

# Learning from Oblivion: Predicting Knowledge-Overflowed Weights via Retrodiction of Forgetting

Jinhyeok Jang, Jaehong Kim, Jung Uk Kim



*Highlight*



# Motivation

- Our work is motivated by the following fundamental questions:

- 
- (i) What defines a better pre-trained weight?*
  - (ii) If practical constraints exist, how can we obtain the better weights?*
- 

- For the first question, we draw inspiration from "Scaling Law."
  - More training data generally leads to better model weights.
  - However, obtaining additional training data is challenging.

# Hypothesis

- Our Hypothesis

*Weights trained as if on more training data can be predicted by reversing the forgetting trajectory.*

## Preliminary

Consider a dataset  $D^0$  and a subset  $D^1 \subset D^0$ . Let  $\Theta^0$  be the weights obtained from pre-training on  $D^0$ . When fine-tuning  $\Theta^0$  on  $D^1$ , the resulting weights  $\Theta^1$  exhibit a loss of knowledge about  $D^0 \setminus D^1$  due to forgetting.

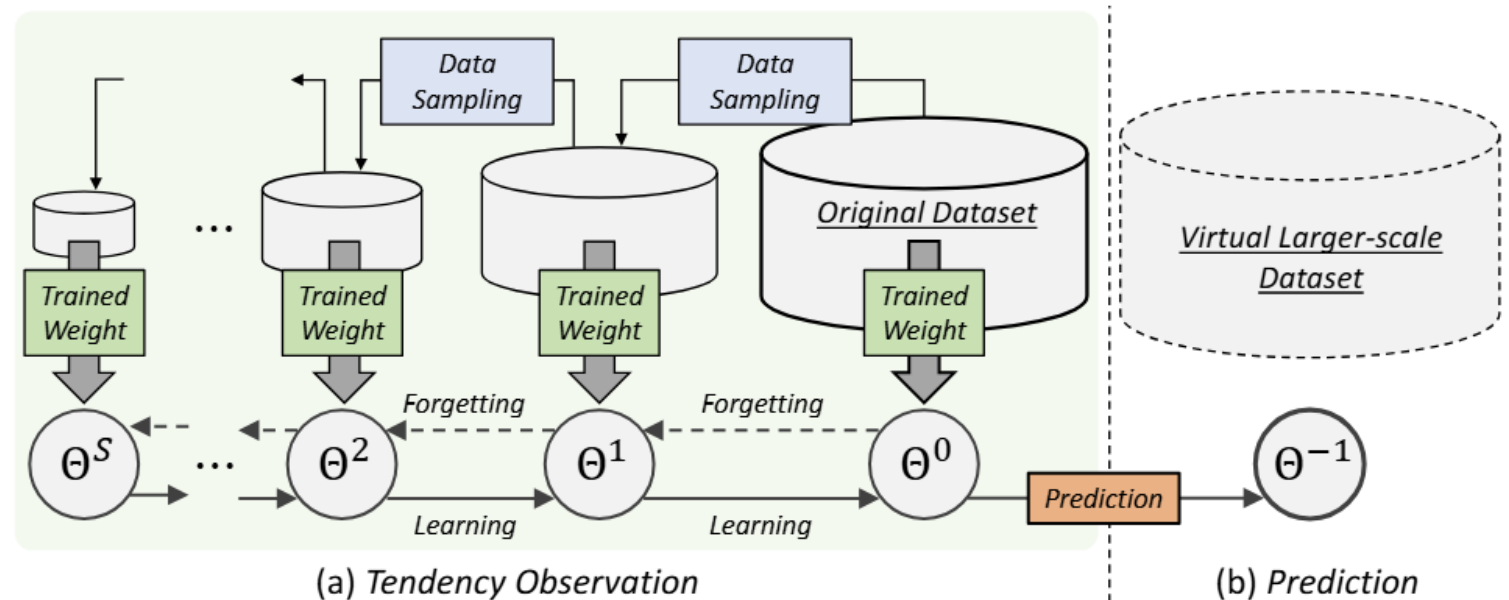
## Hypothesis

Given  $\Theta^0$  trained on  $D^0$ , there exist ideal weights  $\Theta^{-1}$  satisfying the following conditions:

1.  $\Theta^{-1}$  is trained on  $D^{-1}$ , where  $D^0 \subset D^{-1}$ .
2. Fine-tuning  $\Theta^{-1}$  on  $D^0$  results in  $\Theta^0$ .

