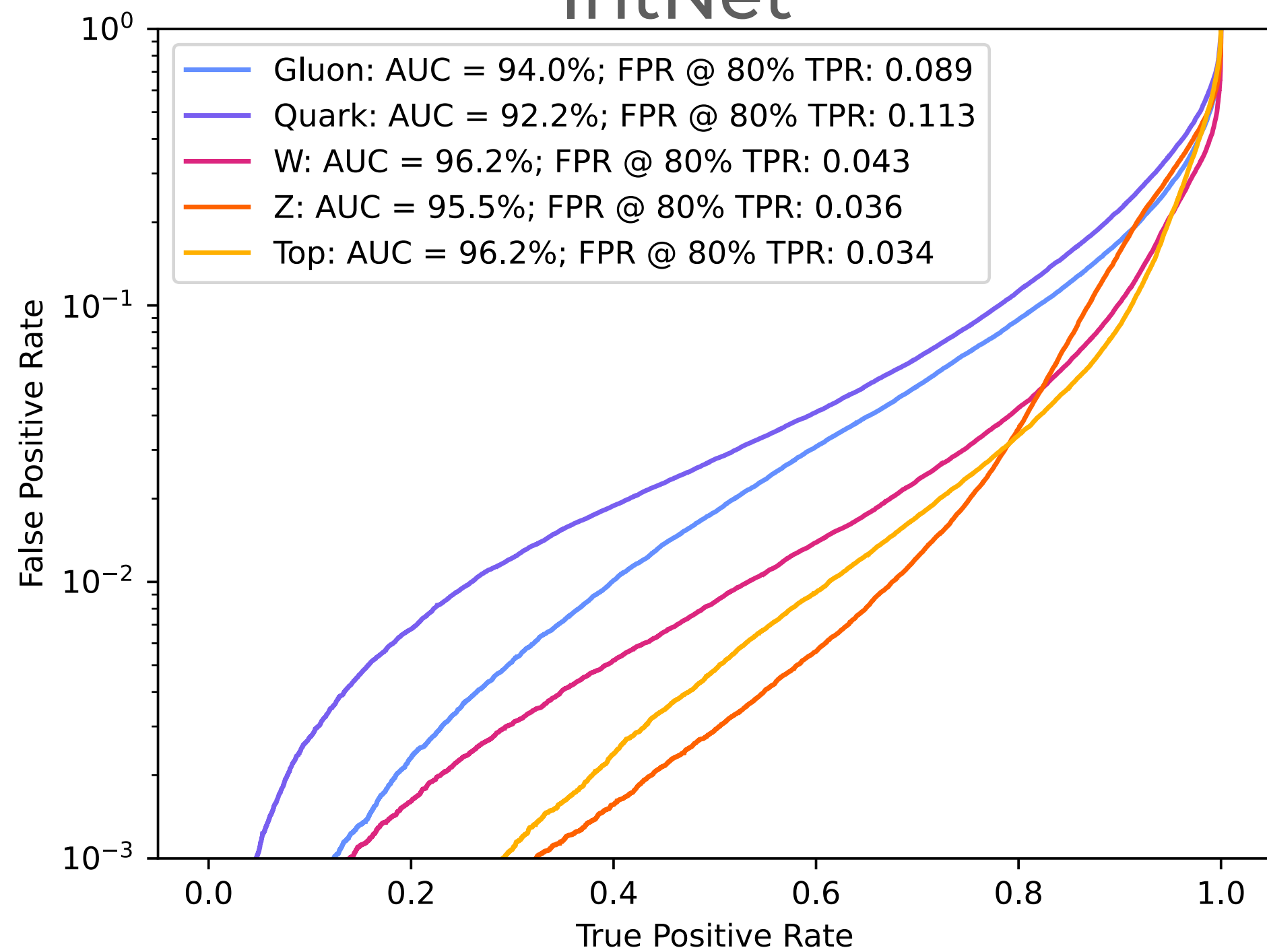
IntNet



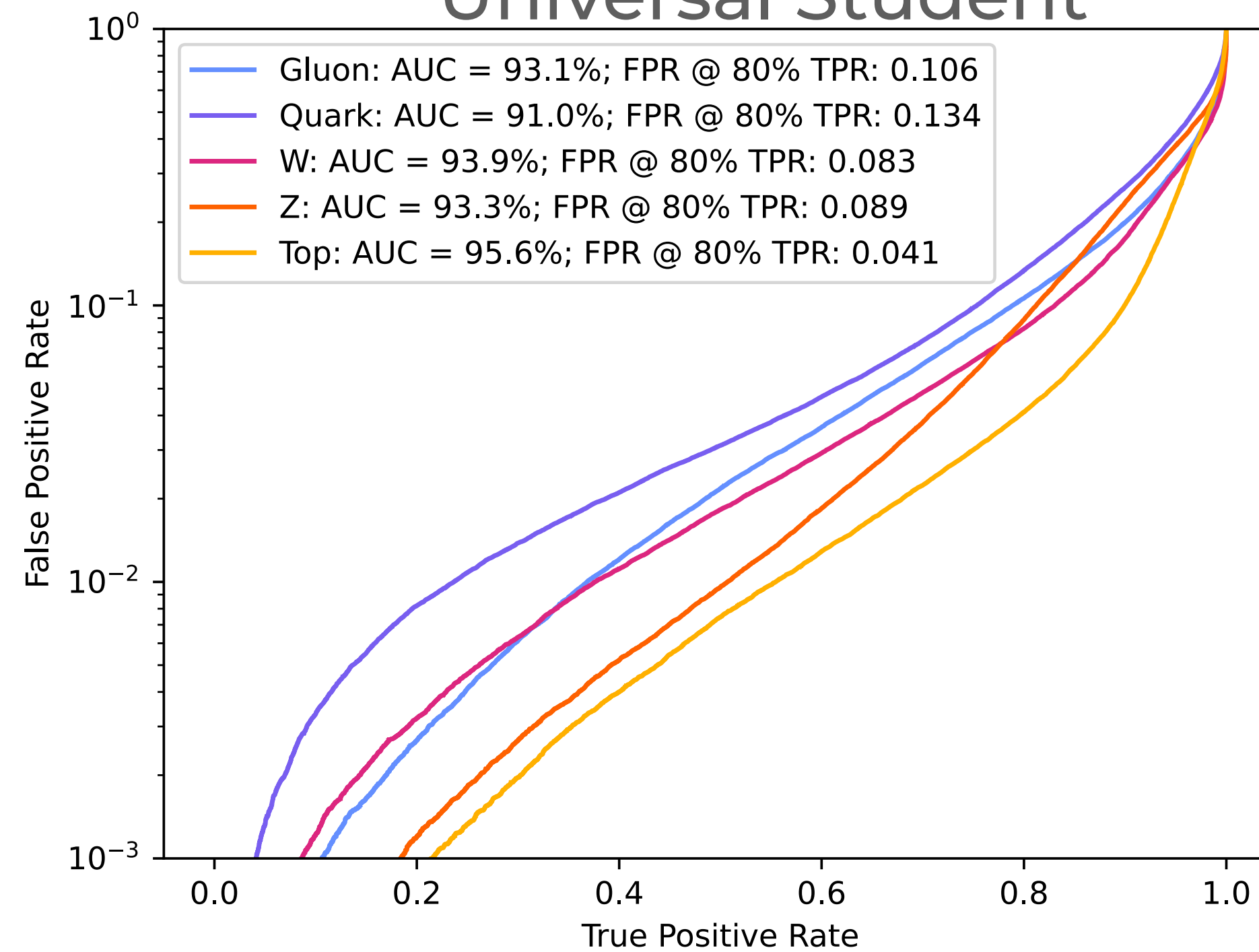
JEDIdnn



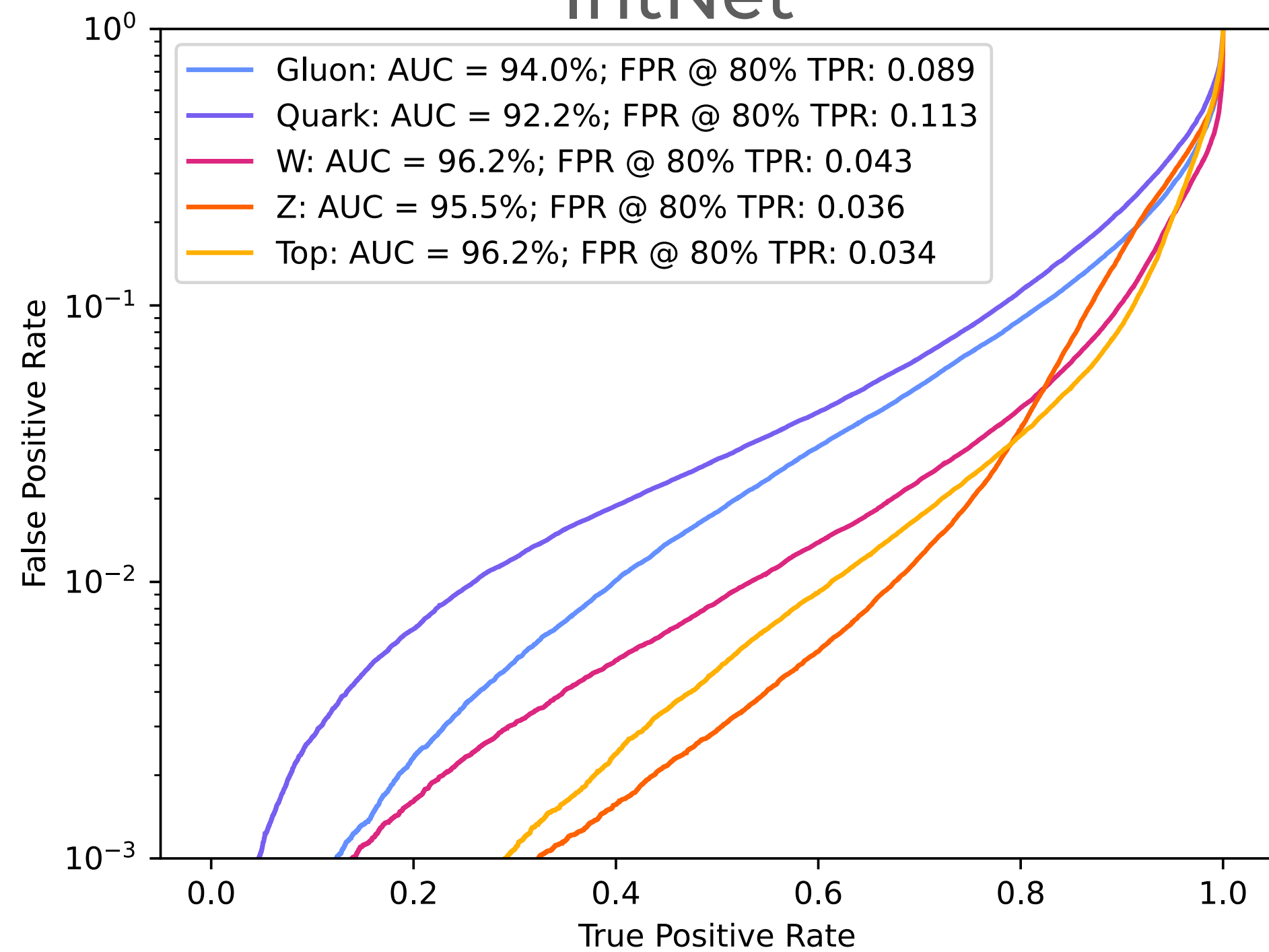
IntNet



