

LoTUS: Large-Scale Machine Unlearning with a Taste of Uncertainty

Christoforos N. Spartalis Theodoros Semertzidis Efstratios Gavves Petros Daras







UNIVERSITEIT VAN AMSTERDAM



Machine Unlearning Definition & Applications



Challenges & Evaluation





Challenges & Evaluation



Challenges & Evaluation



Runtime of Unlearning Algorithm & Computational Complexity Analysis

Contributions

Large-Scale Unlearning Benchmark Simulates Real-World Conditions (ImageNet1k & Limited Data Access)

Retrain-Free Jensen-Shannon Divergence (RF-JSD) Metric: Ensures Evaluation without the Gold Standard Model

✓ LoTUS

Increases Entropy In Predictions up to an Information-Theoretical Bound JSD (Unlearned & Gold Standard)

$$\mathcal{JS}(f_{un}(D_f) \mid\mid f_{gold}(D_f))$$

	Dataset $\left(\frac{\text{num. of forget samples}}{\text{num. of training samples}} \times 100\%\right)$	PCC (†)	p-value (\downarrow)
	CIFAR-100 (10%)	0.84	0.0043
	CIFAR-10 (10%)	0.92	0.0005
ΥïΤ	MUFAC	0.93	0.0003
-	CIFAR-100 (50%)	0.94	0.0001
	CIFAR-10 (50%)	0.99	0.0000
	CIFAR-100 (10%)	0.97	0.0000
t18	CIFAR-10 (10%)	0.90	0.0011
Nei	MUFAC	0.88	0.0018
Ses	CIFAR-100 (50%)	0.91	0.0006
-	CIFAR-10 (50%)	0.89	0.0013
	Mean \pm Std	$0.92_{\pm 0.04}$	$0.0010_{\pm 0.0010}$

Table 7.Retrain Free-JSD (RF-JSD) and JSD Correlationmeasured with the Pearson correlation coefficient (PCC).

RF-JSD (Unlearned & Pre-trained)

$$\begin{aligned} \mathcal{JS}\!\left(f_{\text{un}}(D_f) \,||\, f_{\text{orig}}(D_u)\right) &= \frac{1}{k} \sum_{c=1}^k \mathcal{JS}(P_i \;||\; Q_i) \\ P_i \!=\! \frac{1}{Z_P} \sum_{j=1}^{n_i} f_{\text{un}}(x_j \,|\, y_j \!=\! i) \,, \; Q_i \!=\! \frac{1}{Z_Q} \sum_{j=1}^{n_i} f_{\text{orig}}(x_j \,|\, y_j \!=\! i) \end{aligned}$$

Method	RF -JSD $\times 1e4 (\downarrow)$	Time (\downarrow)	Retain Acc.	MIA Acc.
Original	$1.22_{\pm 0.01}$	(pre-trained)	$0.94_{\pm 0.00}$	$0.71_{\pm 0.00}$
Finetuning	$2.22_{\pm 0.02}$	$16.24_{\pm 0.03}$	$0.97_{\pm 0.00}$	$0.78_{\pm 0.00}$
NegGrad+	$2.17_{\pm 0.02}$	18.10 ± 0.03	$0.97_{\pm 0.00}$	$0.80_{\pm 0.00}$
Rnd Labeling	$1.80_{\pm 0.09}$	$19.37_{\pm 0.03}$	$0.95_{\pm 0.01}$	$0.74_{\pm 0.01}$
Bad Teacher	$3.16_{\pm 3.25}$	11.66 ± 0.03	$0.77_{\pm 0.21}$	$0.52_{\pm 0.18}$
SCRUB	$1.24_{\pm 0.01}$	24.49 ± 0.03	$0.94_{\pm 0.00}$	$0.71_{\pm 0.00}$
SSD	$1.23_{\pm 0.01}$	22.61 ± 0.10	$0.94_{\pm 0.00}$	$0.71_{\pm 0.00}$
UNSIR	$2.54_{\pm 0.03}$	$33.12_{\pm 0.03}$	$0.99_{\pm 0.00}$	$0.77_{\pm 0.01}$
SalUn	$1.83_{\pm 0.03}$	$59.27_{\pm 0.37}$	$0.95_{\pm 0.00}$	$0.74_{\pm 0.01}$
LoTUS	$1.11_{\pm 0.01}$	$10.72_{\pm0.01}$	$0.94_{\pm 0.00}$	$0.61_{\pm 0.01}$

Table 4. Large-Scale Unlearning with ImageNet1k: LoTUS outperforms state-of-the-art approaches in both unlearning effectiveness (RF-JSD) and efficiency (Time). While other metrics lack concrete validation due to the absence of a Gold Standard, they provide additional insights: LoTUS uniquely preserves the Retain Accuracy of the pre-trained model while reducing MIA Accuracy.