## Extension to Progressive Forgetting

Given a sequence of datasets  $D^S \subset D^{S-1} \subset \dots \subset D^1 \subset D^0$  constructed with a sampling ratio  $r \in [0, 1]$ , we obtain a corresponding sequence of weights  $[\Theta^S, \Theta^{S-1}, \dots, \Theta^0]$ , where each  $\Theta^i$  is obtained by fine-tuning  $\Theta^{i-1}$  on  $D^i$ .



# Observation

- **Visualizing the Forgetting Trajectory as a Loss Landscape**

- **Dataset:** CIFAR10

- **Method:**

- 1) **Base trajectory collection**

Fine-tune the model on progressively downsampled datasets.

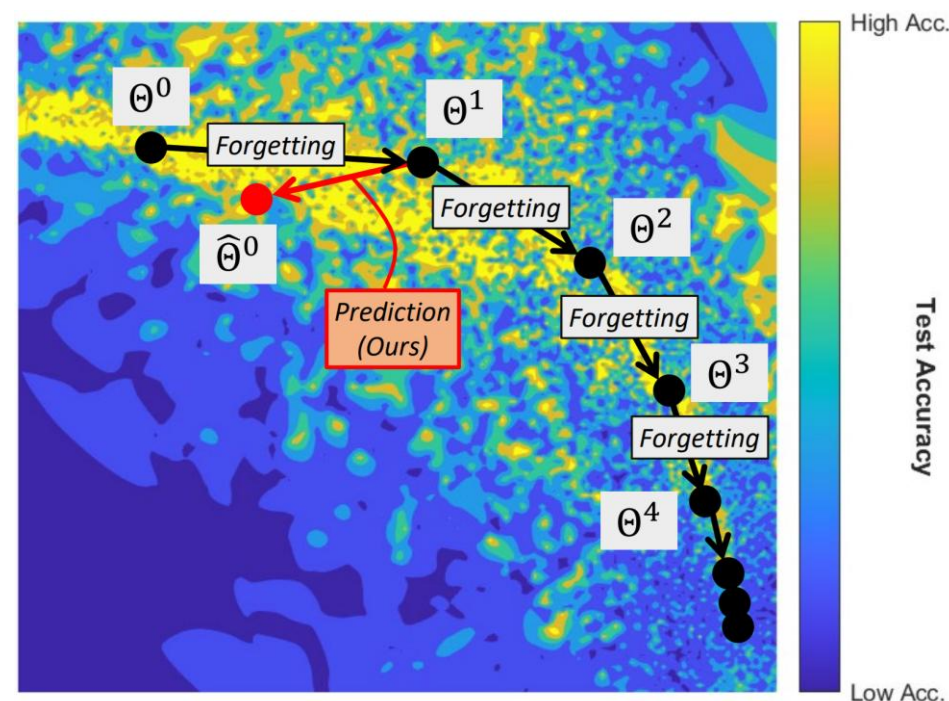
- 2) **Neighborhood exploration**

Perform exhaustive fine-tuning with random noise injection.