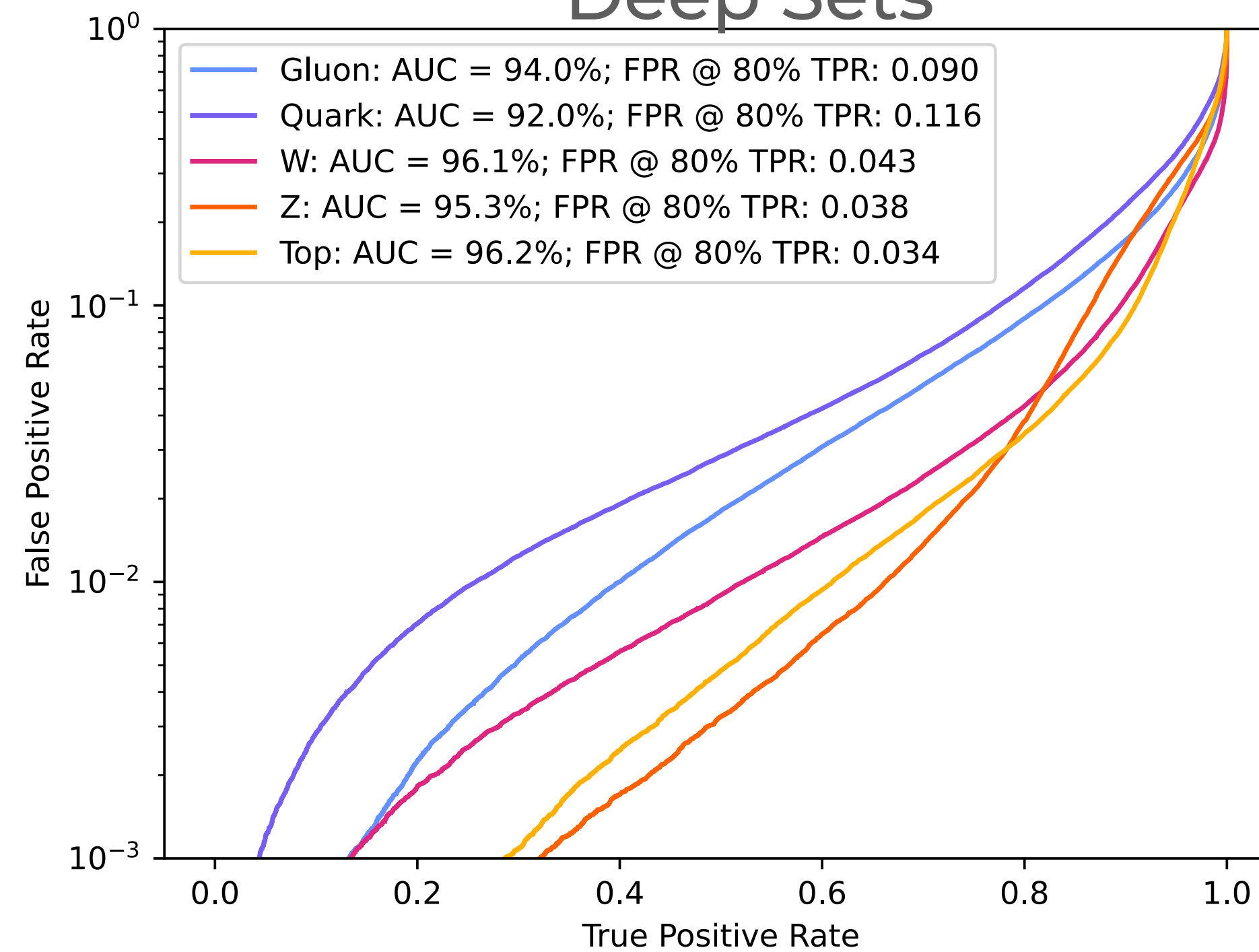
Universal Student



IntNet



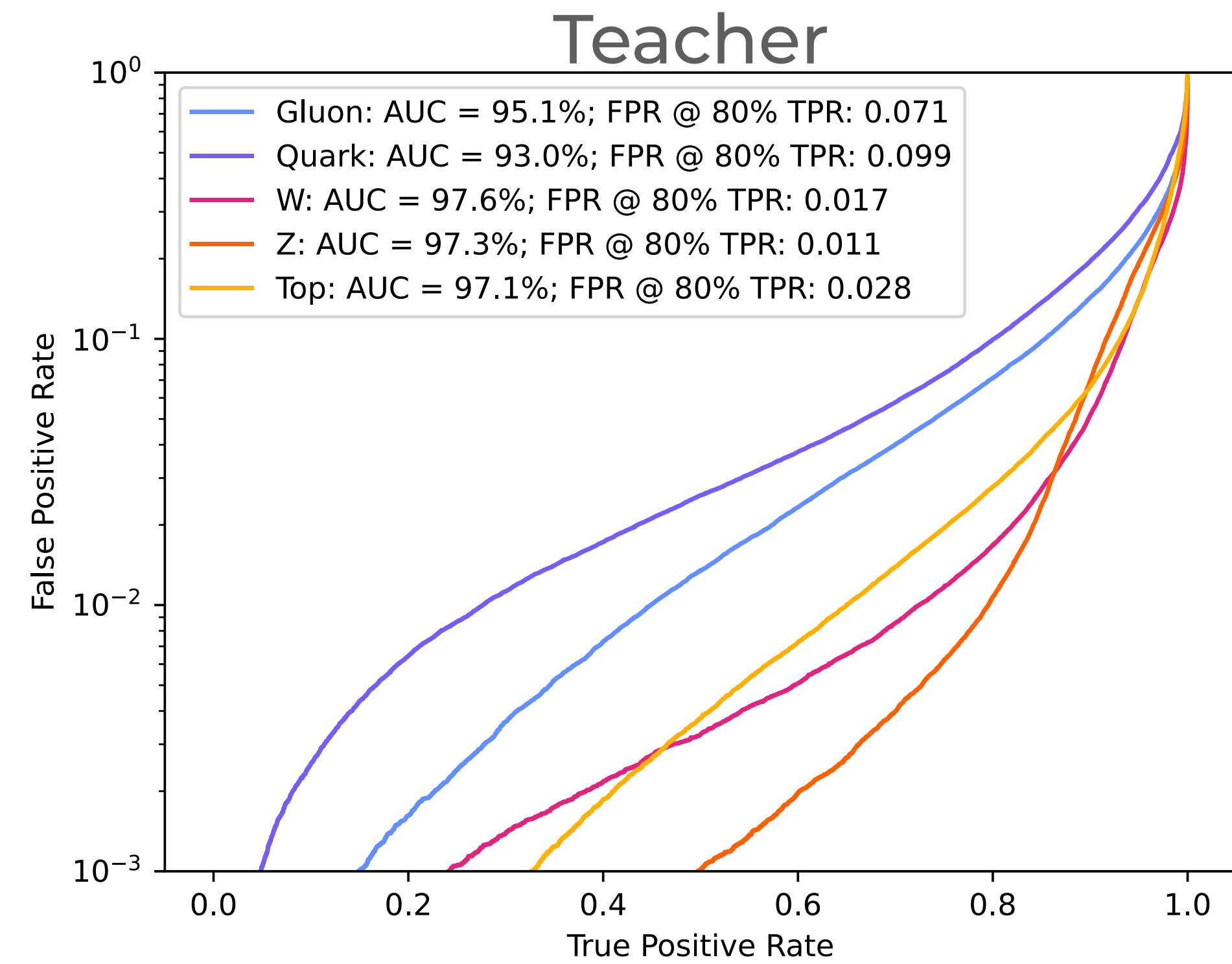
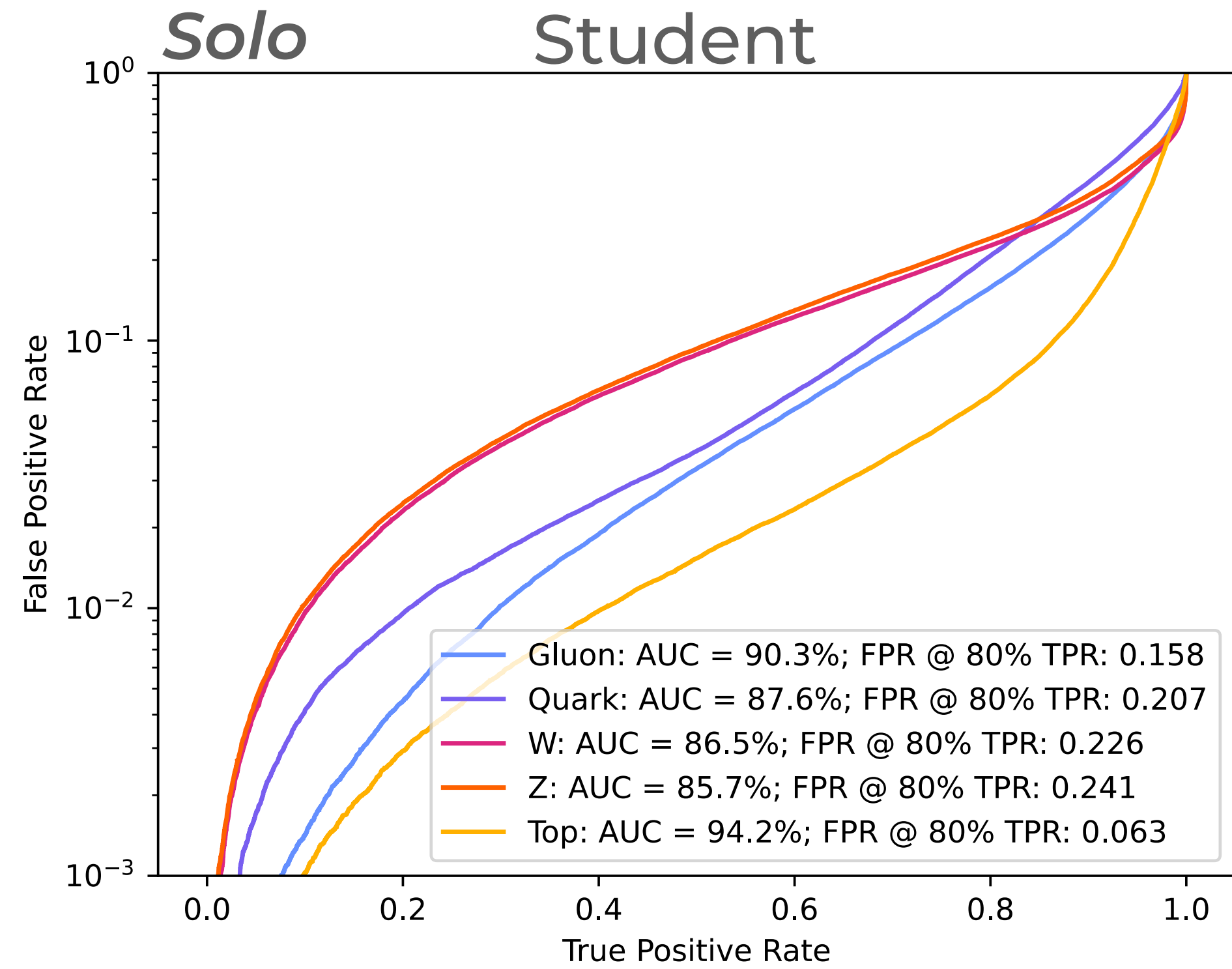
Deep Sets