Results

	Metric (\downarrow)	Gold Std	Finetuning	NegGrad+ [23]	RndLbl [17]	BadT [9]	SCRUB [23]	SSD [14]	UNSIR [34]	SalUn [12]	LoTUS
er (ViT) TinyIN	Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c} 0.0000 \\ 0.00_{\pm 0.00} \\ 228.9_{\pm 6.49} \end{array}$	$\begin{array}{c} 0.0175 \\ 0.05_{\pm 0.00} \\ 22.64_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0400 \\ 0.10_{\pm 0.00} \\ 25.20_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.2925 \\ 0.64_{\pm 1.03} \\ 25.19_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0775 \\ 0.18_{\pm 0.01} \\ 16.91_{\pm 0.05} \end{array}$	$\begin{array}{c} 0.0225 \\ 0.04_{\pm 0.00} \\ 33.25_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0225 \\ 0.04_{\pm 0.00} \\ 27.27_{\pm 0.06} \end{array}$	$\begin{array}{c} 0.0225 \\ 0.06_{\pm 0.00} \\ 21.17_{\pm 0.22} \end{array}$	$\begin{array}{c} 0.0925 \\ 0.25_{\pm 0.59} \\ 76.97_{\pm 1.72} \end{array}$	$\begin{array}{c} 0.0150 \\ 0.03_{\pm 0.00} \\ 13.41_{\pm 0.04} \end{array}$
Transforn C-100	Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c} 0.0000 \\ 0.00 {\pm} 0.00 \\ 112.25 {\pm} 0.13 \end{array}$	$\begin{array}{c} 0.0275 \\ 0.07_{\pm 0.00} \\ 11.35_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.0325 \\ 0.13_{\pm 0.01} \\ 12.63_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0175 \\ 0.06_{\pm 0.00} \\ 12.79_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0375 \\ 0.17_{\pm 0.01} \\ 9.18_{\pm 0.27} \end{array}$	$\begin{array}{c} 0.0200 \\ \textbf{0.04}_{\pm \textbf{0.00}} \\ 16.74_{\pm 0.03} \end{array}$	0.0175 0.04_{±0.00} 13.67 _{±0.02}	$\begin{array}{c} 0.0250 \\ 0.08_{\pm 0.01} \\ 10.69_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0200 \\ 0.06_{\pm 0.01} \\ 38.15_{\pm 0.04} \end{array}$	$\begin{array}{c} 0.0125\\ 0.04_{\pm 0.02}\\ 7.02_{\pm 0.01}\end{array}$
Vision MUFAC	Avg Gap JSD×1e-4 Time (min.)	$\begin{array}{c} 0.0000 \\ 0.00_{\pm 0.00} \\ 13.83_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0400 \\ 0.27_{\pm 0.02} \\ 1.40_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0475 \\ 0.39_{\pm 0.04} \\ 1.67_{\pm 0.00} \end{array}$	$\begin{array}{c} \textbf{0.0200} \\ 0.35_{\pm 0.09} \\ 1.76_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.1750 \\ 1.89_{\pm 1.01} \\ 2.09_{\pm 0.25} \end{array}$	$\begin{array}{c} \textbf{0.0200} \\ \textbf{0.05}_{\pm \textbf{0.02}} \\ 2.21_{\pm 0.01} \end{array}$	$\begin{array}{c} \textbf{0.0200} \\ 0.17_{\pm 0.17} \\ 1.91_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.0475 \\ 0.85_{\pm 0.06} \\ 3.21_{\pm 0.01} \end{array}$	$\begin{array}{c} \textbf{0.0200} \\ 0.29_{\pm 0.01} \\ 8.15_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.0200 \\ 0.05_{\pm 0.01} \\ 1.09_{\pm 0.00} \end{array}$