- 3) **Trajectory visualization**

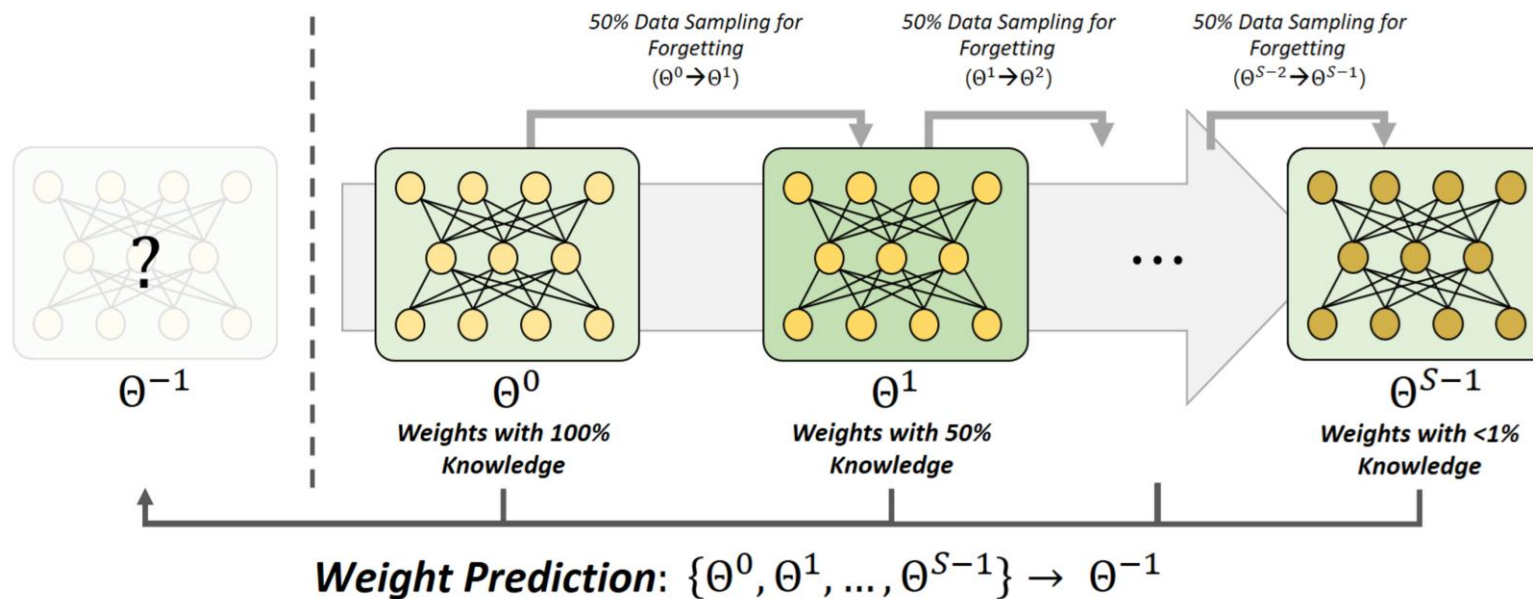
Project the collected weight trajectories onto a PCA space.

- **The forgetting trajectory exhibits a pattern in the weight space.**



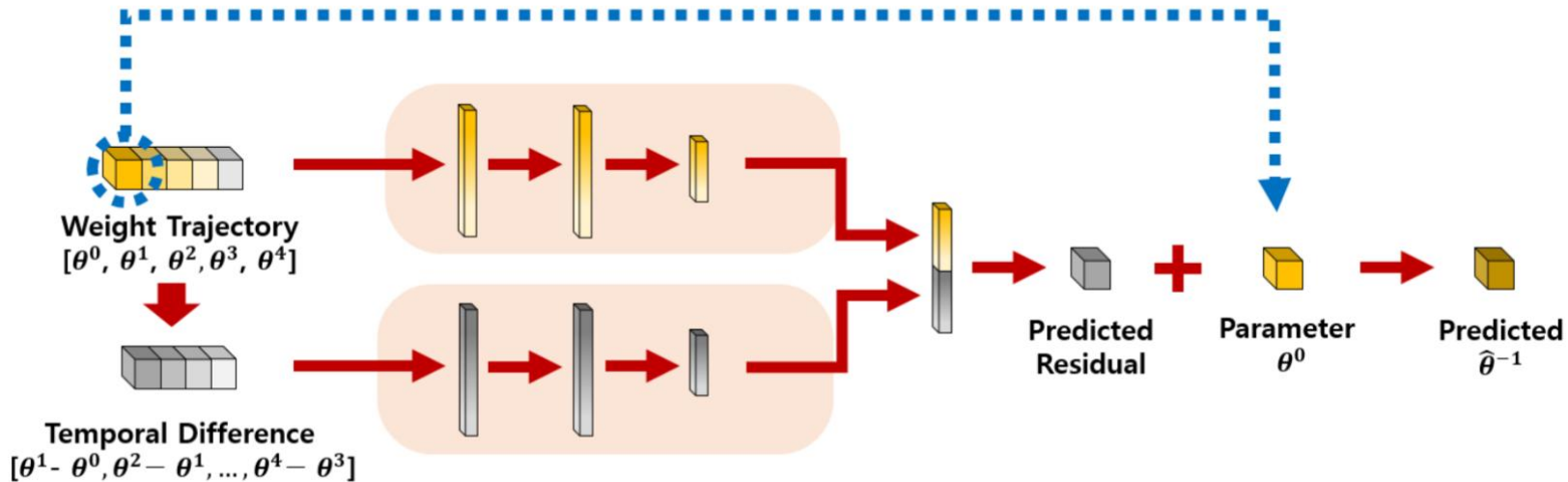
# Knowledge-Overflowed Weight Prediction