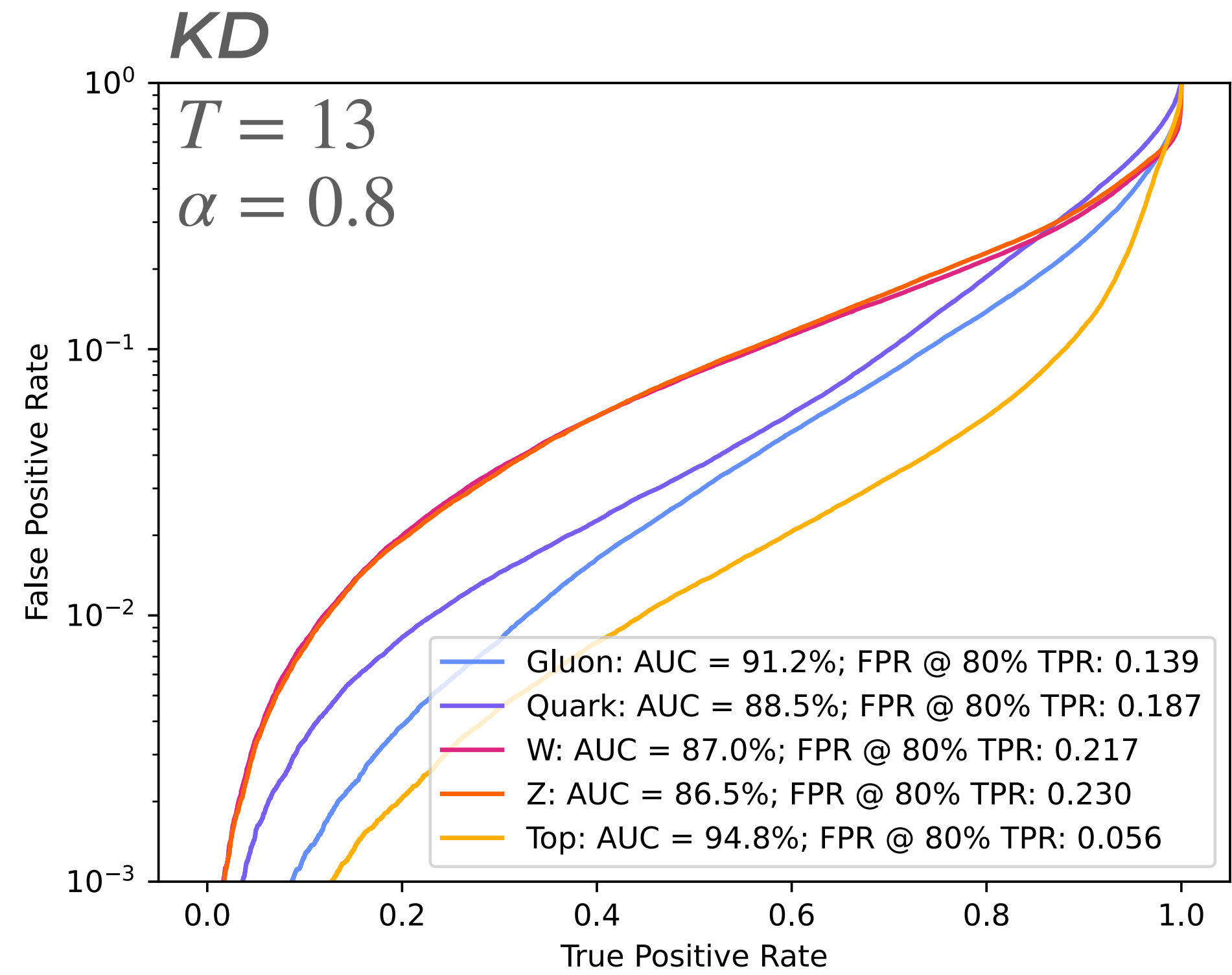
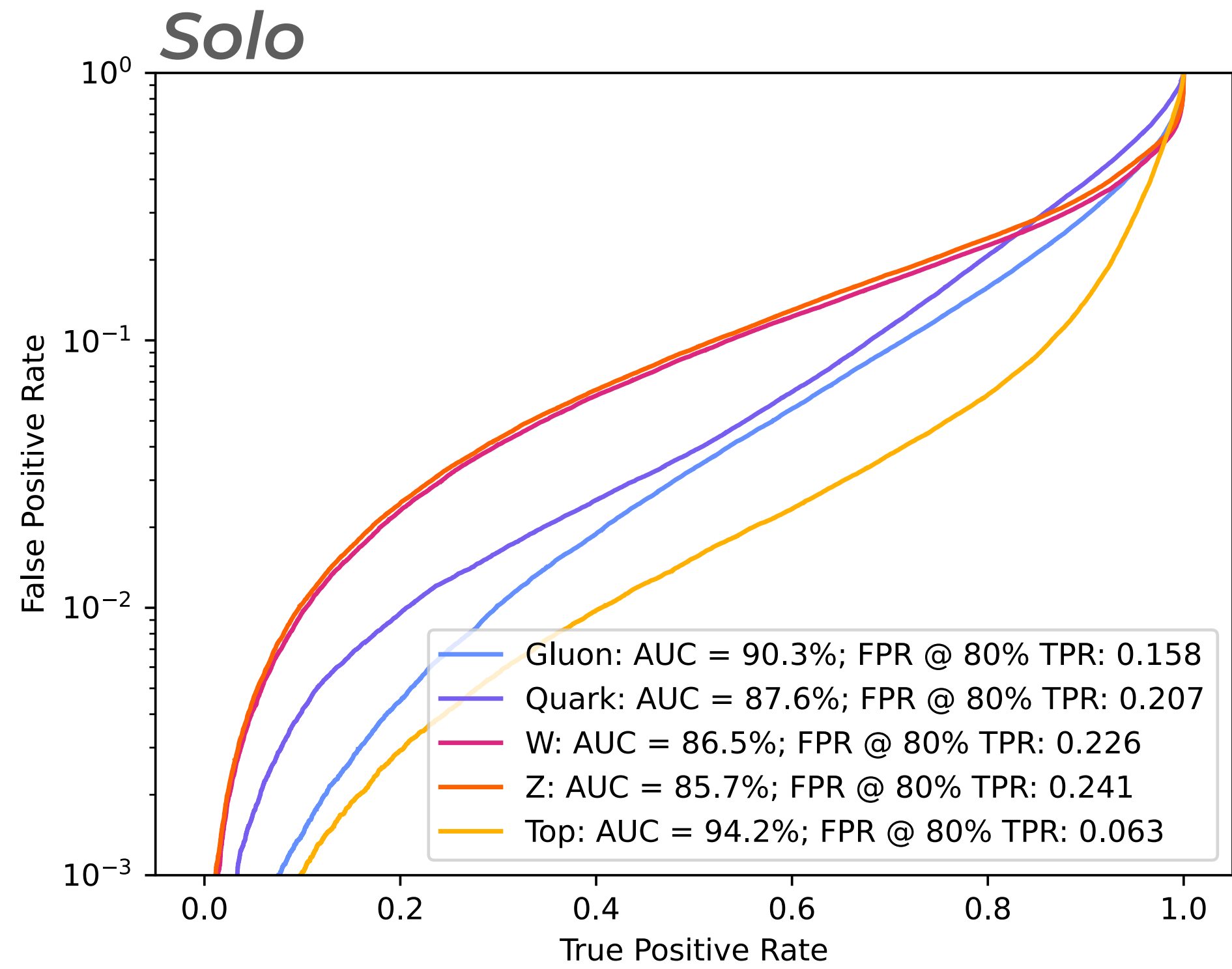