N18) TinyIN	Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c c} 0.0000 \\ 0.00_{\pm 0.00} \\ 46.81_{\pm 0.57} \end{array}$	$0.2200 \\ 1.80_{\pm 0.04} \\ 2.85_{\pm 0.00}$	$\begin{array}{c} 0.2250 \\ 1.82_{\pm 0.07} \\ 3.17_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1925 \\ 1.71_{\pm 0.11} \\ 3.44_{\pm 0.01} \end{array}$	$0.2850 \\ 1.81_{\pm 0.04} \\ 1.91_{\pm 0.01}$	$\begin{array}{c} 0.2725 \\ 0.98_{\pm 0.00} \\ 3.97_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.2700 \\ 0.96_{\pm 0.02} \\ 3.47_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.2375 \\ 1.76_{\pm 0.05} \\ 5.00_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.2025 \\ 1.93_{\pm 0.09} \\ 6.06_{\pm 0.06} \end{array}$	$\begin{array}{c} 0.1675 \\ 0.62_{\pm 0.01} \\ 1.62_{\pm 0.00} \end{array}$
SNet18 (RN C-100	Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c c} 0.0000\\ 0.00_{\pm 0.00}\\ 3.39_{\pm 0.30}\end{array}$	$\begin{array}{c} 0.3600 \\ 6.88_{\pm 0.59} \\ 0.43_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.3575 \\ 6.87_{\pm 0.62} \\ 0.49_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.4025 \\ 5.84_{\pm 0.98} \\ 0.57_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.3675 \\ 4.30_{\pm 0.49} \\ 0.34_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.1650 \\ 1.87_{\pm 0.08} \\ 0.58_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.2125 \\ 3.04 {\pm} 1.55 \\ 0.54 {\pm} 0.00 \end{array}$	$\begin{array}{c} 0.3625 \\ 3.05_{\pm 0.32} \\ 0.45_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.3650 \\ 6.50_{\pm 0.60} \\ 1.55_{\pm 0.01} \end{array}$	$\begin{array}{c} \textbf{0.1200} \\ \textbf{1.67}_{\pm 0.37} \\ \textbf{0.30}_{\pm 0.01} \end{array}$
Re MUFAC	Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c} 0.0000 \\ 0.00_{\pm 0.00} \\ 7.34_{\pm 0.77} \end{array}$	$\begin{array}{c} 0.1525 \\ 19.52_{\pm 6.23} \\ 0.76_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1550 \\ 19.16_{\pm 5.31} \\ 0.91_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1300 \\ 9.51_{\pm 2.39} \\ 1.06_{\pm 0.00} \end{array}$	$\begin{array}{c} \textbf{0.1025} \\ 9.41_{\pm 0.04} \\ 0.66_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1625 \\ 10.53_{\pm 2.31} \\ 1.20_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1600 \\ 10.30_{\pm 2.28} \\ 1.07_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1450 \\ 16.32_{\pm 4.82} \\ 1.68_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.1400 \\ 15.25_{\pm 5.20} \\ 2.72_{\pm 0.02} \end{array}$	0.1250 6.90 _{±1.49} 0.62 _{±0.00}