- We establish three steps:
  - 1) Pretraining on the full dataset
  - 2) Progressive fine-tuning on sequential subsets (*intentional forgetting*)
  - 3) Predicting weights trained as if on more data by reversing the forgetting trajectory and extrapolation



# KNOWN Nowcaster (KNOWN)

- We meta-train a dedicated hypernetwork, called KNOWN.
  - 1) **Trajectory Collection:** Collect weight trajectories across diverse settings (e.g., architectures, datasets, and training recipes)
  - 2) **KNOWN Meta-Training:** Train KNOWN as a lightweight meta-learned predictor
  - 3) After meta-training, KNOWN requires no additional data collection or meta-learning.



# Experiments

- **KNOW prediction is effective across diverse settings.**
  - Weight-space prediction methods (e.g., curve fitting and task arithmetic) can improve performance.
  - Our meta-learned **KNOWN** provides the largest performance gain.

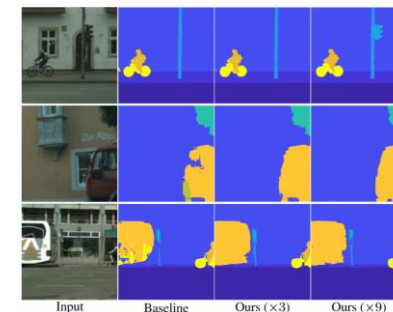
Pred. ( $\times n$ )	Methods	Amount of pre-training Data (%)					
		100%	50%	25%	12.5%	6.25%	
$\times 1$	Naïve Transfer (Baseline)	92.40 $\pm$ 0.11	92.08 $\pm$ 0.18	91.90 $\pm$ 0.24	91.49 $\pm$ 0.22	91.51 $\pm$ 0.12	
	Incremental Pretraining (Baseline)	92.29 $\pm$ 0.10	91.32 $\pm$ 0.14	90.79 $\pm$ 0.48	90.07 $\pm$ 0.10	89.51 $\pm$ 0.08	
$\times 2$	KNOW <sup>s</sup>	LinearFit	92.70 $\pm$ 0.16	92.28 $\pm$ 0.16	91.89 $\pm$ 0.24	91.42 $\pm$ 0.15	91.37 $\pm$ 0.16
		LogFit	92.83 $\pm$ 0.25	92.41 $\pm$ 0.07	91.89 $\pm$ 0.28	91.39 $\pm$ 0.17	91.33 $\pm$ 0.24
		ExpFit	79.58 $\pm$ 0.41	86.79 $\pm$ 0.23	89.26 $\pm$ 0.04	89.85 $\pm$ 0.08	90.60 $\pm$ 0.16
		TaskVector [24]	92.69 $\pm$ 0.09	92.39 $\pm$ 0.12	<u>92.22 <math>\pm</math> 0.25</u>	91.45 $\pm$ 0.05	91.46 $\pm$ 0.08
		ConsensusTA <sup>†</sup> [55]	92.69 $\pm$ 0.09	92.39 $\pm$ 0.12	92.22 $\pm$ 0.25	91.45 $\pm$ 0.05	91.46 $\pm$ 0.08
		MagMax <sup>‡</sup> [44]	92.69 $\pm$ 0.09	92.39 $\pm$ 0.12	92.22 $\pm$ 0.25	91.45 $\pm$ 0.05	91.46 $\pm$ 0.08
		TSV [11]	92.82 $\pm$ 0.14	92.36 $\pm$ 0.16	92.19 $\pm$ 0.11	91.64 $\pm$ 0.19	91.74 $\pm$ 0.25