JEDInet DNN

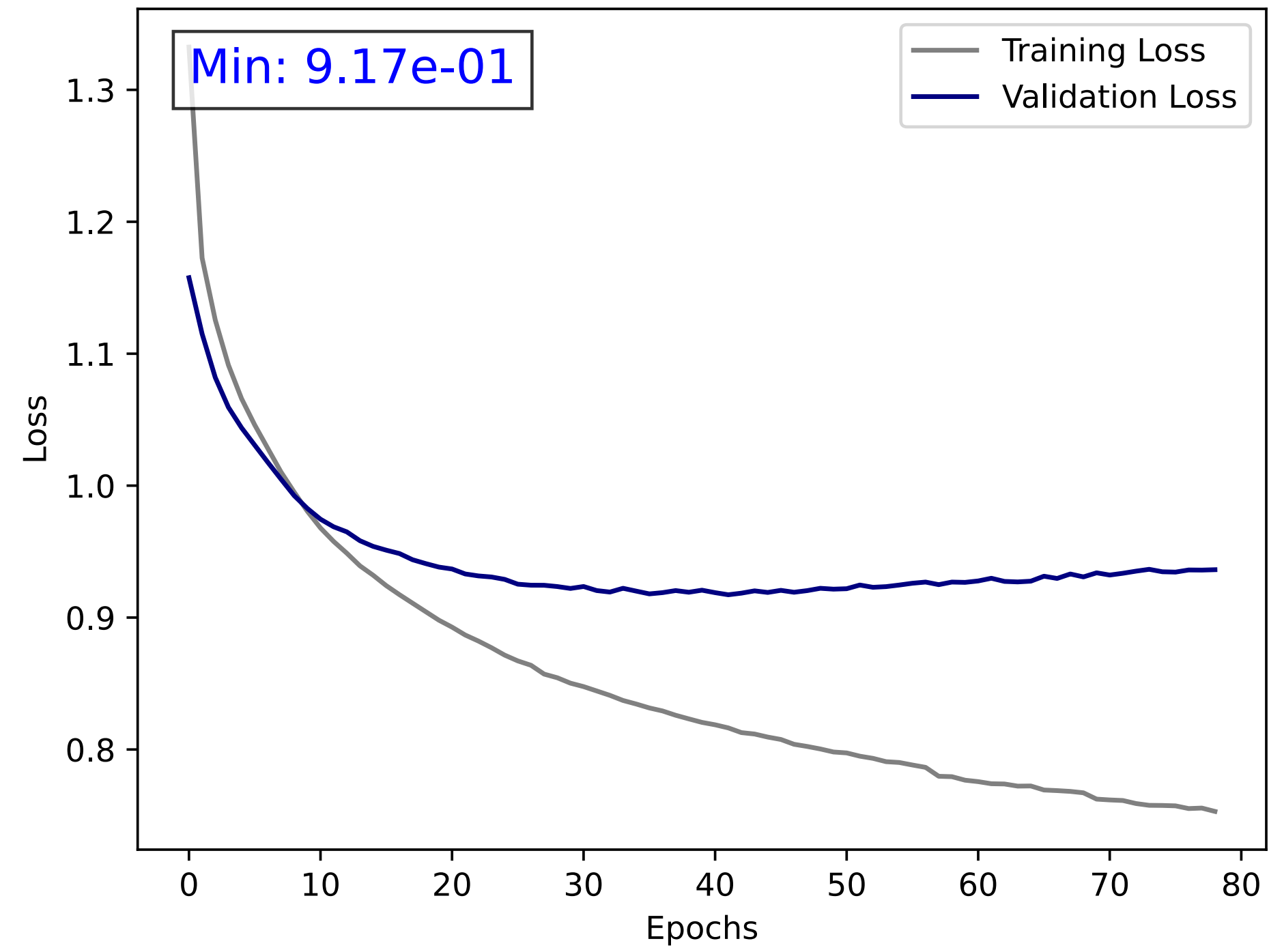
arxiv:1908.05318v3

150 constituents

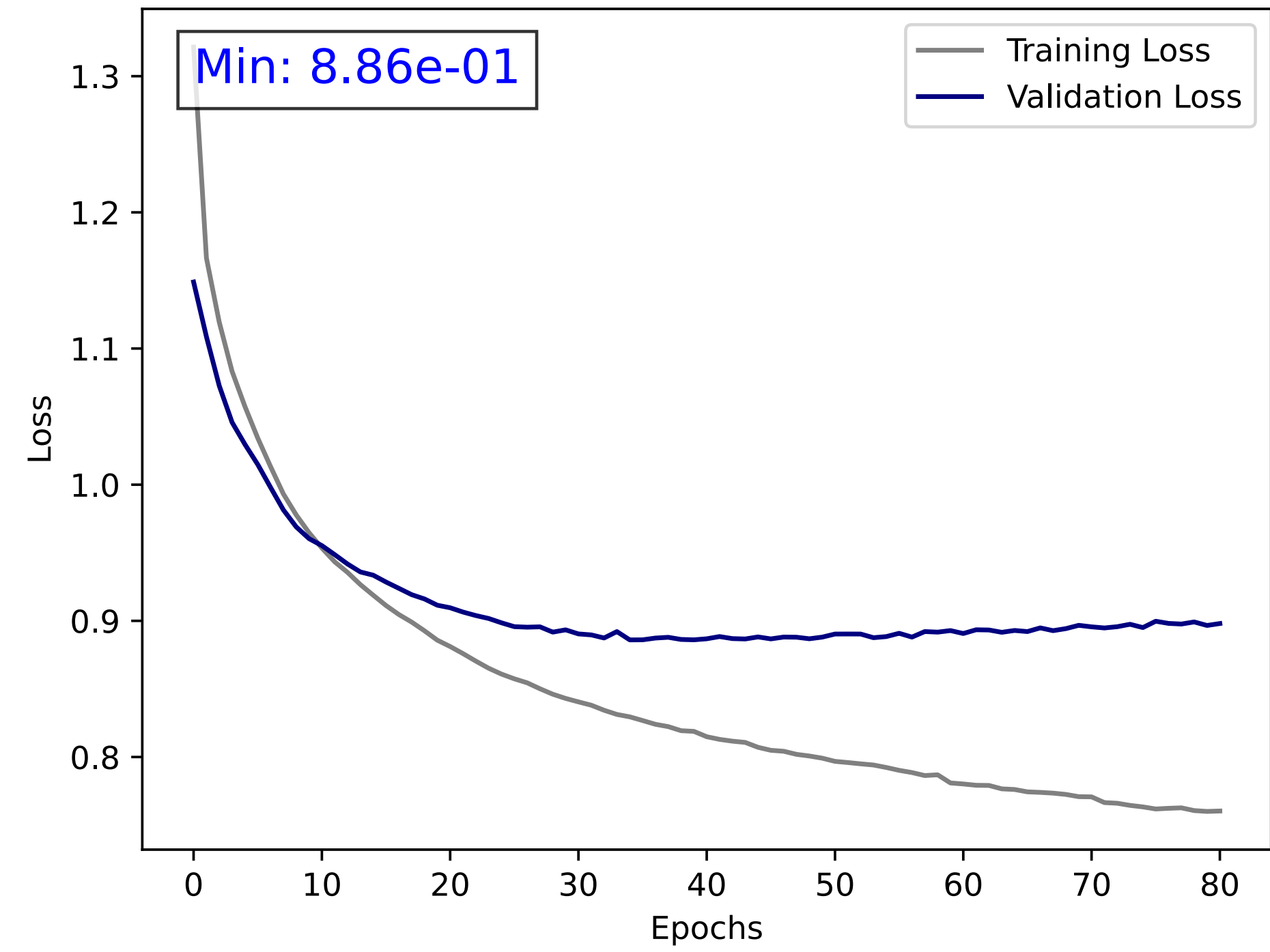




Solo

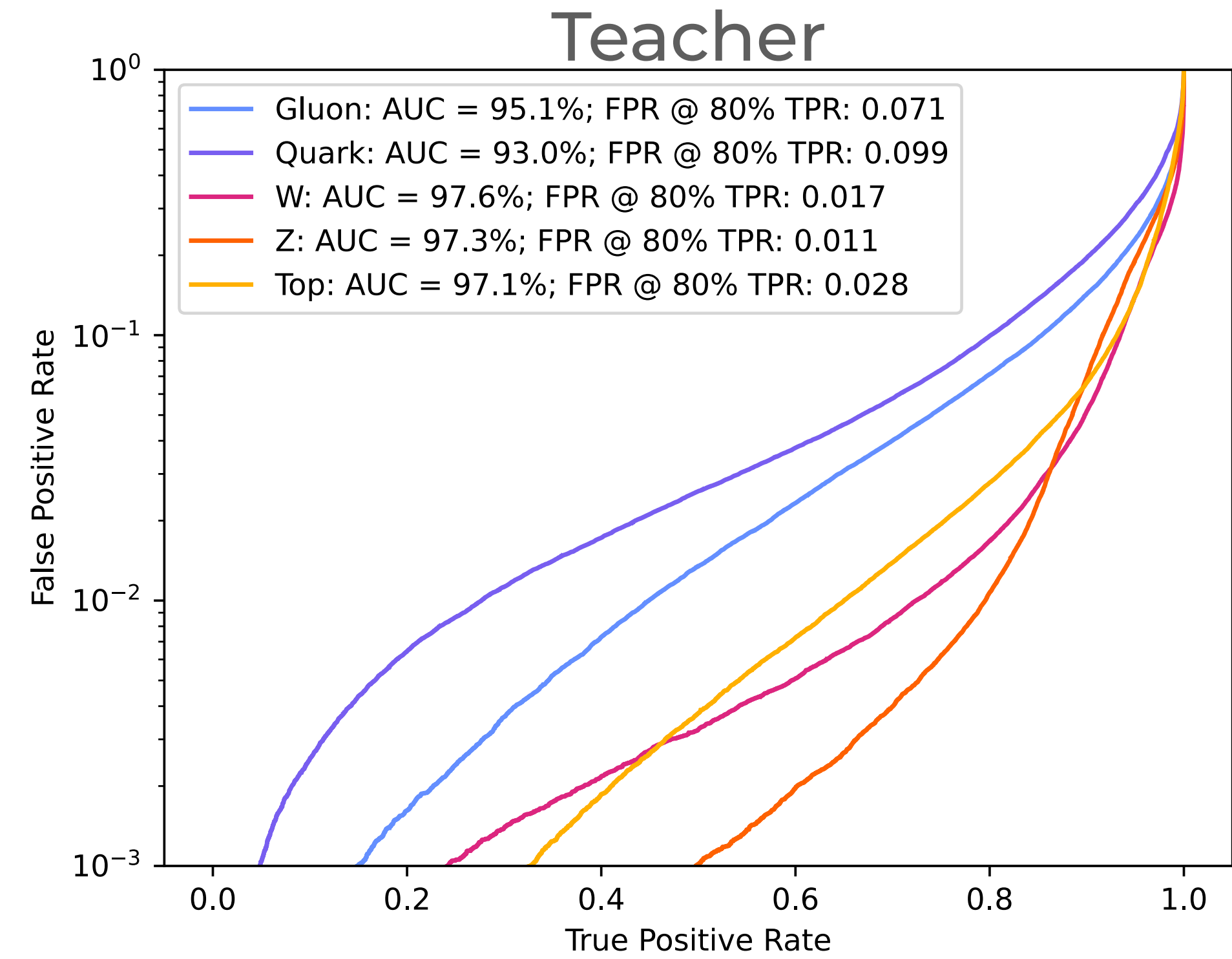
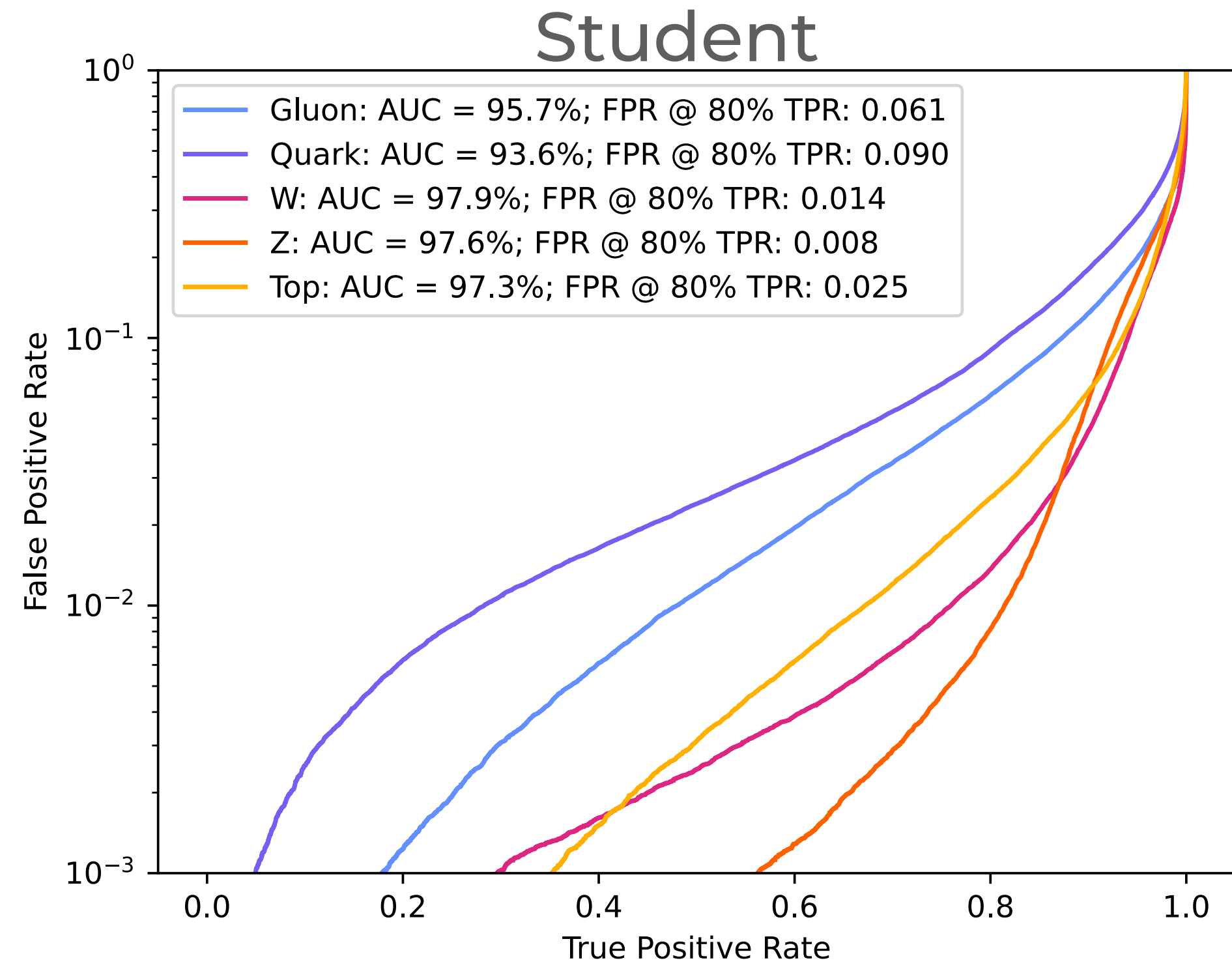


