

Google Research

# ResMem:

Learn what you can and memorize the rest

---

*Speaker: Zitong Yang*

# Collaborators



*Michal Lukasik*



*Vaishnavh Nagarajan*



*Zonglin Li*



*Ankit Rawat*



*Manzil Zaheer*

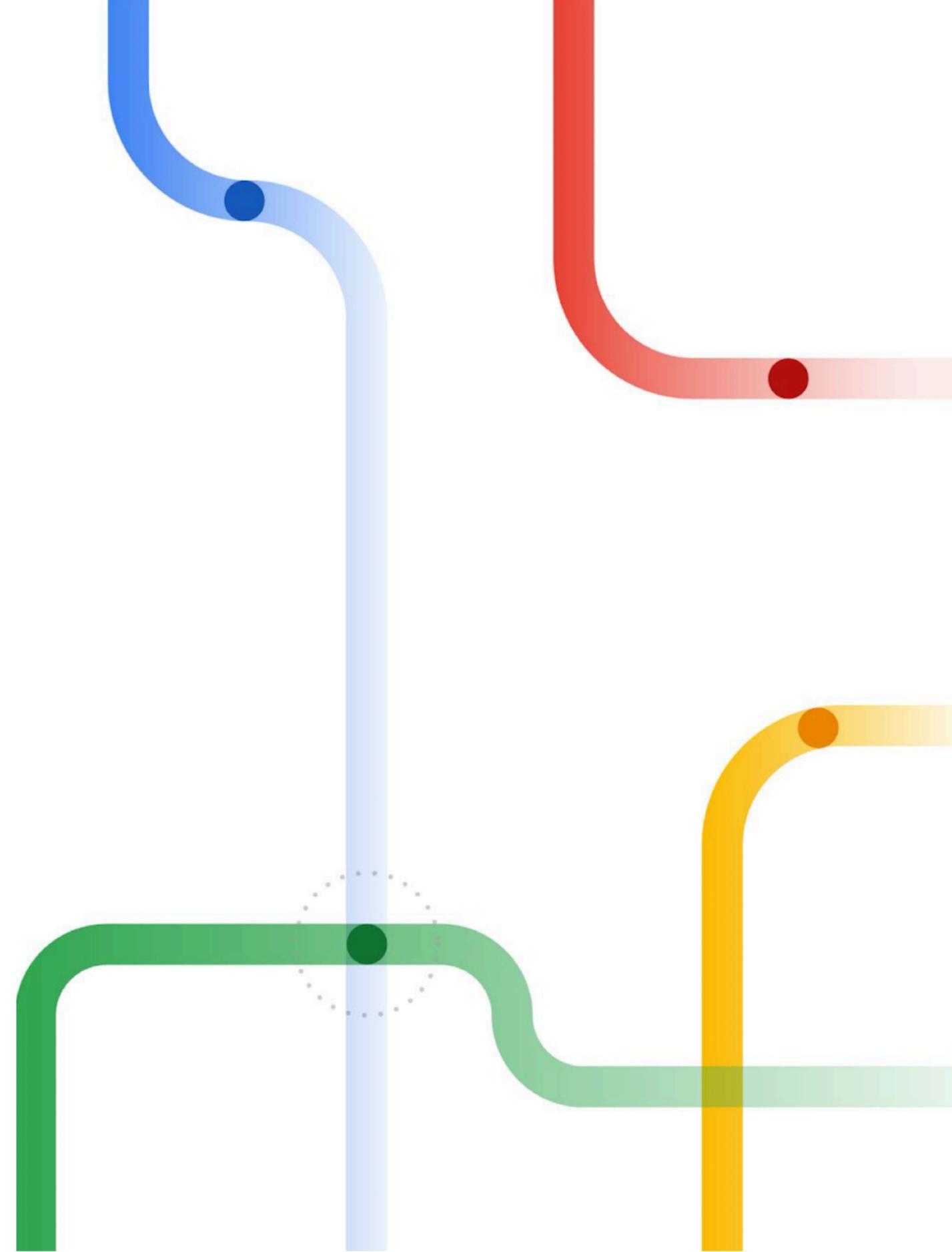


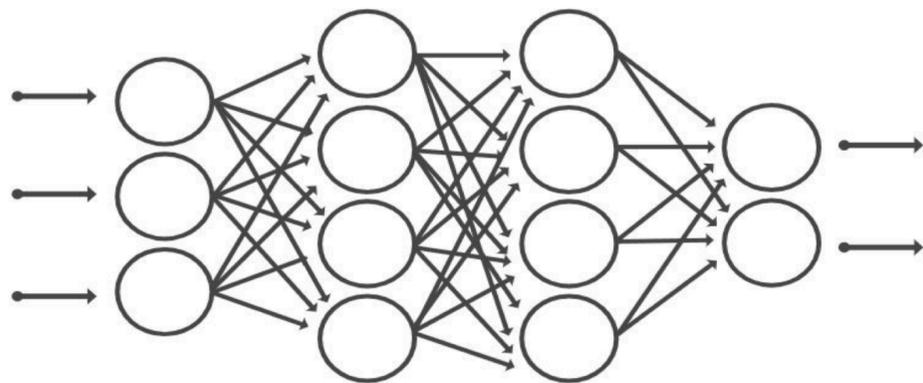
*Aditya Menon*



*Sanjiv Kumar*

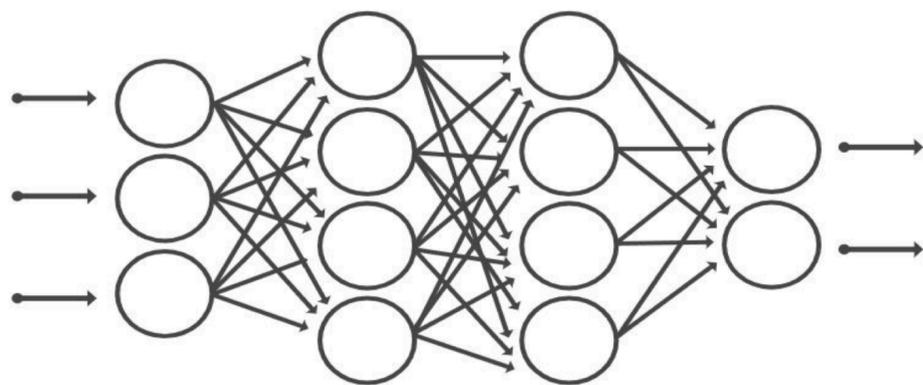
**Google** Research





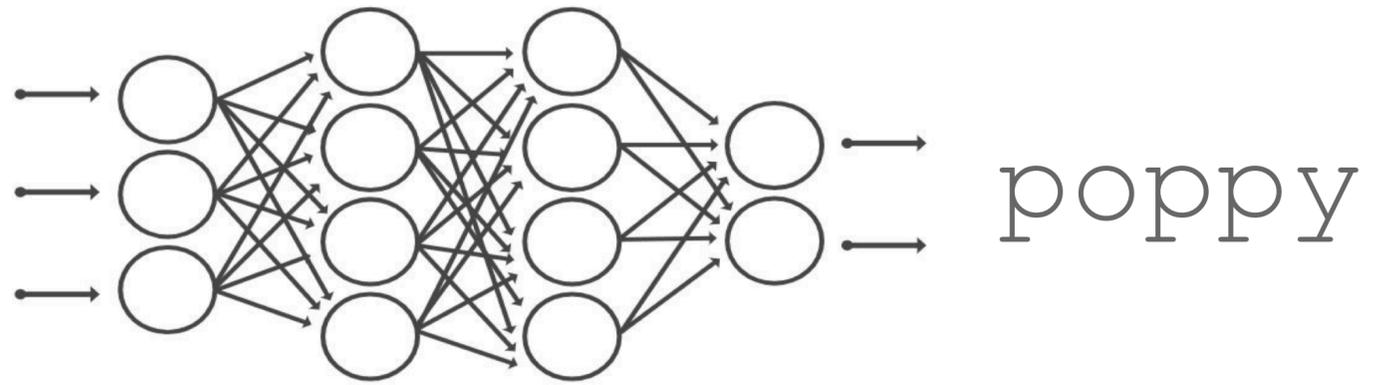


rose



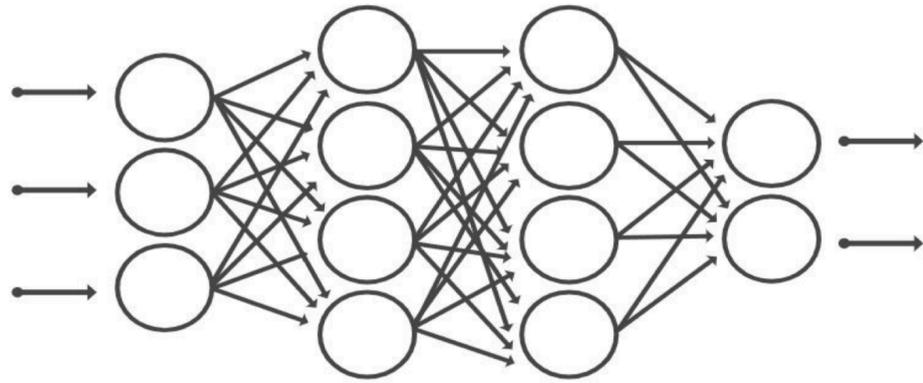


rose

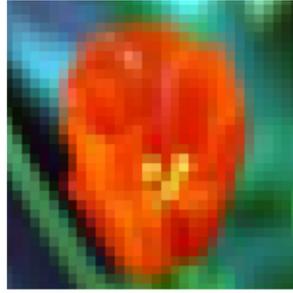




rose

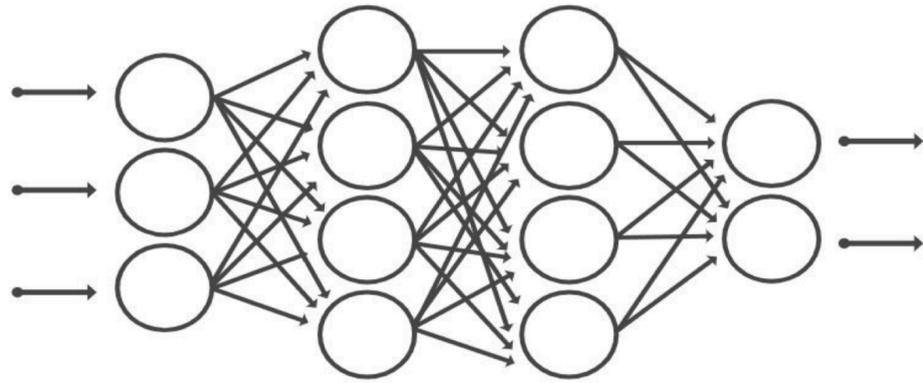


poppy

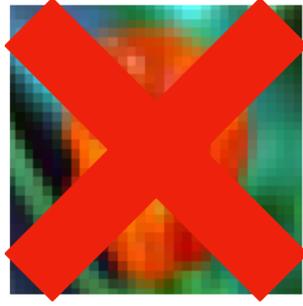




rose



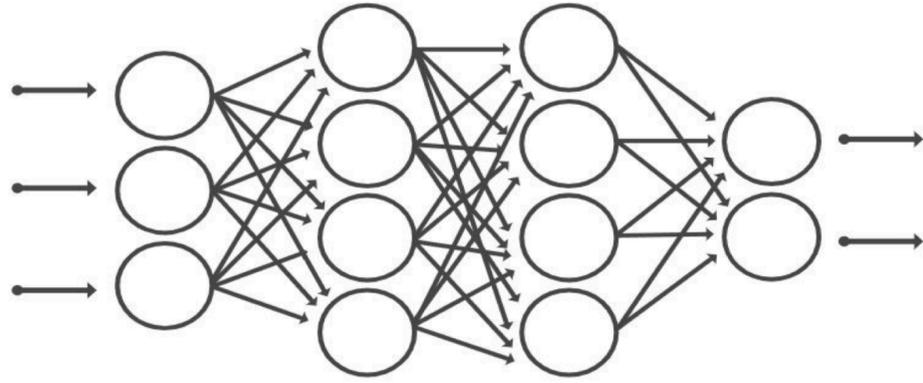
poppy



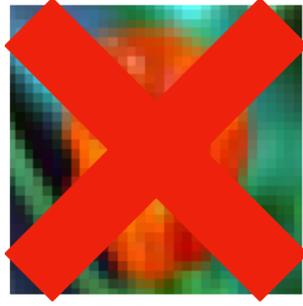
flowers



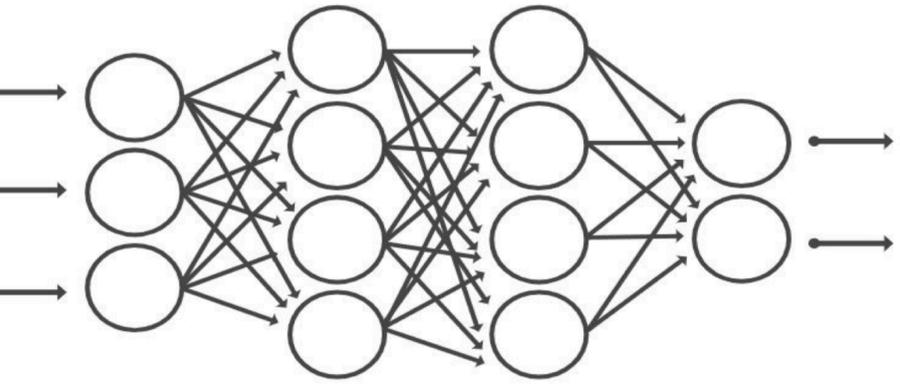
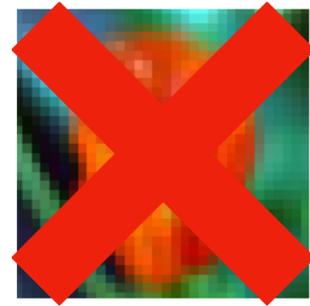
rose



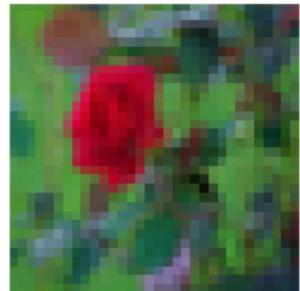
poppy



flowers



poppy



rose

+