		<b>KNOWN</b>	<b>93.00 <math>\pm</math> 0.11</b>	<b>92.58 <math>\pm</math> 0.14</b>	<b>92.29 <math>\pm</math> 0.04</b>	<b>92.11 <math>\pm</math> 0.10</b>	<b>91.90 <math>\pm</math> 0.13</b>
$\times 4$	KNOW <sup>s</sup>	LinearFit	92.40 $\pm$ 0.08	91.99 $\pm$ 0.22	91.63 $\pm$ 0.14	90.52 $\pm$ 0.12	90.64 $\pm$ 0.09
		LogFit	92.82 $\pm$ 0.10	92.33 $\pm$ 0.13	91.87 $\pm$ 0.13	91.05 $\pm$ 0.15	91.08 $\pm$ 0.23
		ExpFit	72.70 $\pm$ 0.64	80.75 $\pm$ 0.55	85.16 $\pm$ 0.04	87.46 $\pm$ 0.27	88.80 $\pm$ 0.29
		TaskVector [24]	92.70 $\pm$ 0.09	92.48 $\pm$ 0.10	92.19 $\pm$ 0.12	91.45 $\pm$ 0.14	91.36 $\pm$ 0.18
		ConsensusTA <sup>‡</sup> [55]	92.70 $\pm$ 0.09	92.48 $\pm$ 0.10	92.19 $\pm$ 0.12	91.45 $\pm$ 0.14	91.36 $\pm$ 0.18
		MagMax [44]	92.44 $\pm$ 0.28	92.19 $\pm$ 0.22	92.12 $\pm$ 0.19	91.29 $\pm$ 0.16	91.15 $\pm$ 0.14
		TSV [11]	92.74 $\pm$ 0.13	92.40 $\pm$ 0.16	92.17 $\pm$ 0.18	91.64 $\pm$ 0.20	91.63 $\pm$ 0.16
		<b>KNOWN</b>	<b>93.27 <math>\pm</math> 0.09</b>	<b>92.62 <math>\pm</math> 0.25</b>	<b>92.88 <math>\pm</math> 0.11</b>	<b>92.40 <math>\pm</math> 0.06</b>	<b>91.98 <math>\pm</math> 0.14</b>
$\times 8$	KNOW <sup>s</sup>	LinearFit	92.13 $\pm$ 0.10	91.71 $\pm$ 0.22	90.77 $\pm$ 0.18	90.15 $\pm$ 0.40	90.42 $\pm$ 0.44
		LogFit	92.65 $\pm$ 0.10	91.93 $\pm$ 0.25	91.59 $\pm$ 0.05	90.49 $\pm$ 0.25	90.68 $\pm$ 0.24
		ExpFit	66.54 $\pm$ 0.29	73.89 $\pm$ 0.91	81.05 $\pm$ 0.14	84.34 $\pm$ 0.56	85.16 $\pm$ 0.45
		TaskVector [24]	92.65 $\pm$ 0.13	92.24 $\pm$ 0.15	92.03 $\pm$ 0.15	91.13 $\pm$ 0.21	91.28 $\pm$ 0.34
		ConsensusTA [55]	92.43 $\pm$ 0.14	92.02 $\pm$ 0.12	91.93 $\pm$ 0.17	91.11 $\pm$ 0.29	90.63 $\pm$ 0.32
		MagMax [44]	92.39 $\pm$ 0.20	92.16 $\pm$ 0.10	92.07 $\pm$ 0.09	91.43 $\pm$ 0.20	90.85 $\pm$ 0.33
		TSV [11]	92.72 $\pm$ 0.11	92.37 $\pm$ 0.15	92.25 $\pm$ 0.15	91.51 $\pm$ 0.17	91.48 $\pm$ 0.15
		<b>KNOWN</b>	<b>93.55 <math>\pm</math> 0.05</b>	<b>93.11 <math>\pm</math> 0.19</b>	<b>92.92 <math>\pm</math> 0.15</b>	<b>92.22 <math>\pm</math> 0.37</b>	<b>92.07 <math>\pm</math> 0.18</b>