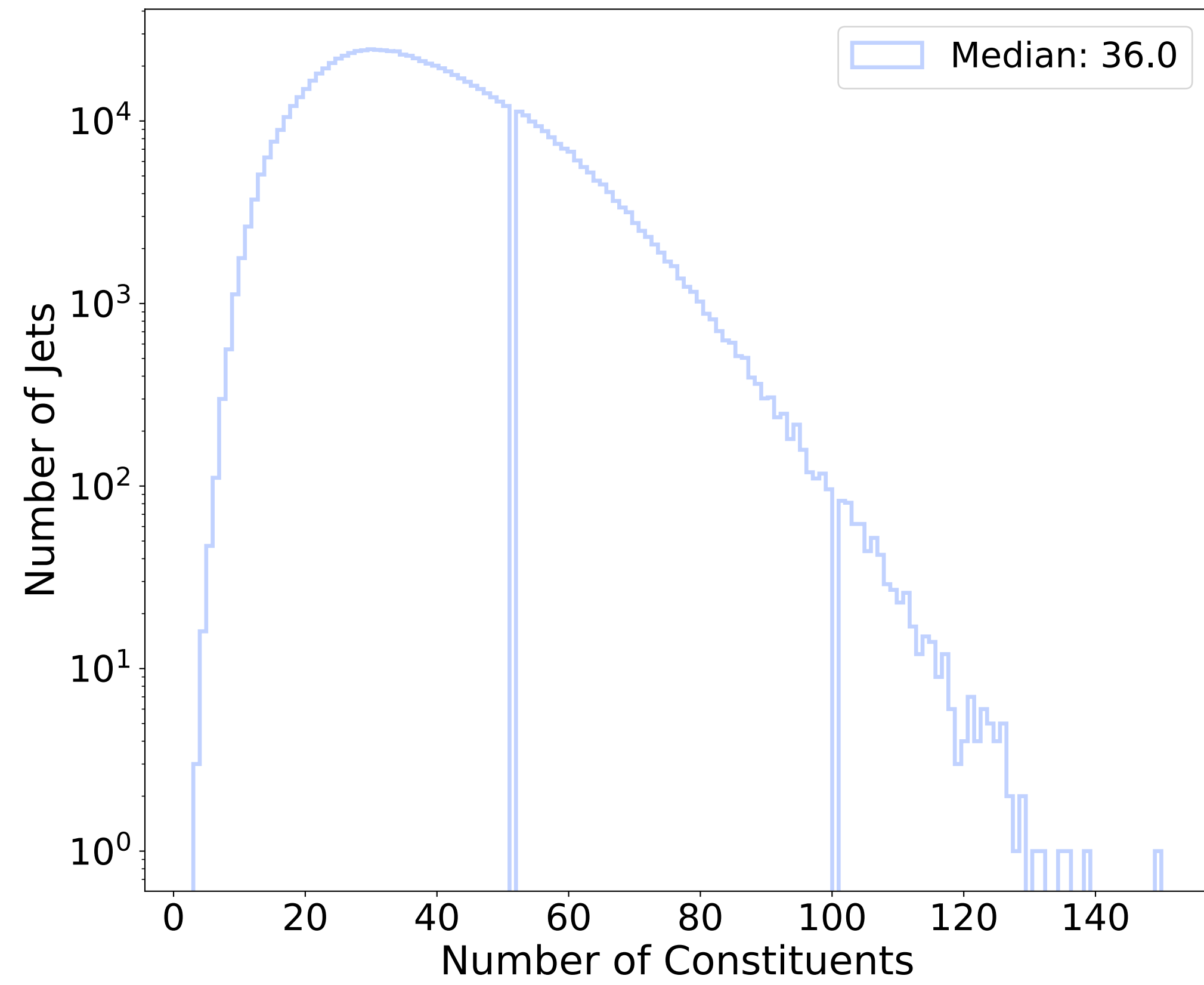
KD



Deep Sets

150 constituents 16 features

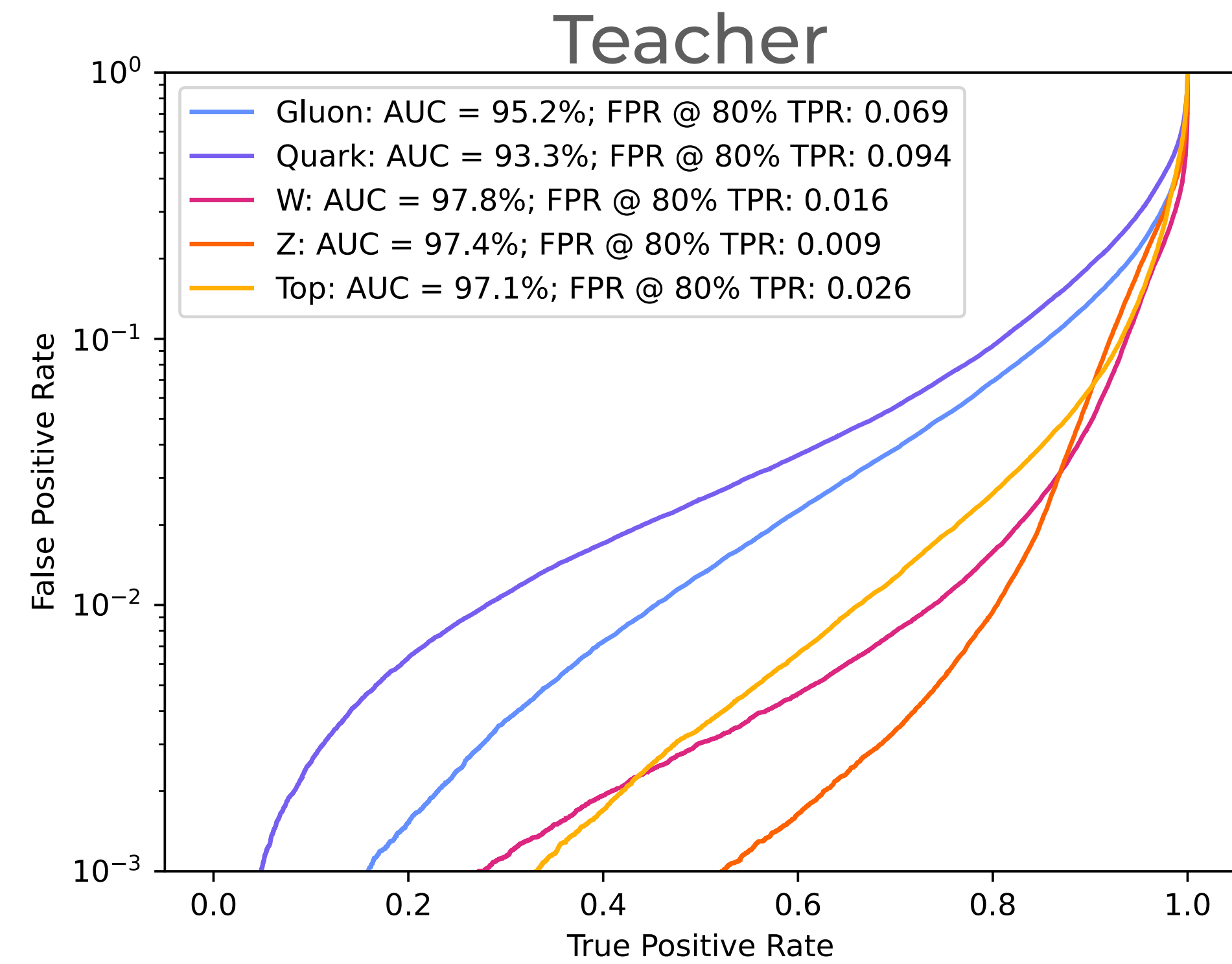
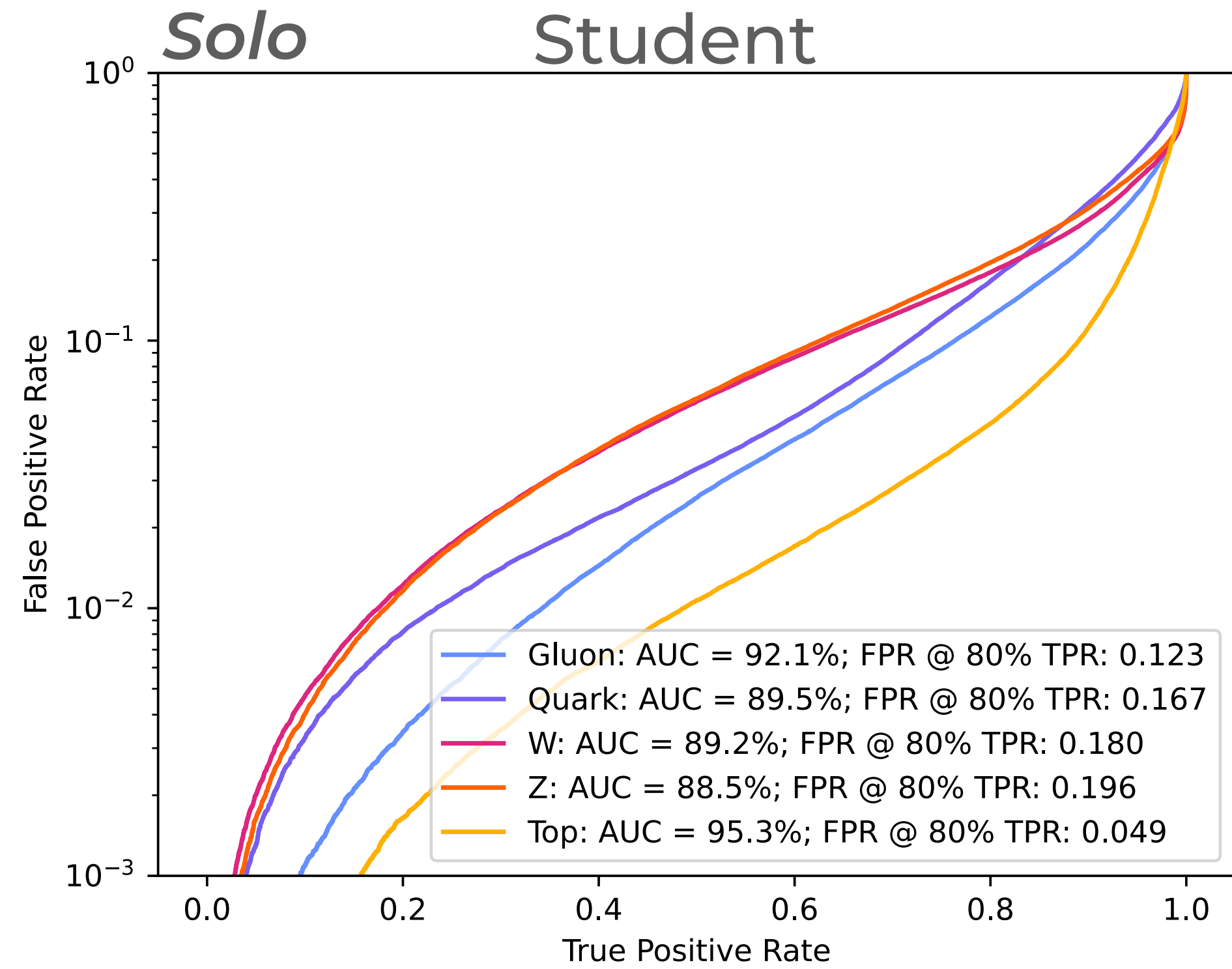