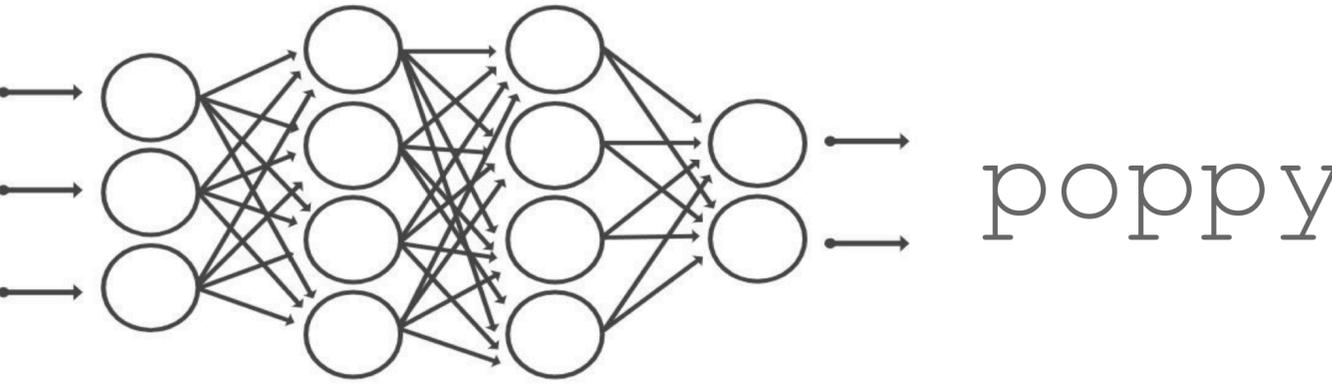
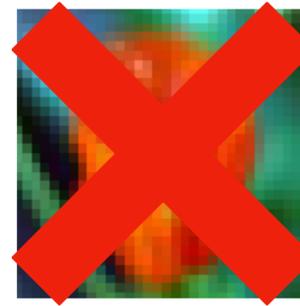


Small perturbation for correct prediction



rose

flowers



rose

+



Small perturbation for correct prediction

**ResMem**



rose



# Background: Memorization $\iff$ Generalization

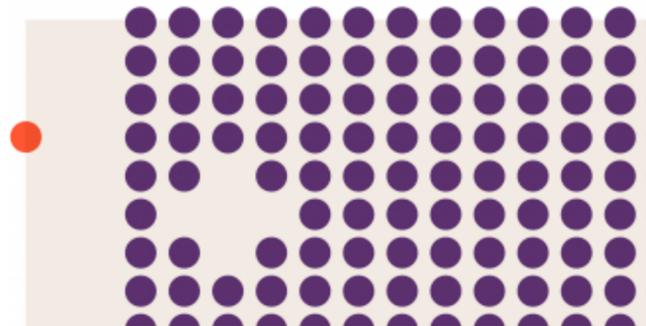
- [**Memorization**  $\implies$  **Generalization**] SoTA large neural networks generalizes well, despite memorizing the training data. (*Zhang et al., '17; Bartlett et al., '20; Kaplan et al., '20*)

# Background: Memorization $\iff$ Generalization

- [**Memorization**  $\implies$  **Generalization**] SoTA large neural networks generalizes well, despite memorizing the training data. (*Zhang et al., '17; Bartlett et al., '20; Kaplan et al., '20*)
- [**Generalization**  $\implies$  **Memorization**] More surprisingly, memorization can be necessary for generalization. (*Feldman et al., '19; Feldman et al., '20; Chen et al., '20*)