<CIFAR100→CIFAR10>

Pred. ( $\times n$ )	Methods	Downstream Dataset $\mathcal{D}$ (ImageNet (Pre-train) $\rightarrow \mathcal{D}$ )					
		CIFAR100	TinyImageNet	Car	CUB	Flowers	
$\times 1$	Naïve Transfer (Baseline)	82.03 $\pm$ 0.25	76.17 $\pm$ 0.33	88.12 $\pm$ 0.35	70.49 $\pm$ 0.58	87.98 $\pm$ 0.38	
$\times 3$	KNOW	LogFit	81.73 $\pm$ 0.37 (-0.30)	77.32 $\pm$ 0.14 (+1.15)	88.49 $\pm$ 0.68 (+0.37)	70.45 $\pm$ 0.71 (-0.04)	88.24 $\pm$ 0.20 (+0.26)
		TaskVector [24]	82.15 $\pm$ 0.29 (+0.12)	77.49 $\pm$ 0.21 (+1.32)	88.31 $\pm$ 0.24 (+0.19)	71.00 $\pm$ 0.25 (+0.51)	88.18 $\pm$ 0.42 (+0.20)
		TSV [11]	82.14 $\pm$ 0.46 (+0.11)	76.06 $\pm$ 0.17 (-0.11)	<b>88.78 <math>\pm</math> 0.12</b> (+0.80)	70.84 $\pm$ 0.48 (+0.35)	86.37 $\pm$ 0.37 (-1.61)
$\times 9$	KNOW	<b>KNOWN</b>	<b>82.46 <math>\pm</math> 0.20</b> (+0.43)	<b>77.53 <math>\pm</math> 0.11</b> (+1.36)	88.57 $\pm$ 0.18 (+0.45)	<b>71.18 <math>\pm</math> 0.45</b> (+0.69)	<b>88.65 <math>\pm</math> 0.69</b> (+0.67)
		LogFit	81.75 $\pm$ 0.35 (-0.28)	77.29 $\pm$ 0.14 (+1.12)	88.42 $\pm$ 0.30 (+0.30)	71.09 $\pm$ 0.65 (+0.60)	88.43 $\pm$ 0.57 (+0.45)
		TaskVector [24]	82.12 $\pm$ 0.45 (+0.09)	77.29 $\pm$ 0.23 (+1.12)	88.18 $\pm$ 0.17 (+0.06)	70.38 $\pm$ 0.34 (-0.11)	88.37 $\pm$ 0.43 (+0.39)
		MagMax [44]	82.27 $\pm$ 0.24 (+0.24)	75.98 $\pm$ 0.18 (-0.19)	88.61 $\pm$ 0.31 (+0.49)	70.95 $\pm$ 0.39 (+0.46)	86.51 $\pm$ 0.76 (-1.47)
$\times 9$	KNOW	TSV [11]	82.31 $\pm$ 0.20 (+0.28)	76.18 $\pm$ 0.27 (+0.01)	<b>88.85 <math>\pm</math> 0.09</b> (+0.73)	70.79 $\pm$ 0.32 (+0.30)	86.47 $\pm$ 0.56 (-1.51)
		<b>KNOWN</b>	<b>82.33 <math>\pm</math> 0.15</b> (+0.30)	<b>77.58 <math>\pm</math> 0.12</b> (+1.41)	<b>88.85 <math>\pm</math> 0.23</b> (+0.73)	<b>71.30 <math>\pm</math> 0.28</b> (+0.81)	<b>88.53 <math>\pm</math> 0.27</b> (+0.55)