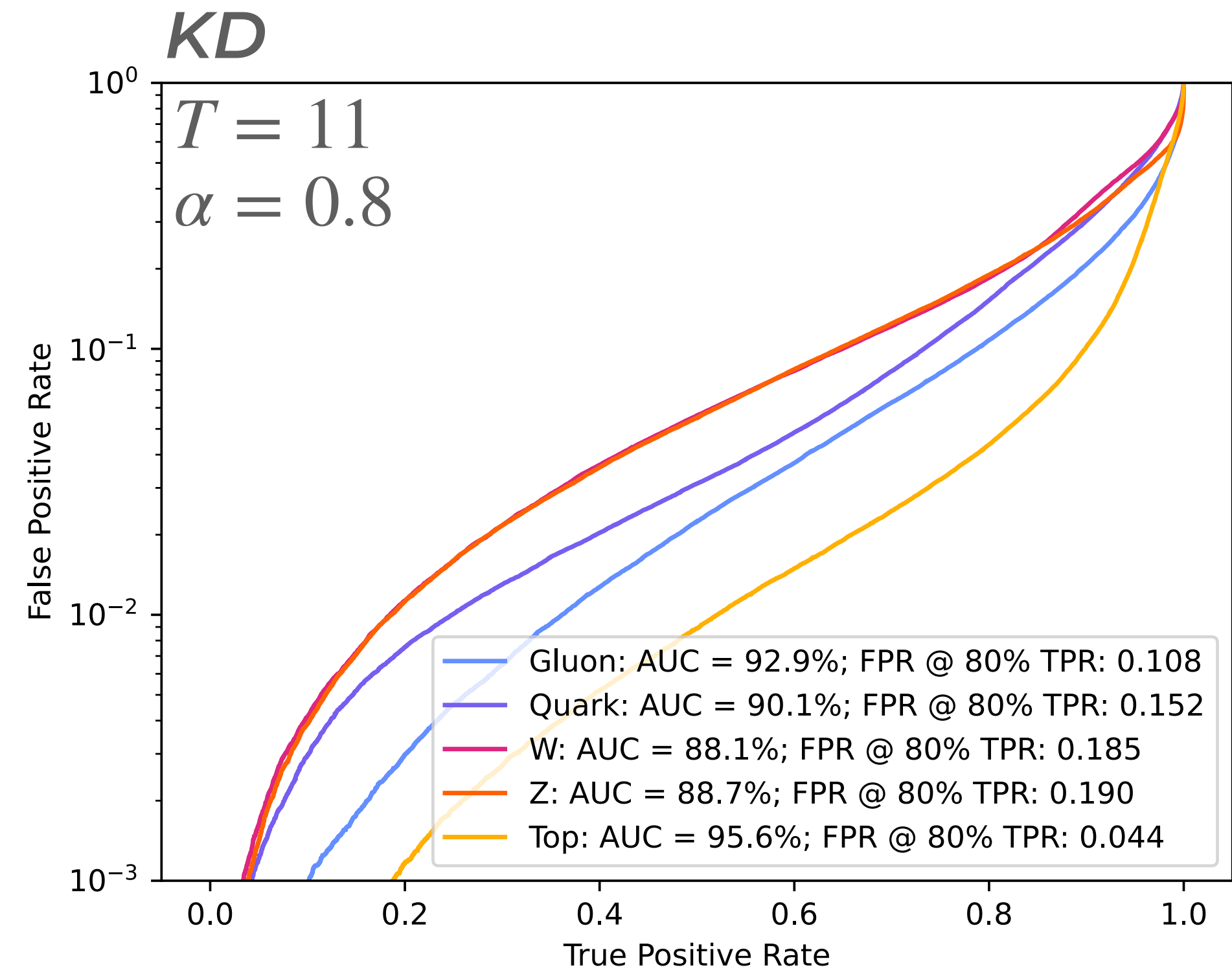
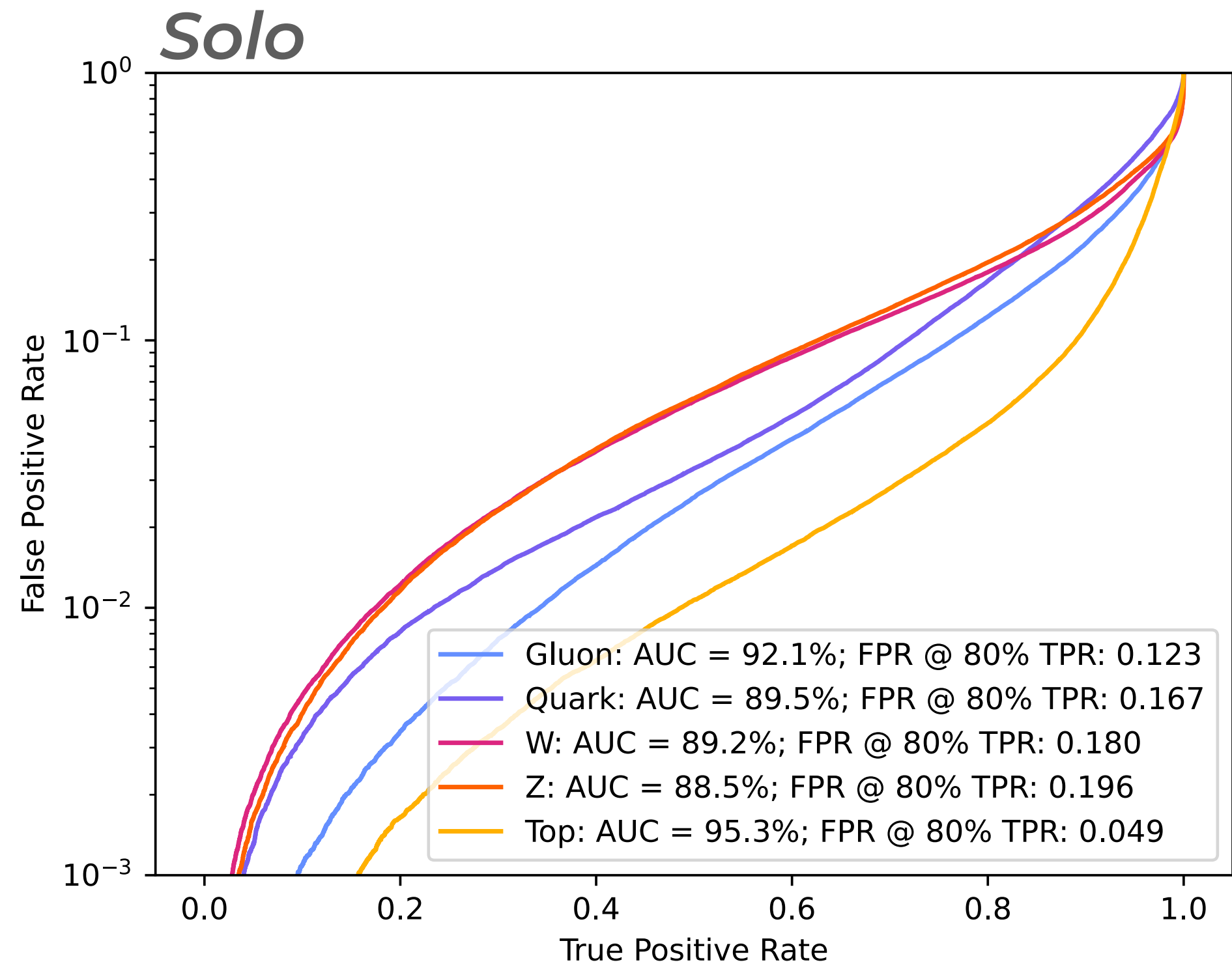