# Background: Memorization $\iff$ Generalization

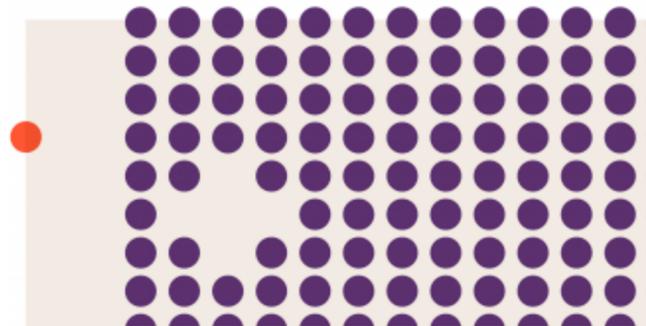
- [**Memorization**  $\implies$  **Generalization**] SoTA large neural networks generalizes well, despite memorizing the training data. (*Zhang et al., '17; Bartlett et al., '20; Kaplan et al., '20*)
- [**Generalization**  $\implies$  **Memorization**] More surprisingly, memorization can be necessary for generalization. (*Feldman et al., '19; Feldman et al., '20; Chen et al., '20*)



*Computationally intractable to distinguish outliers and rare examples.*

# Background: Memorization $\iff$ Generalization

- [**Memorization**  $\implies$  **Generalization**] SoTA large neural networks generalizes well, despite memorizing the training data. (*Zhang et al., '17; Bartlett et al., '20; Kaplan et al., '20*)
- [**Generalization**  $\implies$  **Memorization**] More surprisingly, memorization can be necessary for generalization. (*Feldman et al., '19; Feldman et al., '20; Chen et al., '20*)

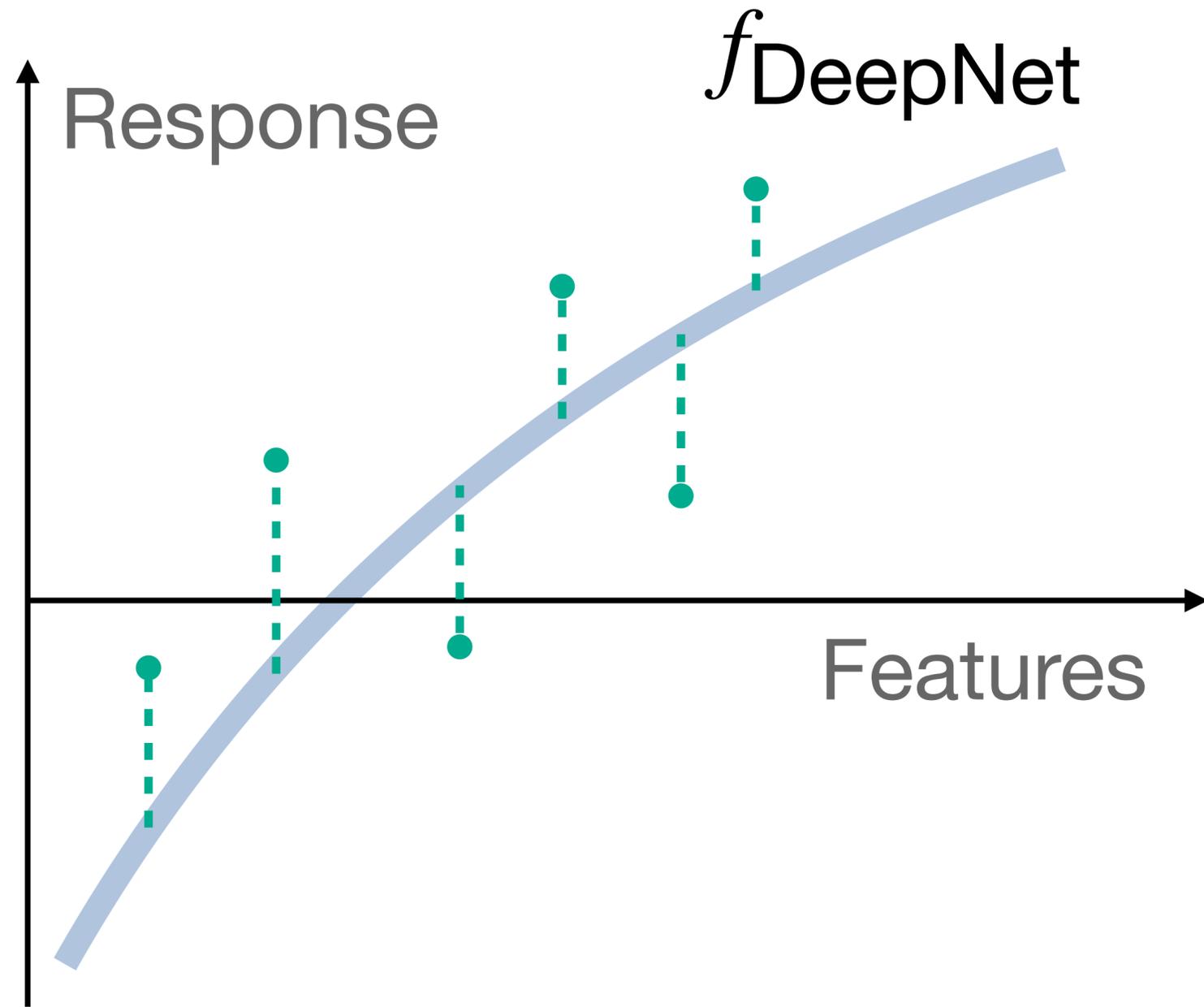


*Computationally intractable to distinguish outliers and rare examples.*

- **Question:** Explicit memorization for generalization?

# ResMem: residual memorization

- **[Step 1]** Train a neural network  $f_{\text{DeepNet}}$  as usual.