<ImageNet→CIFAR100/TinyImageNet/Cars196/CUB/Flowers>

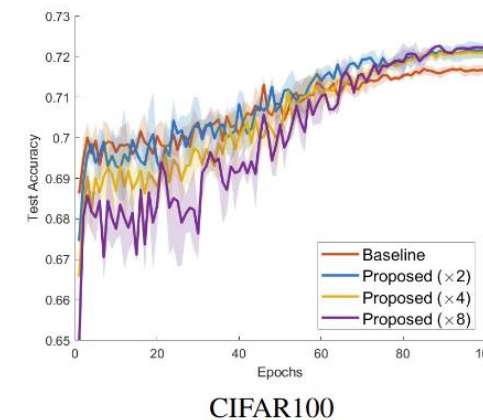
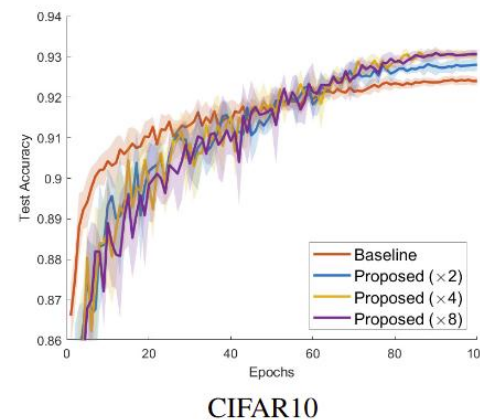
Pred. ( $\times n$ )	Methods	mIoU (%)	
$\times 1$	Naïve Transfer (Baseline)	68.52 $\pm$ 1.34	
$\times 3$	KNOW	LogFit	68.79 $\pm$ 0.68 (+0.27)
		TaskVector [24]	68.83 $\pm$ 1.71 (+0.31)
		TSV [11]	68.64 $\pm$ 0.71 (+0.12)
$\times 9$	KNOW	<b>KNOWN</b>	<b>69.00 <math>\pm</math> 1.04</b> (+0.48)
		LogFit	69.16 $\pm$ 1.95 (+0.64)
		TaskVector [24]	67.99 $\pm$ 0.96 (-0.53)
		MagMax [44]	70.04 $\pm$ 0.81 (+1.52)
$\times 9$	KNOW	TSV [11]	68.98 $\pm$ 1.06 (+0.46)
		<b>KNOWN</b>	<b>71.22 <math>\pm</math> 0.82</b> (+2.70)



<ImageNet→CityScapes>

# Ablation

- KNOWN progresses slowly in the early phase due to prediction error, but reaches higher final performance.
- Longer history improves convergence.
  - However, longer history requires more training data.
- KNOWN achieves a good balance between data efficiency and training efficiency.



	<b>Pred.(×2)</b>	<b>Pred.(×4)</b>	<b>Pred.(×8)</b>
$S = 2$ (Task Vector)	92.69 ± 0.09	92.70 ± 0.09	92.65 ± 0.13
$S = 3$ (KNOWN)	<b>93.01 ± 0.16</b>	93.04 ± 0.06	92.72 ± 0.10
$S = 4$ (KNOWN)	92.97 ± 0.16	93.10 ± 0.13	92.89 ± 0.16
$S = 5$ (KNOWN)	93.00 ± 0.11	<b>93.27 ± 0.09</b>	<b>93.55 ± 0.05</b>

Methods	Dataset	Training Inference		Accuracy (%)
	Amount (%)	Time (s)	Time (s)	
Naïve Transfer (Baseline)	50	3,800	N/A	92.08
	100	7,600	N/A	92.40
KNOWN (×2)	50	7,372	3.01	92.58
KNOWN (×4)	50	7,372	6.02	92.62
KNOWN (×8)	50	7,372	9.03	<b>93.11</b>

# Summary

- We introduce **KNOW prediction**, which predicts weights trained as if on more data by reversing forgetting trajectories.
- We observe that progressive forgetting exhibits a pattern in weight space.
- We propose **KNOWN**, a lightweight hypernetwork for efficient weight prediction across diverse settings.
- Our findings show that forgetting dynamics can serve as a useful signal for data- and training-efficient model improvement.

# Thanks