JEDInet DNN

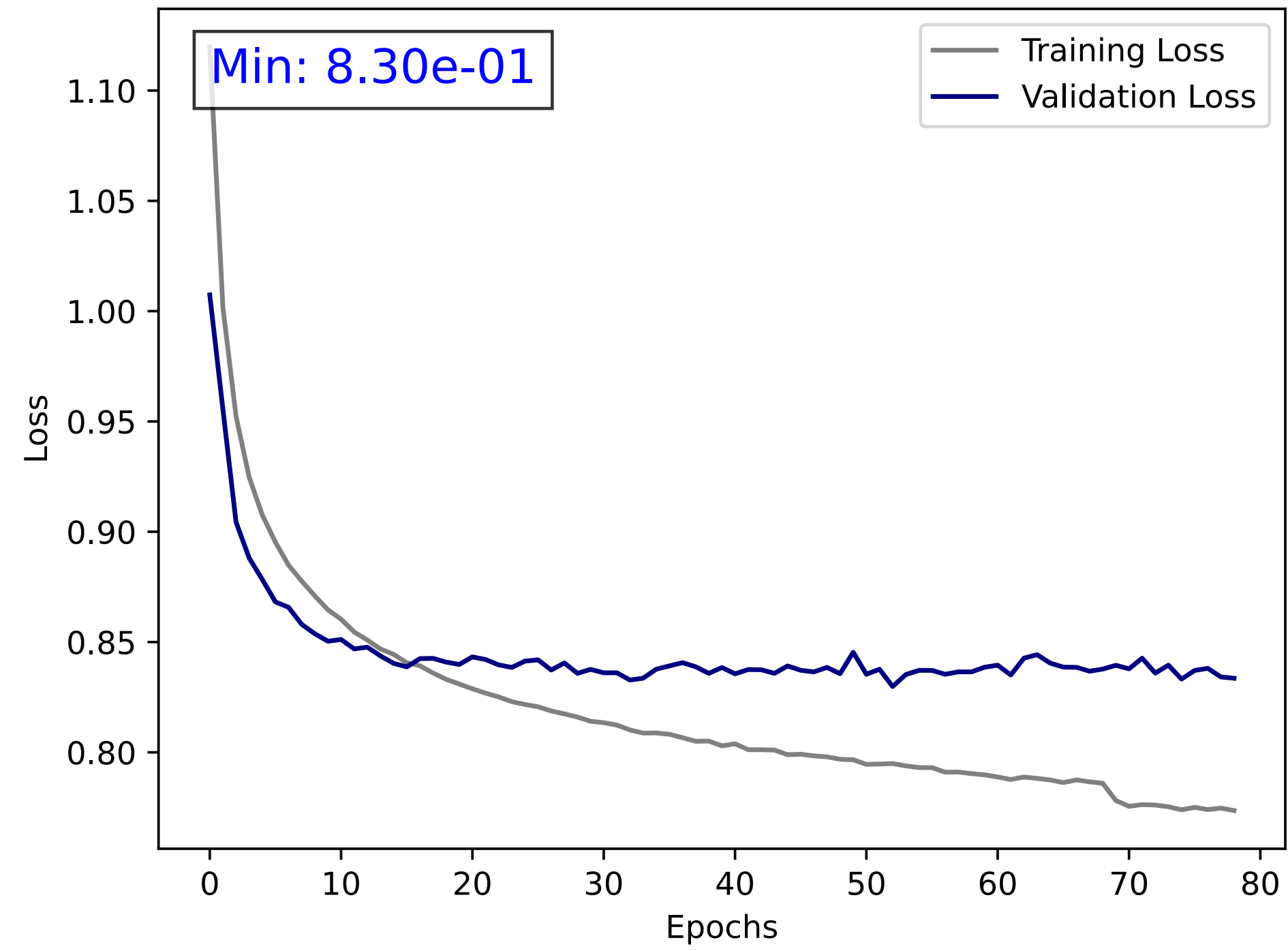
arxiv:1908.05318v3

50 constituents

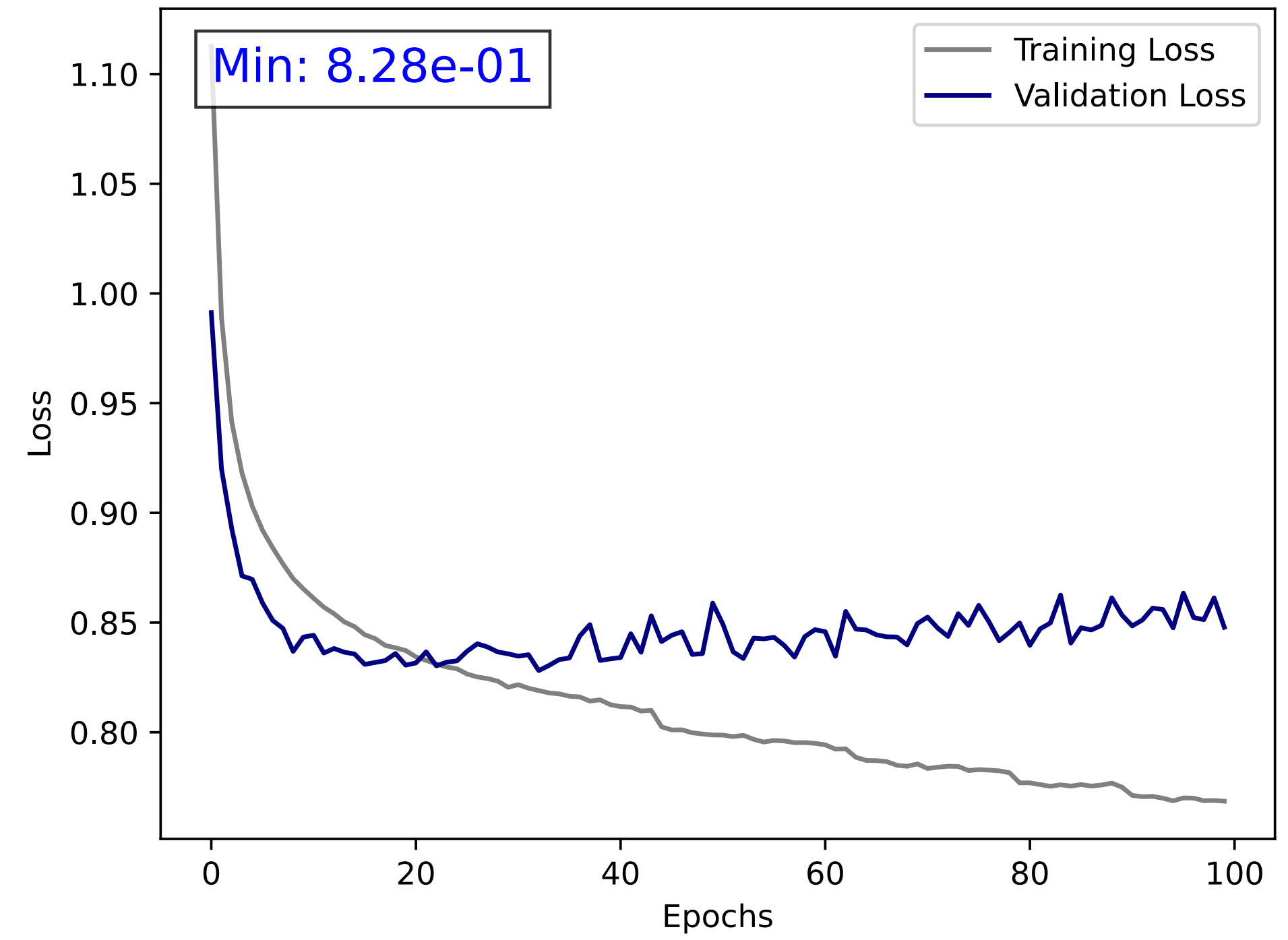




Solo



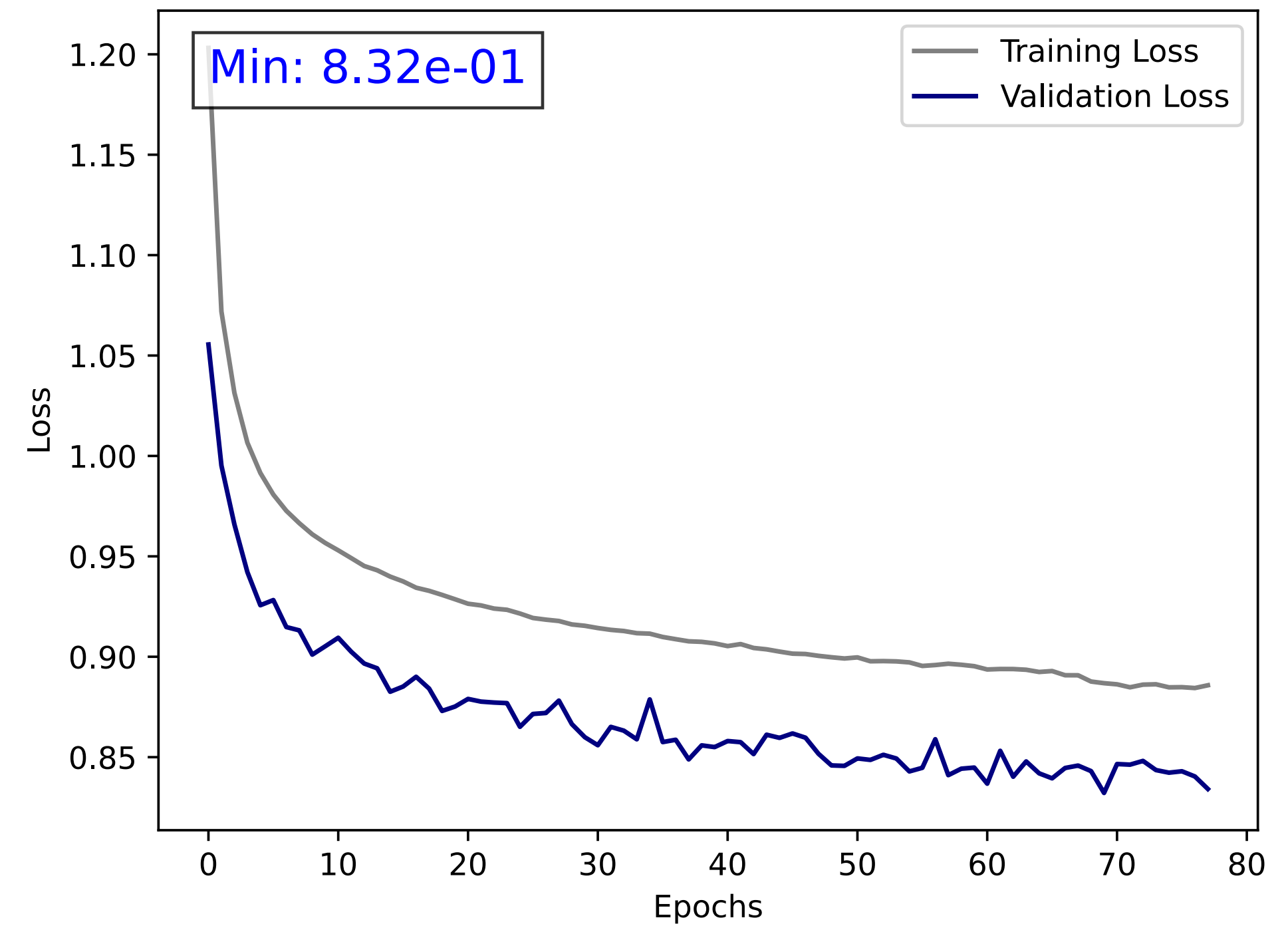
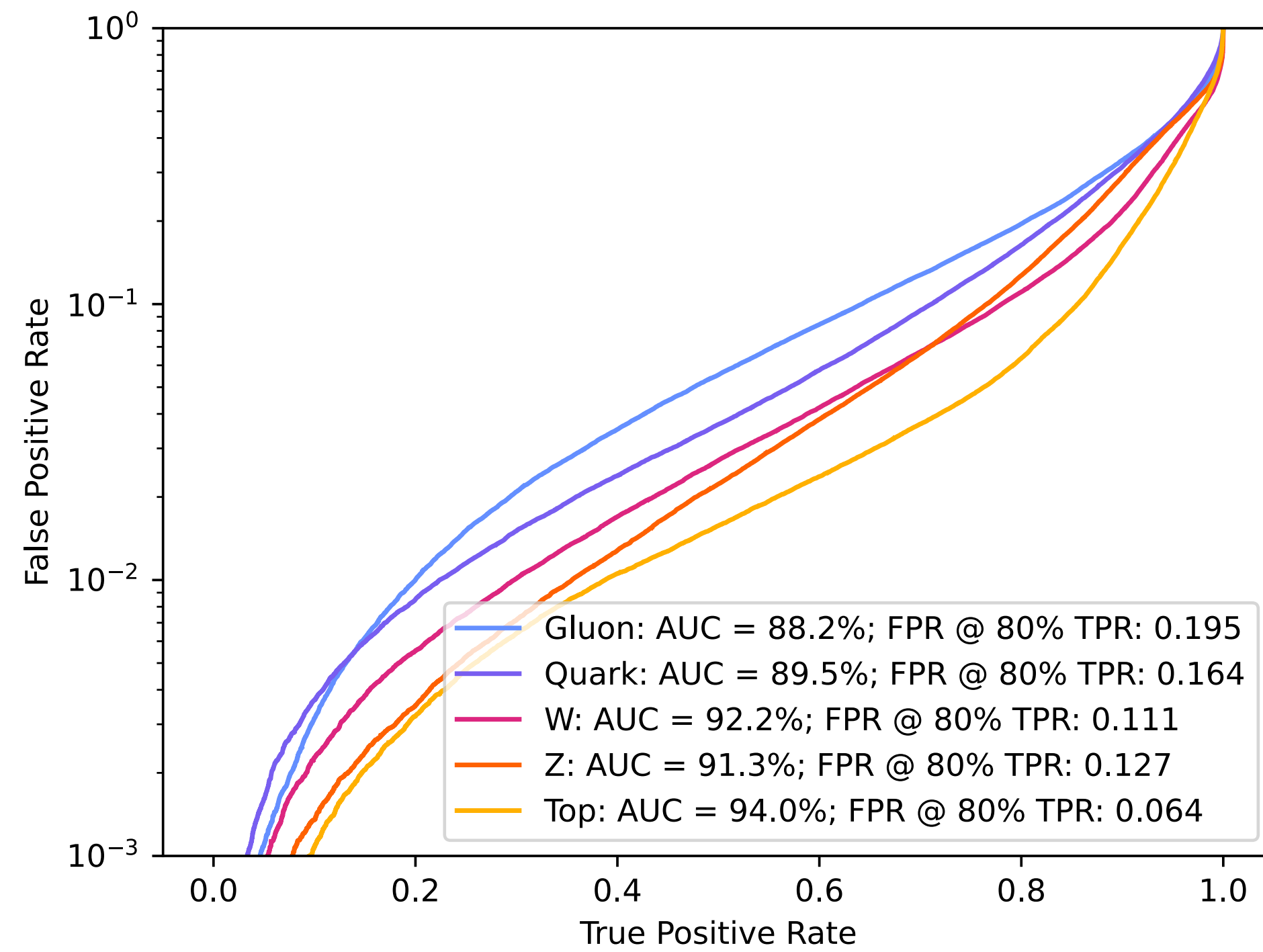
KD