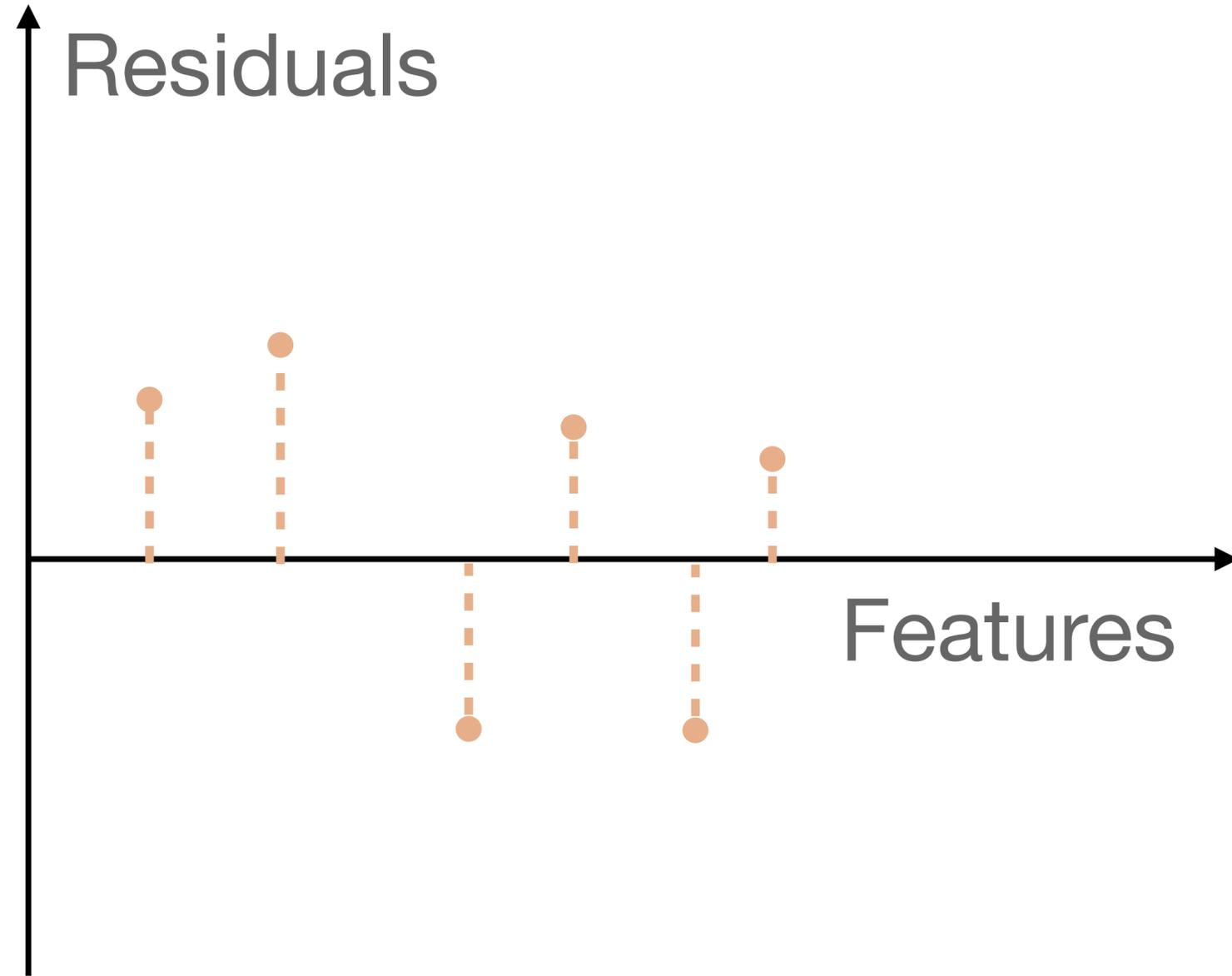


# ResMem: residual memorization

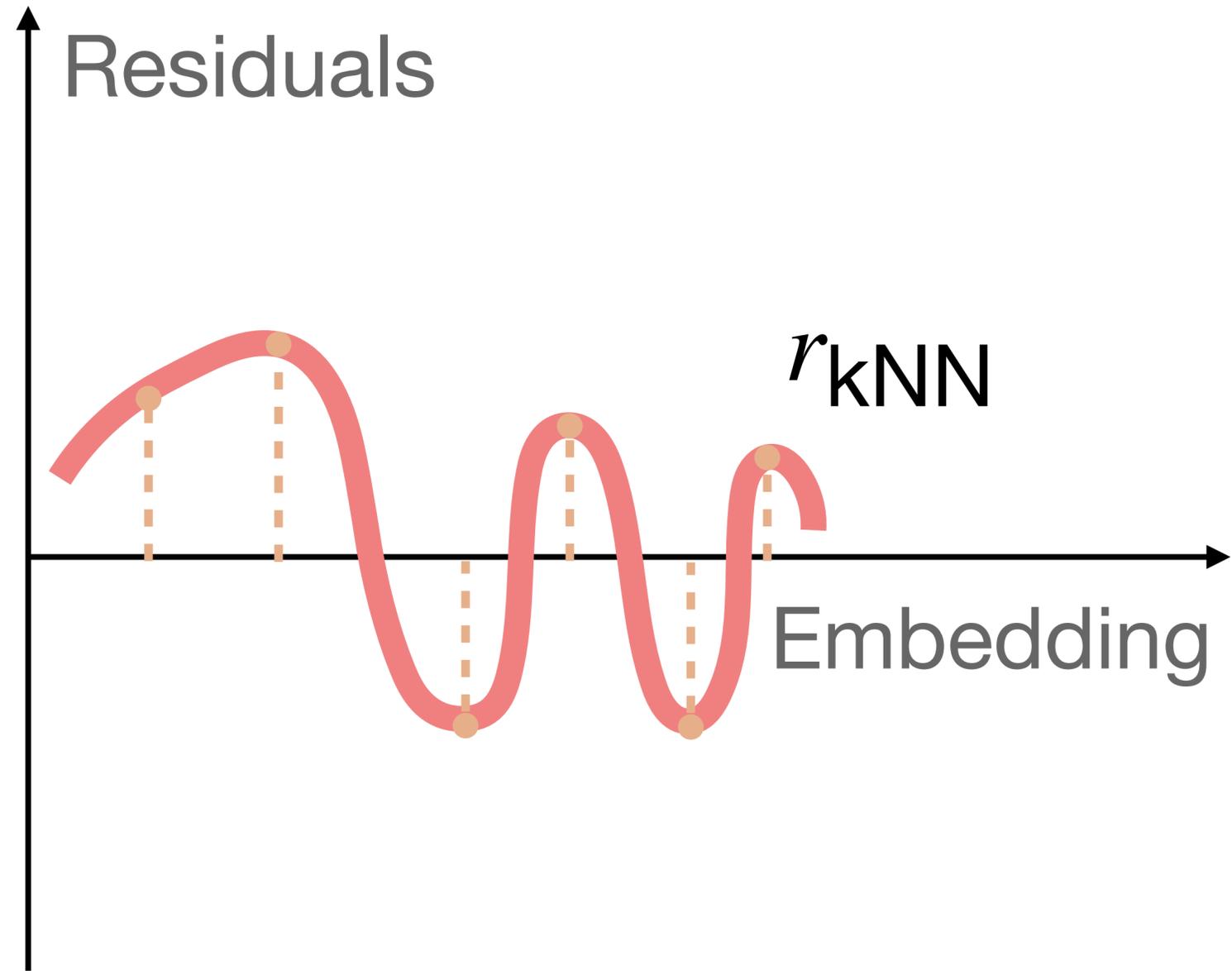
- [Step 1] Train a neural network  $f_{\text{DeepNet}}$  as usual.

- **[Step 2]** Compute the residuals. For classification,

$$r_i = \text{onehot}(y_i) - \text{softmax} \left( f_{\text{DeepNet}}(x_i) \right).$$



# ResMem: residual memorization



- [Step 1] Train a neural network  $f_{\text{DeepNet}}$  as usual.

- [Step 2] Compute the residuals. For classification,

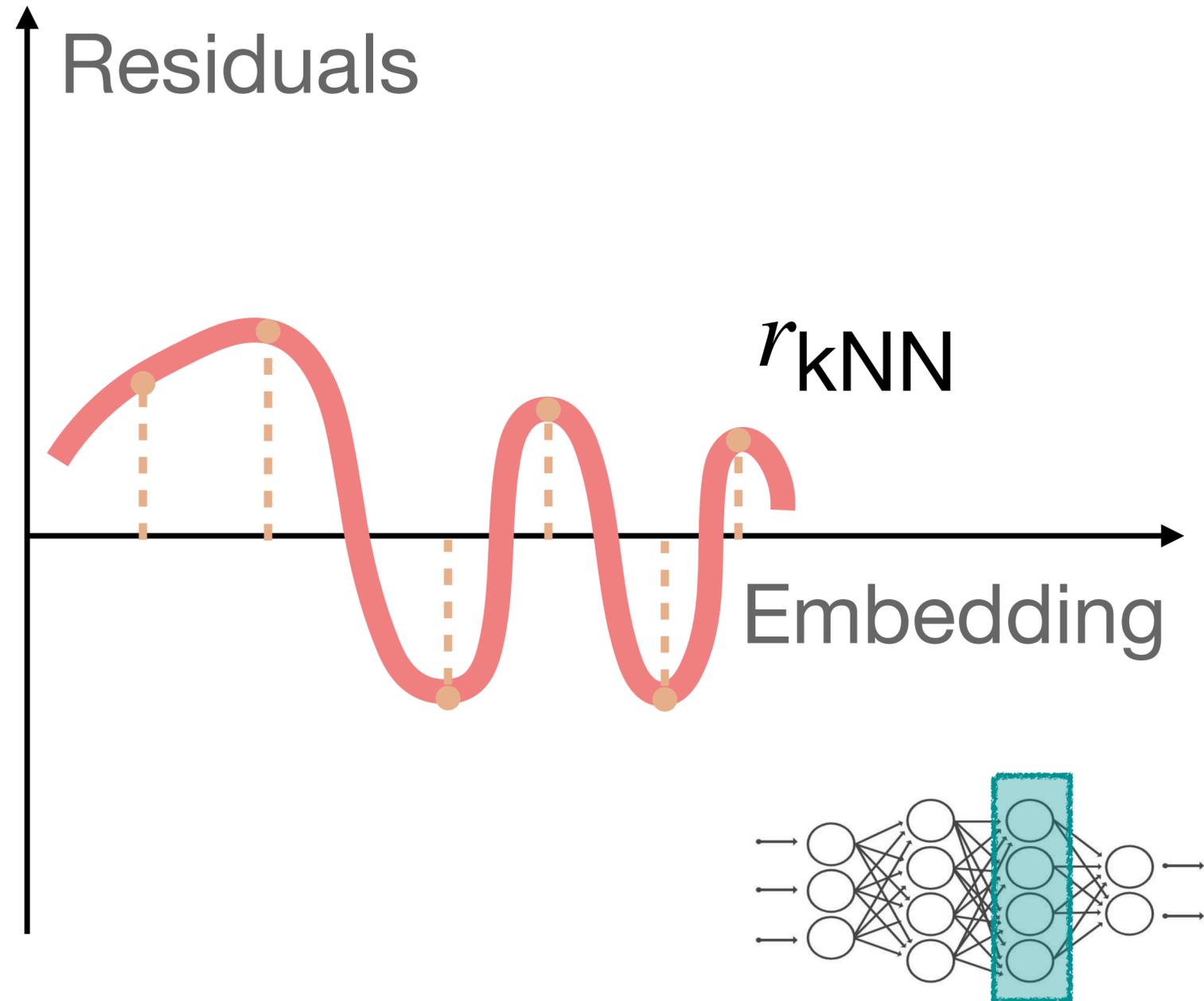
$$r_i = \text{onehot}(y_i) - \text{softmax} \left( f_{\text{DeepNet}}(x_i) \right).$$

- **[Step 3]** Memorize residuals using a  $k$ -NN

$$r_{\text{kNN}}(\tilde{x}) = \sum_i r_i \cdot w_i,$$

where  $w_i \sim -\|\phi(\tilde{x}) - \phi(x_i)\|$  is computed from the intermediate layer of the neural network.

# ResMem: residual memorization



- [Step 1] Train a neural network  $f_{\text{DeepNet}}$  as usual.