JEDInet DNN

arxiv:1908.05318v3

16 constituents 3 features (pT, eta, phi) - pT ordered

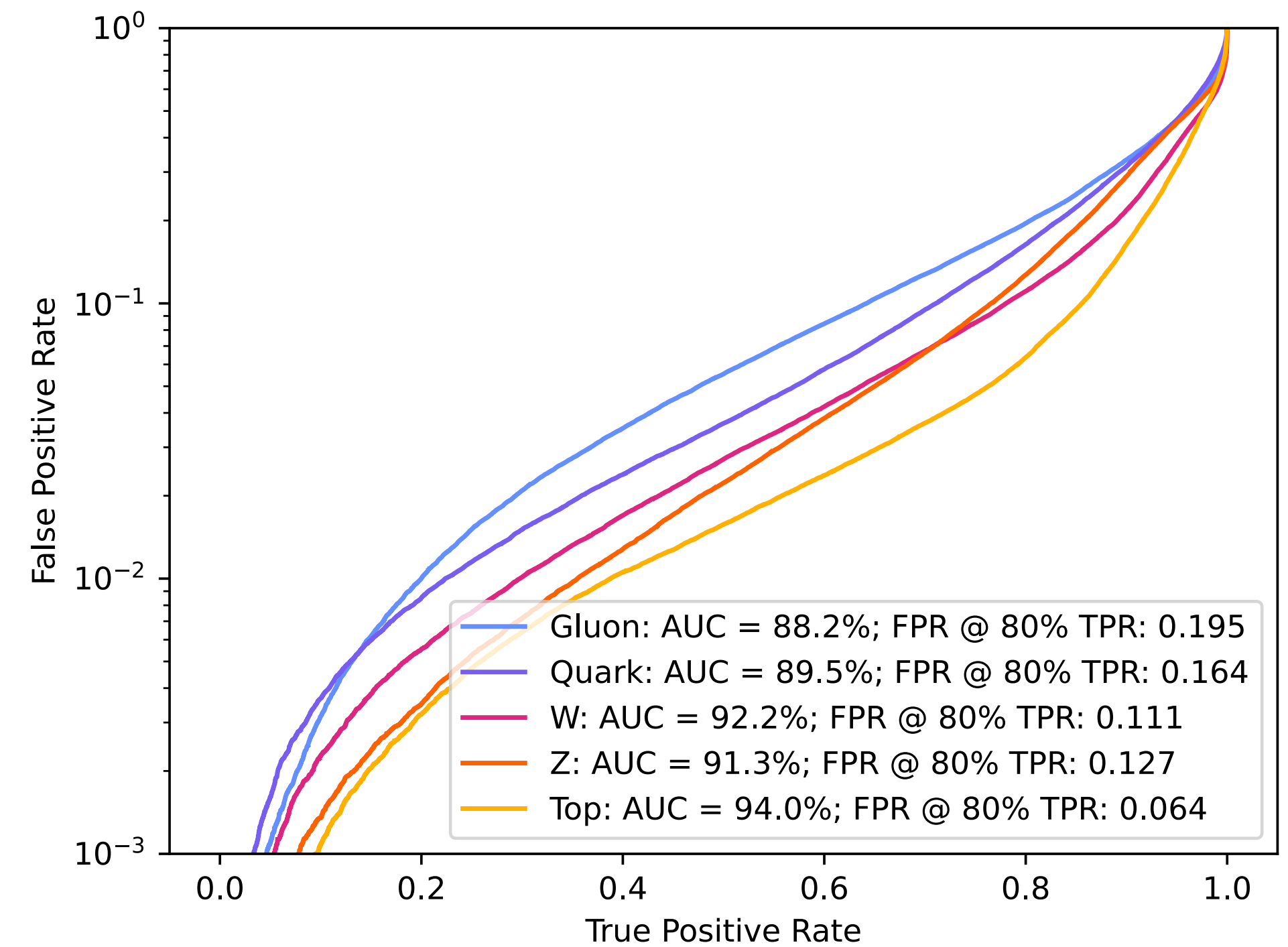
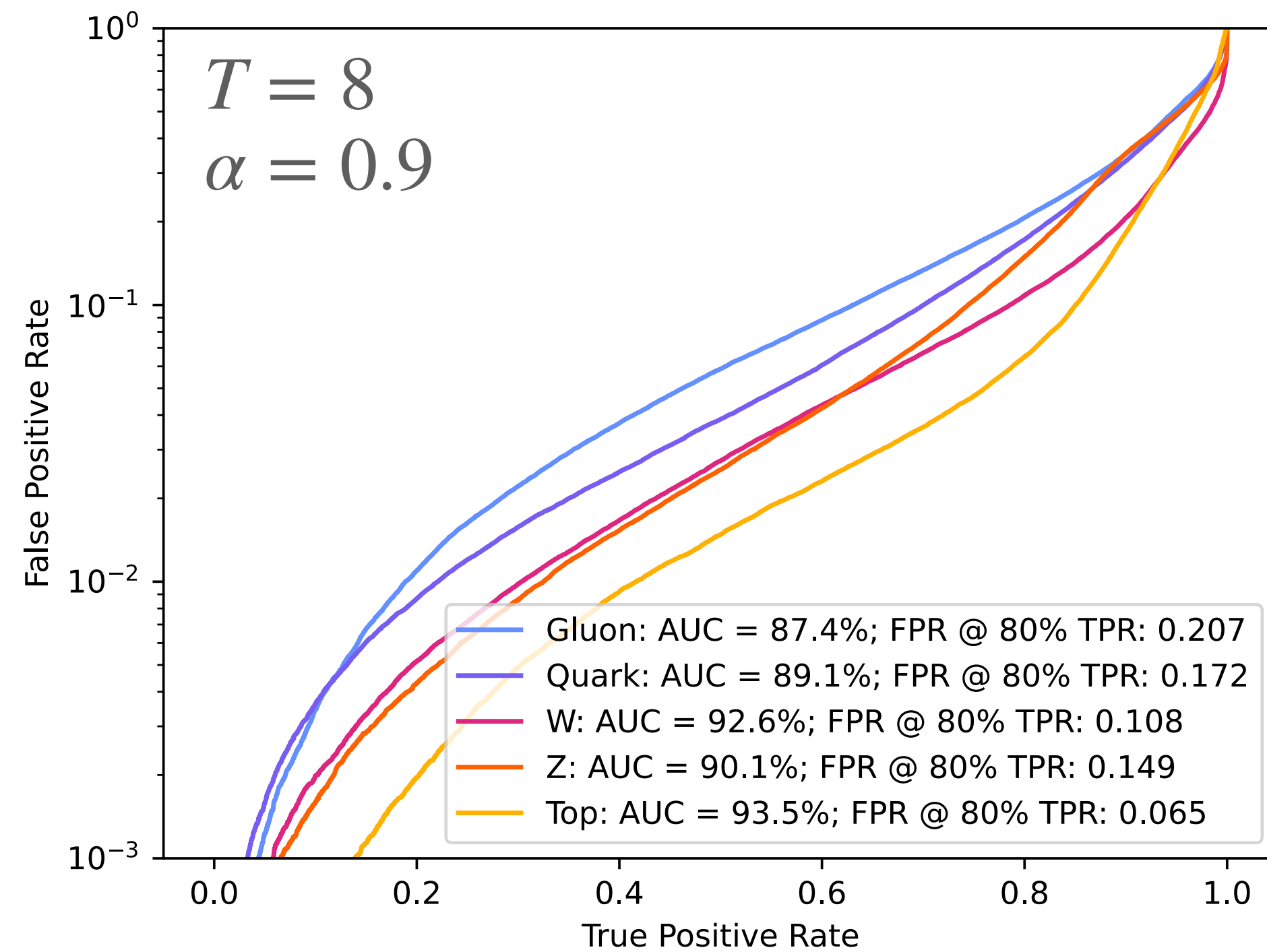


JEDInet DNN

arxiv:1908.05318v3

150 constituents 16 features IntNet teacher

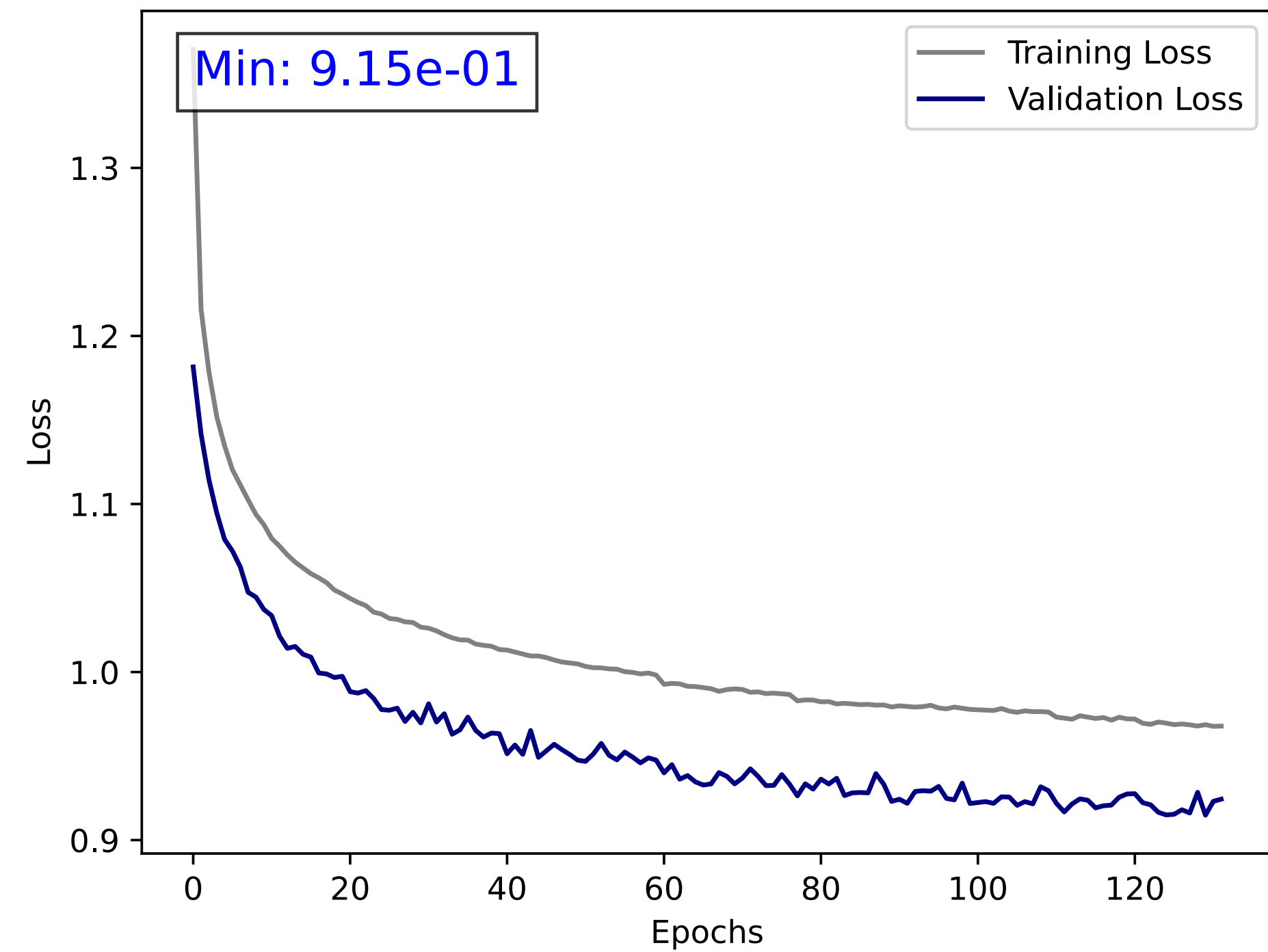
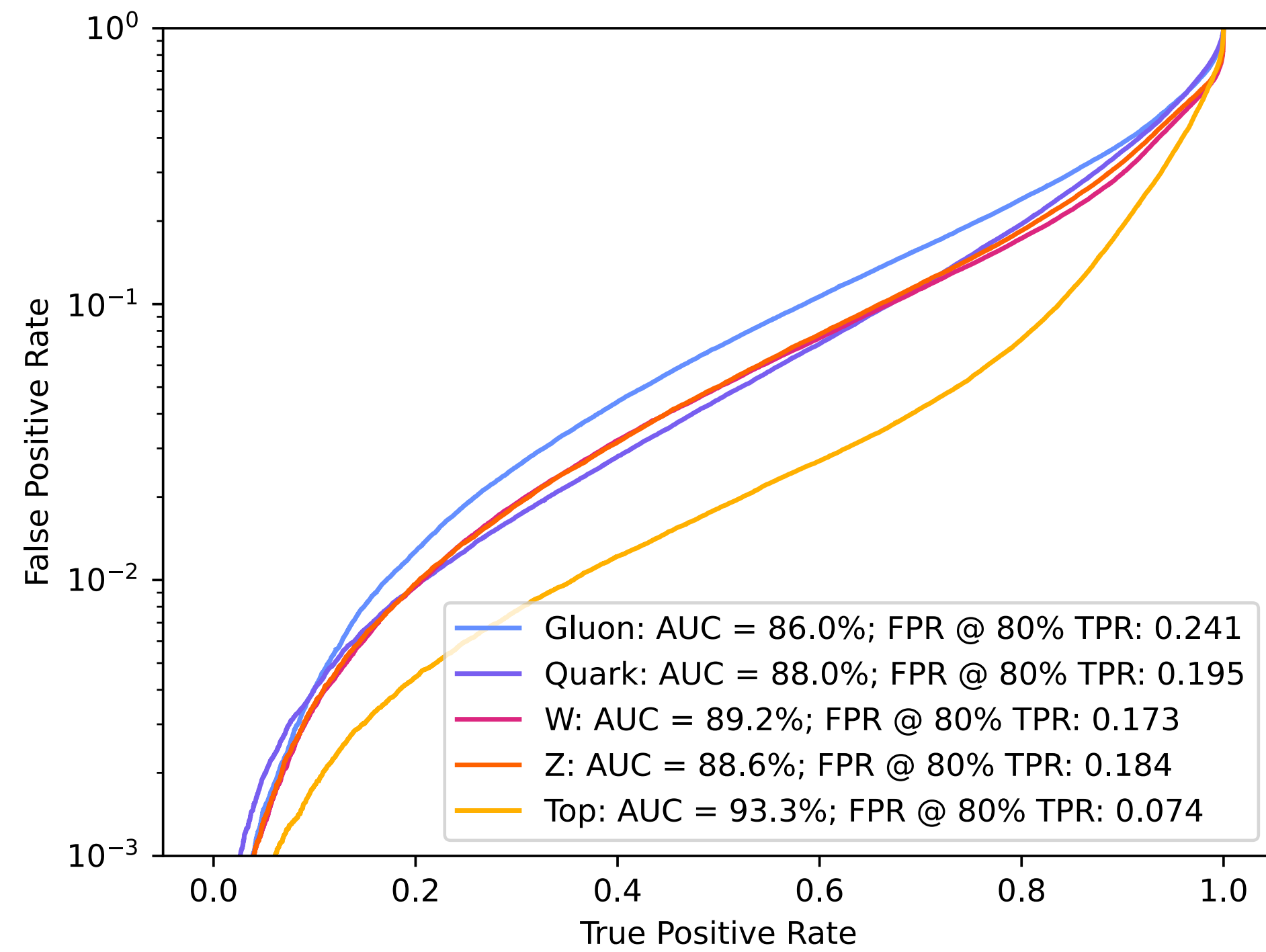
16 constituents 3 feature (pt, eta, phi) - pT ordered



JEDInet DNN

arxiv:1908.05318v3

16 constituents 3 feature (pt, eta, phi) - shuffled



JEDInet DNN

arxiv:1908.05318v3

150 constituents 16 features int net teacher

16 constituents 3 feature (pt, eta, phi) - shuffled