- [Step 2] Compute the residuals. For classification,

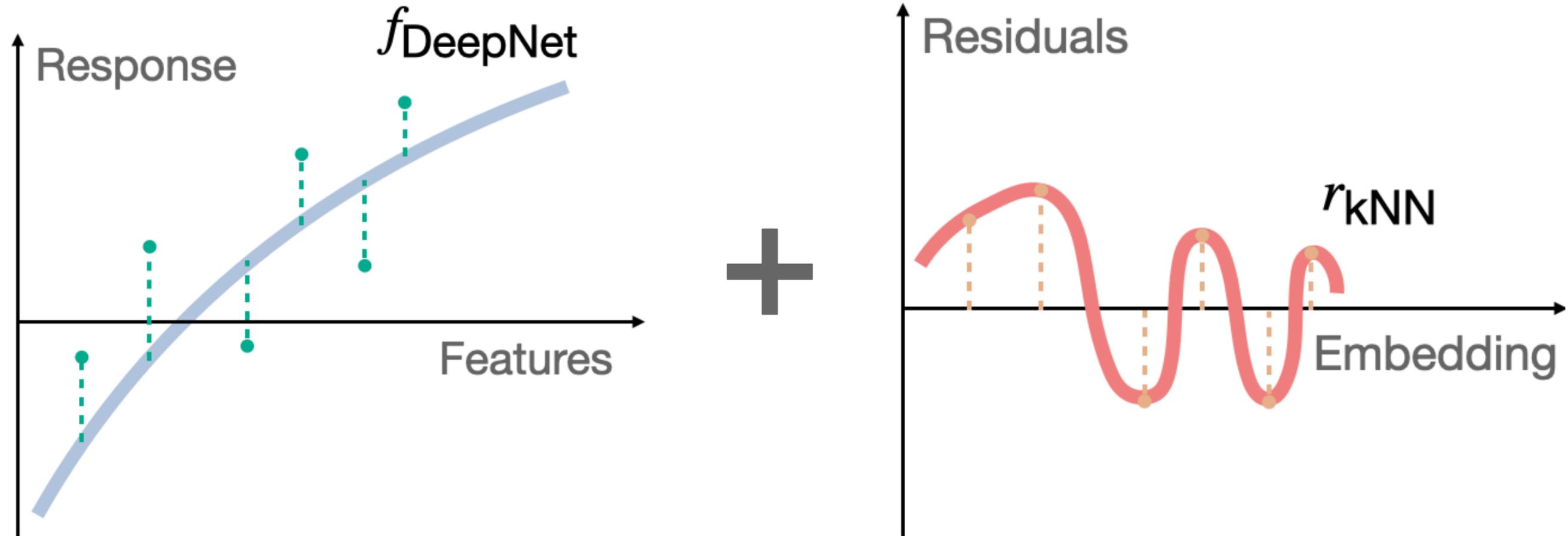
$$r_i = \text{onehot}(y_i) - \text{softmax} \left( f_{\text{DeepNet}}(x_i) \right).$$

- **[Step 3]** Memorize residuals using a  $k$ -NN

$$r_{\text{kNN}}(\tilde{x}) = \sum_i r_i \cdot w_i,$$

where  $w_i \sim -\|\phi(\tilde{x}) - \phi(x_i)\|$  is computed from the intermediate layer of the neural network.

# ResMem: residual memorization



$$\text{ResMem prediction} = f_{\text{DeepNet}} + r_{\text{kNN}}$$

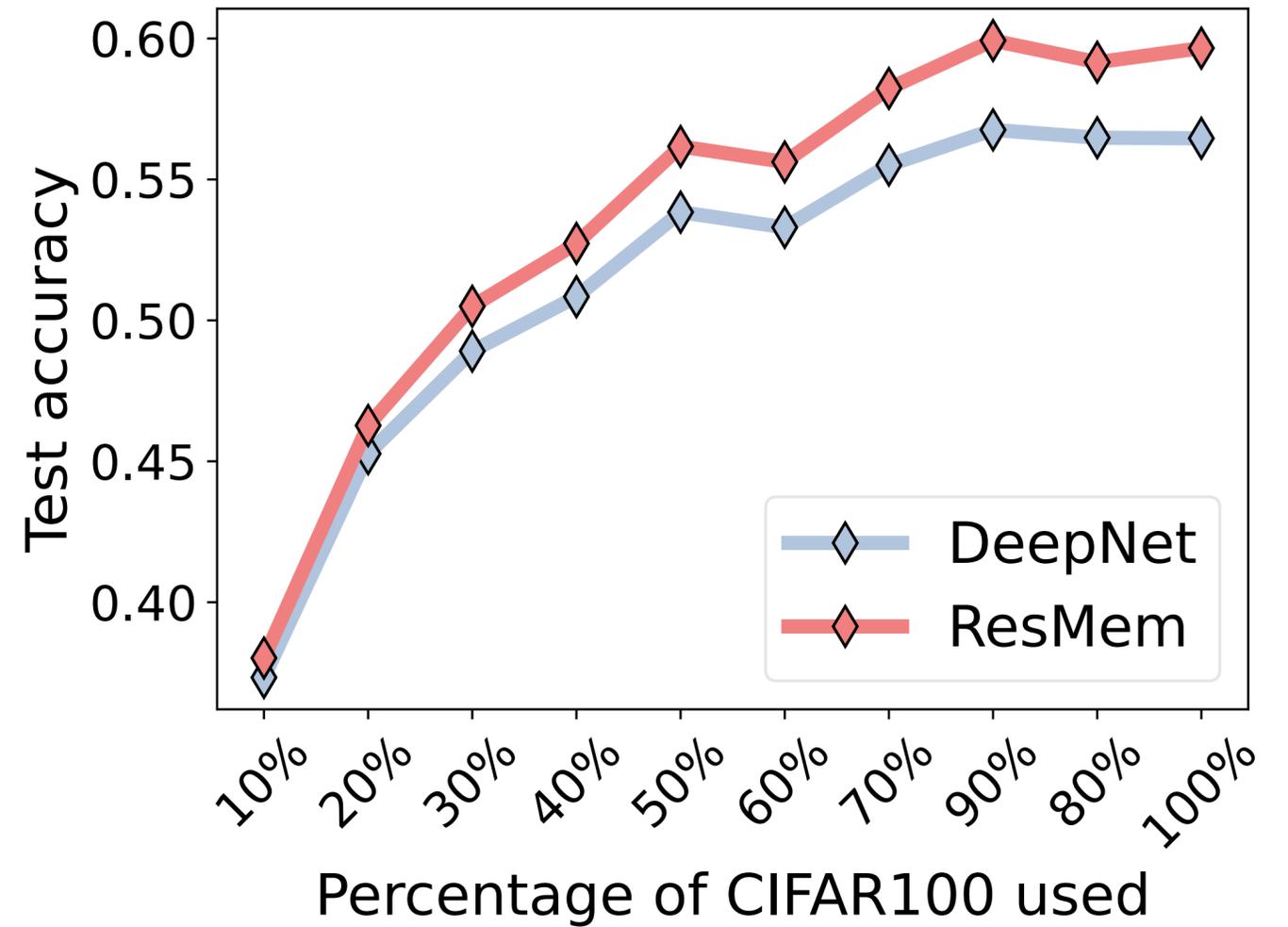
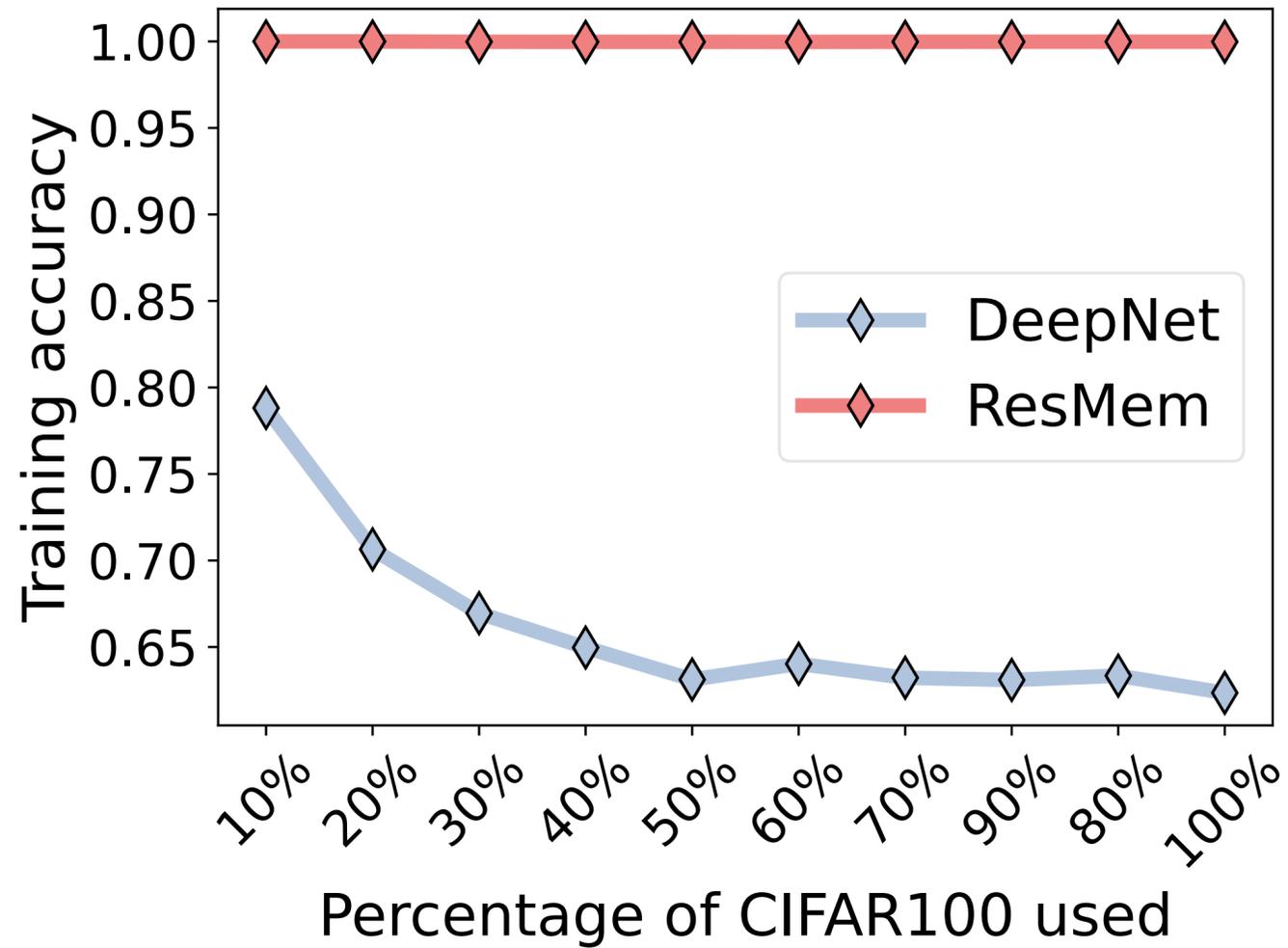
# Empirical results:

Dataset	Architecture	Test accuracy		
		DeepNet	ResMem	Gain

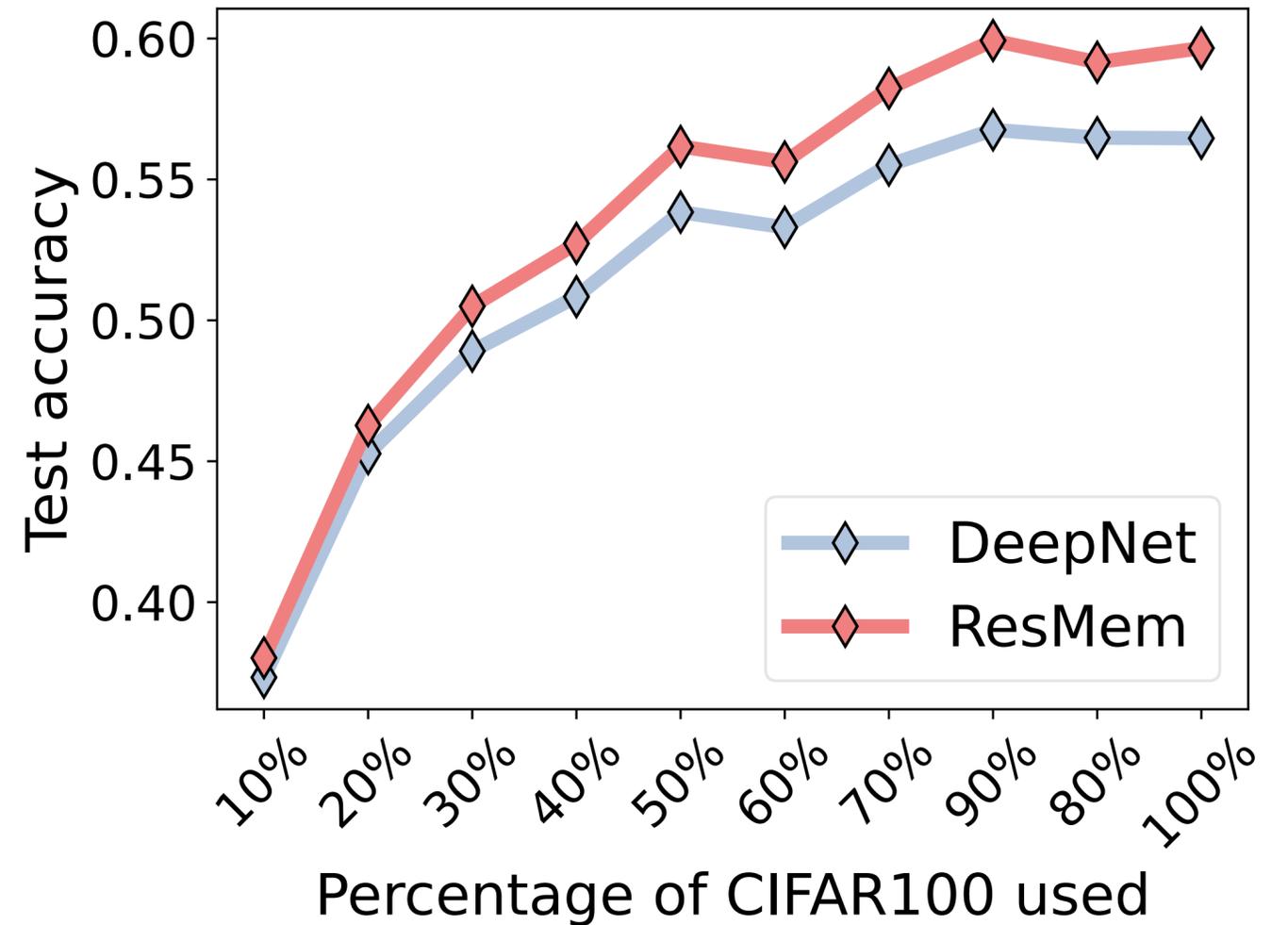
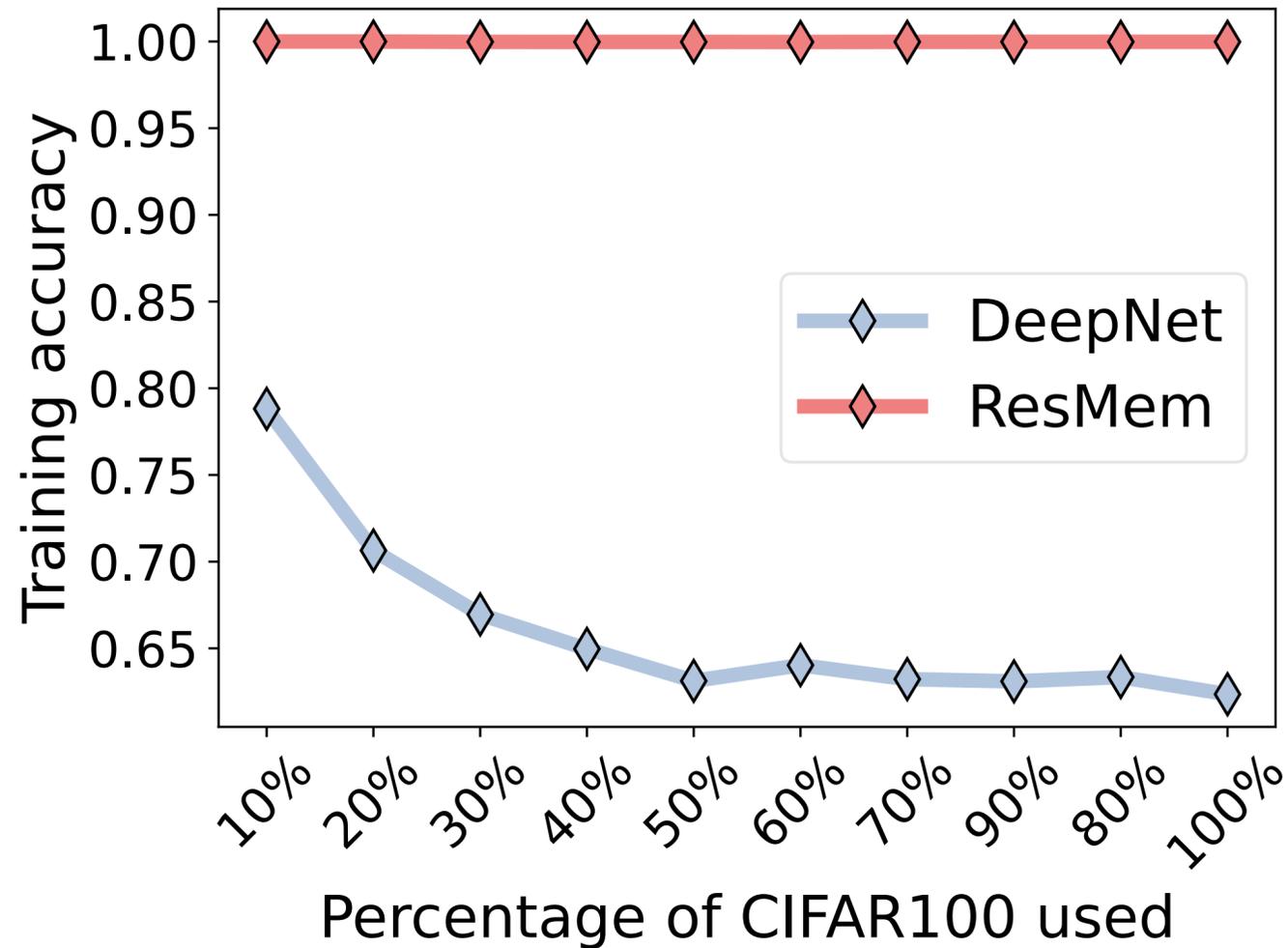
# Empirical results: simple CIFAR experiment

Dataset	Architecture	Test accuracy		
		DeepNet	ResMem	Gain
CIFAR100	ResNet-8	56.46%	<b>59.66%</b>	<b>3.20%</b>

# When does ResMem help?

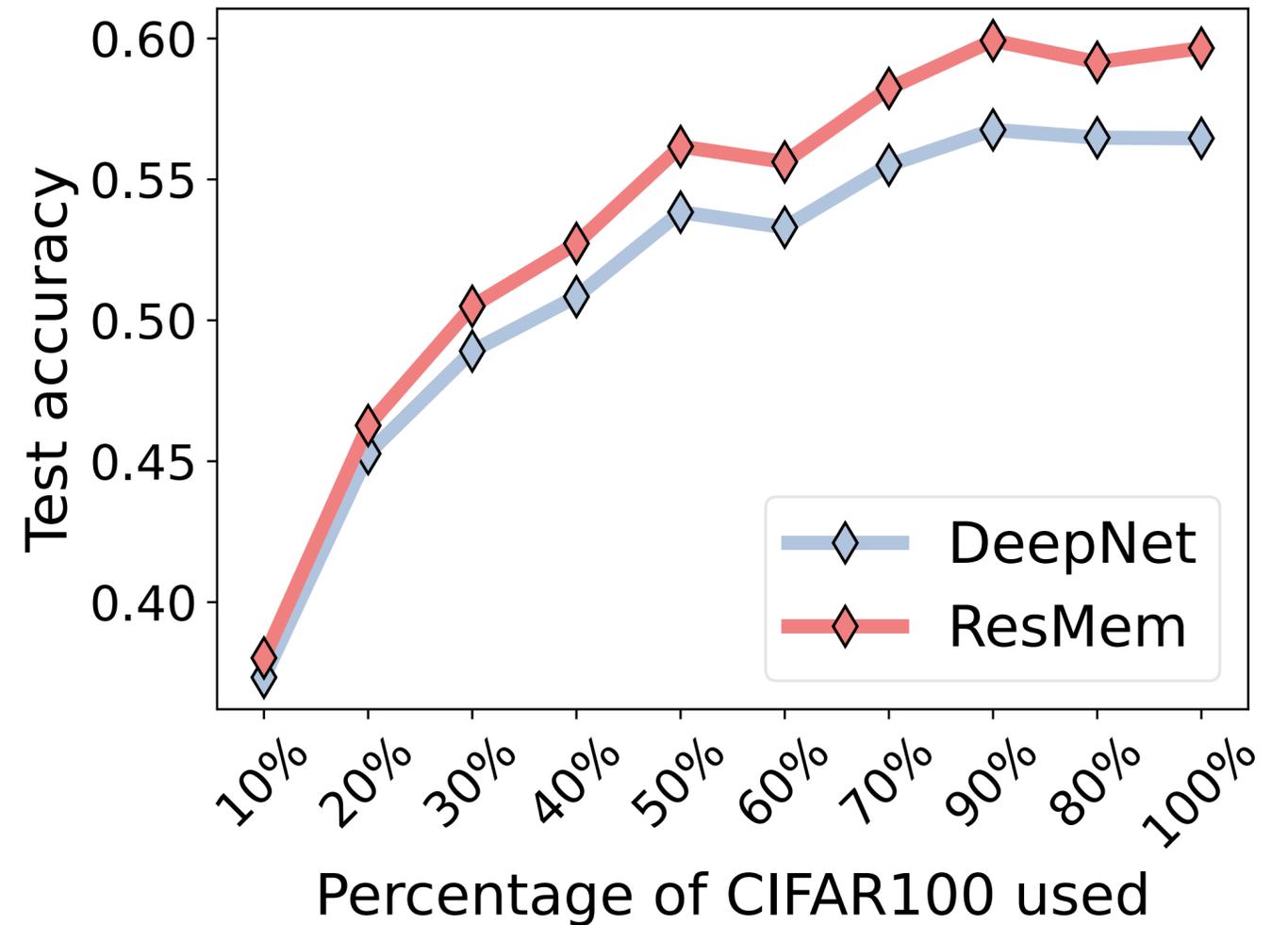
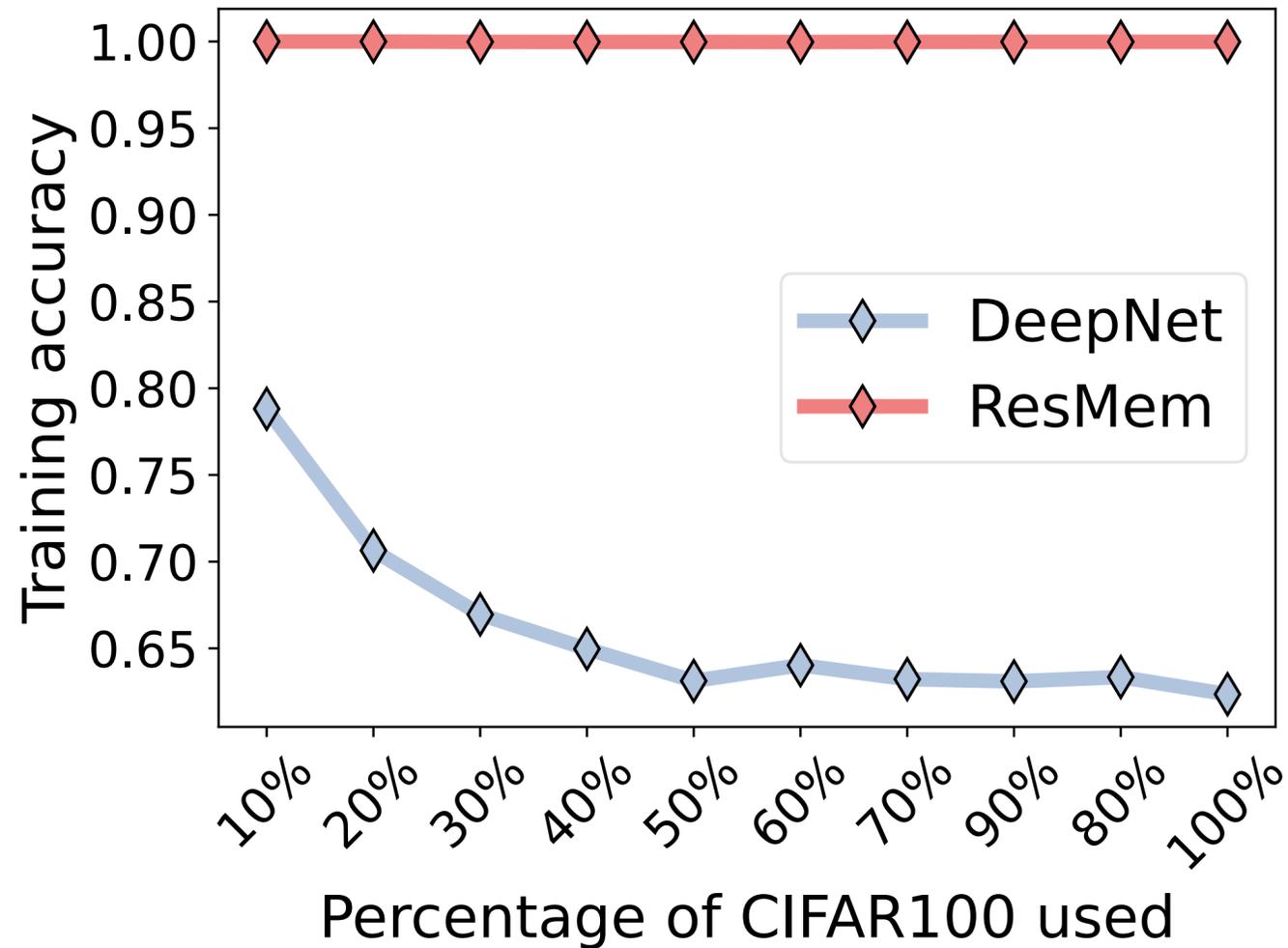


# When does ResMem help?



- When the **training sample is large**, ResMem is particularly effective.

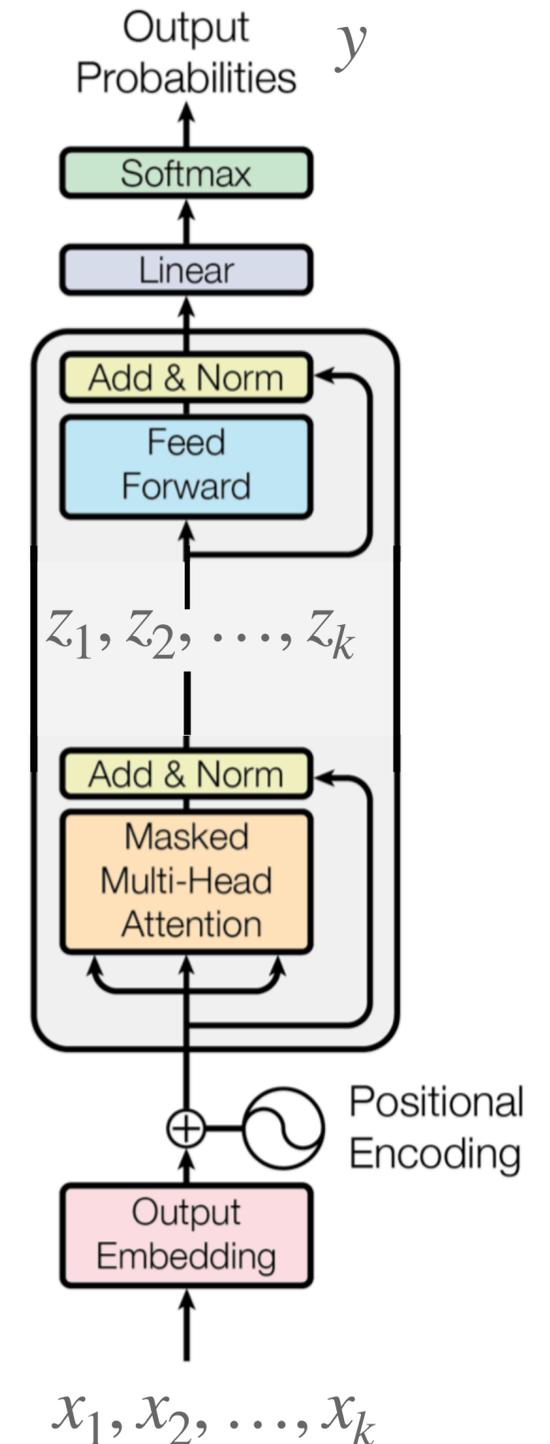
# When does ResMem help?



- When the **training sample is large**, ResMem is particularly effective.
- SOTA LLMs are typically trained for at most a single epoch. (*Google's PaLM '22*)

# ResMem on languages models

- Architecture: Decoder Only T5 (*Raffel et al. '20*).
- Dataset: C4 (text scrapped from internet).
- Residual:  $\text{onehot}(x_{k+1}) - y$ .
- Embedding: Post-attention, pre-feed forward  $z_k$ .
- Nearest neighbor: ScaNN (*Guo et al. '20*) for search over **1.6B** tokens.



# Empirical results: an overview

Dataset	Architecture	Test accuracy		
		DeepNet	ResMem	Gain
CIFAR100	ResNet-8	56.46%	<b>59.66%</b>	<b>3.20%</b>
C4	T5-Small	38.01%	<b>40.87%</b>	<b>2.86%</b>
C4	T5-Large	44.80%	<b>46.60%</b>	<b>1.80%</b>

# Where does the improved accuracy come from?

$y^{\text{ResMem}}$

$y^{\text{DeepNet}}$



rose

poppy



cup

plate



squirrel

rabbit

...allow for plenty of headroom  
inside whilst still being less than  
2.5m in **height**.

Graphic now consists of all cities  
with greater than 30,000 locals.  
Acquiring a **home** in Spain...

Filmed around 7:30-8:30 a.m. on  
Friday, **March** 9, 2012.

$y^{\text{ResMem}}$

$y^{\text{DeepNet}}$

height

length

home

residence

March

June

# Where does the improved accuracy come from?

	<u><math>y^{\text{ResMem}}</math></u>	<u><math>y^{\text{DeepNet}}</math></u>		<u><math>y^{\text{ResMem}}</math></u>	<u><math>y^{\text{DeepNet}}</math></u>
	rose	poppy	...allow for plenty of headroom inside whilst still being less than 2.5m in <b>height</b> .	height	length
	cup	plate	Graphic now consists of all cities with greater than 30,000 locals. Acquiring a <b>home</b> in Spain...	home	residence
	squirrel	rabbit	Filmed around 7:30-8:30 a.m. on Friday, <b>March</b> 9, 2012.	March	June

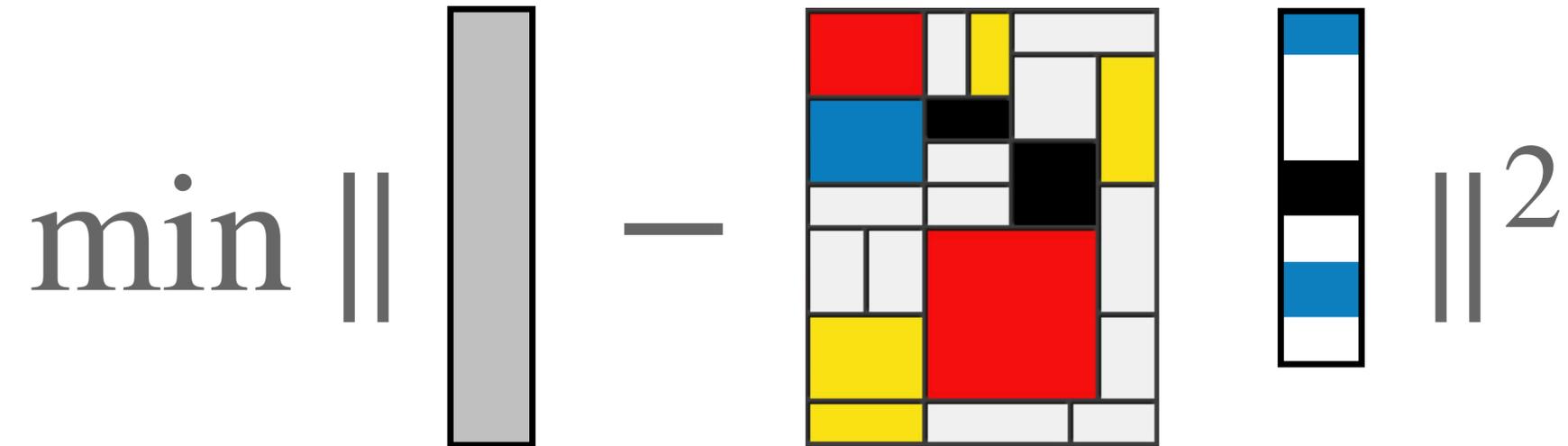
- $y^{\text{DeepNet}}$  learns some **coarse structures**. (the `flowers`)
- $y^{\text{ResMem}}$  memorizes the **fine-grained details**. (the `roses`)

# Theoretical analysis: linear regression

$$\min \| \text{gray bar} - \text{matrix} \cdot \text{vector} \|_2$$

- Why don't we memorize the **labels  $y$  directly?**

# Theoretical analysis: assumptions



# Theoretical analysis: assumptions

$$\min \| \text{gray bar} - X \text{ (matrix) } \|_2$$

$X \in \mathbb{R}^{n \times d}$

- $X \in \mathbb{R}^{n \times d}$  is the usual design matrix with row  $x_i \stackrel{i.i.d.}{\sim}$  some distribution.

# Theoretical analysis: assumptions

$$\min \| \begin{matrix} y \\ \text{gray bar} \end{matrix} - \begin{matrix} \text{matrix } X \\ \text{with colored blocks} \end{matrix} \begin{matrix} \text{blue bar} \\ \text{black bar} \\ \text{white bar} \\ \text{blue bar} \end{matrix} \|_2$$

$X \in \mathbb{R}^{n \times d}$

- $X \in \mathbb{R}^{n \times d}$  is the usual design matrix with row  $x_i \stackrel{i.i.d.}{\sim}$  some distribution.
- $y = X\theta_\star$  is generated by some true  $\theta_\star$  with  $\|\theta_\star\| = 1$

# Theoretical analysis: assumptions

$$\theta_n = \min_{\|\theta\| < L} \| y - X\theta \|^2$$

$X \in \mathbb{R}^{n \times d}$

- $X \in \mathbb{R}^{n \times d}$  is the usual design matrix with row  $x_i \stackrel{i.i.d.}{\sim}$  some distribution.
- $y = X\theta_\star$  is generated by some true  $\theta_\star$  with  $\|\theta_\star\| = 1$
- Empirical risk minimization performed over the functions class

$$\mathcal{F} = \{x \mapsto \langle x, \theta \rangle, \|\theta\| < L\}, L < 1$$

# Theoretical analysis: results

Main theoretical result [Theorem 5.3, Yang et al., '23]

stays as a irreducible error without ResMem

$$\text{Test Risk of ResMem} \lesssim \underbrace{d^2 L^2 n^{-2/3}}_{\text{parametric rate from linear regression}} + \overbrace{d^2 (1-L)^2 \left[ \frac{\log(n^{1/d})}{n} \right]^{1/d}}^{\text{stays as a irreducible error without ResMem}} \longrightarrow 0 \text{ as } n \rightarrow \infty$$

# Theoretical analysis: results

Main theoretical result [Theorem 5.3, Yang et al., '23]

stays as a irreducible error without ResMem

$$\text{Test Risk of ResMem} \lesssim \underbrace{d^2 L^2 n^{-2/3}}_{\text{parametric rate from linear regression}} + \overbrace{d^2 (1-L)^2 \left[ \frac{\log(n^{1/d})}{n} \right]^{1/d}}^{\text{stays as a irreducible error without ResMem}} \longrightarrow 0 \text{ as } n \rightarrow \infty$$

- The nearest neighbors component **expands the capacity** of the linear regressor by adding a non parametric component.
- If we memorize the labels  $y$  directly, we end up with a **slow rate of  $n^{-1/d}$** .

*Learn **the flowers**, remember **the roses**.*