Table 1. Performance Summary of unlearning 10% of Tiny-ImageNet (TinyIN), CIFAR-100 (C-100), and MUFAC training sets: LoTUS outperforms state-of-the-art approaches in balancing forgetting and retention (measured by Avg Gap), unlearning effectiveness and resilience to the Streisand effect (indicated by JSD), and efficiency (reflected in Time, measured in minutes).

	Metric (\downarrow)	Gold Std	Finetuning	NegGrad+	RndLbl	BadT	SCRUB	SSD	UNSIR	SalUn	LoTUS
yIN	Avg Gap	0.0000	0.2975	0.3250	<u>0.2925</u> 92.27 + 6 + 2	0.3125	0.4200	0.3650	0.5075	0.2925 91.01 + 9.59	0.0925
Tin	Time (min.)	$42.15_{\pm 16.05}$	$3.23_{\pm 0.01}$	$3.24_{\pm 0.03}$	3.27 ± 6.43 3.27 ± 0.03	$1.59_{\pm 0.01}$	$4.05_{\pm 0.03}$	$\frac{34.90\pm14.21}{3.19\pm0.03}$	$1.01_{\pm 0.01}$	$3.98_{\pm 0.01}$	$\frac{1.30_{\pm 0.02}}{1.30_{\pm 0.02}}$
100 aver	Avg Gap JSD×1e4	0.0000 0.00+0.00	$\frac{0.2825}{101.48_{+2.87}}$	0.3725 108.50+2.59	0.2925 102.66+311	0.3000 78.65+312	0.3225 64.09+8.71	0.4325 45.19+9 19	0.4050 76.28 _{+6.88}	$\frac{0.2850}{100.93_{\pm 2,44}}$	0.1200 25.46+1.41
င် B	Time (min.)	4.00 ± 0.11	$0.43_{\pm 0.00}$	$0.44_{\pm 0.01}$	$0.45_{\pm 0.00}$	$0.26_{\pm 0.01}$	$0.55_{\pm 0.00}$	$0.83_{\pm 0.03}$	$0.20_{\pm 0.01}$	$1.16_{\pm 0.01}$	$0.23_{\pm 0.01}$

Table 9. Class Unlearning with ResNet18 models and the TinyImageNet (TinyIN) and CIFAR-100 (C-100) datasets. We highlight the best and second-best scores.

Metric (\downarrow)	BadT	SCRUB	SSD	UNSIR	SalUn	LoTUS
Avg. Gap JSD×1e4 U Time (min)	$\begin{array}{c} 0.0575 \\ 0.06_{\pm 0.01} \\ 15.04_{\pm 0.03} \end{array}$	$\begin{array}{c} 0.0350 \\ \textbf{0.01}_{\pm \textbf{0.00}} \\ 16.82_{\pm 0.03} \end{array}$	$\begin{array}{c} 0.0350 \\ \textbf{0.01}_{\pm \textbf{0.00}} \\ 18.69_{\pm 0.06} \end{array}$	$\begin{array}{c} 0.0375 \\ 0.02_{\pm 0.00} \\ 18.33_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0375 \\ 0.10_{\pm 0.01} \\ 38.08_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0225 \\ 0.01_{\pm 0.00} \\ 13.79_{\pm 0.02} \end{array}$
$\frac{1}{10000000000000000000000000000000000$	$\begin{array}{c} 0.0600 \\ 0.04_{\pm 0.01} \\ 15.10_{\pm 0.20} \end{array}$	$\begin{array}{c} 0.0125 \\ \textbf{0.00}_{\pm \textbf{0.00}} \\ 16.99_{\pm 0.35} \end{array}$	$\begin{array}{c} 0.0150 \\ \textbf{0.00}_{\pm \textbf{0.00}} \\ 19.03_{\pm 0.54} \end{array}$	$\begin{array}{c} 0.0150 \\ \textbf{0.00}_{\pm \textbf{0.00}} \\ 18.33_{\pm 0.02} \end{array}$	$\begin{array}{c} \textbf{0.0050} \\ 0.02_{\pm 0.00} \\ 37.93_{\pm 0.18} \end{array}$	$\begin{array}{c} 0.0050 \\ 0.00_{\pm 0.00} \\ 14.09_{\pm 0.53} \end{array}$
$ \begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\$	$\begin{array}{c} 0.3050 \\ 0.55_{\pm 0.04} \\ 0.58_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.2225 \\ 0.44_{\pm 0.02} \\ 0.62_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.2225 \\ 0.44_{\pm 0.02} \\ 1.29_{\pm 0.03} \end{array}$	$\begin{array}{c} 0.2925 \\ 0.65_{\pm 0.23} \\ 0.72_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.3300 \\ 6.25_{\pm 0.45} \\ 1.49_{\pm 0.01} \end{array}$	$\begin{array}{c} 0.1725 \\ 0.28_{\pm 0.00} \\ 0.57_{\pm 0.01} \end{array}$
$ \overset{\textbf{Trime (min)}}{\overset{\textbf{O}}{\overset{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}}{\overset{\textbf{O}}{\overset{\textbf{O}}{\overset{\textbf{O}}}}}}}}}}$	0.0625 0.18±0.02 0.57±0.02	$\begin{array}{c} 0.1075 \\ 0.14_{\pm 0.00} \\ 0.61_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.1025 \\ 0.13_{\pm 0.00} \\ 1.27_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.1025 \\ 0.21_{\pm 0.05} \\ 0.73_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1300 \\ 1.40_{\pm 0.02} \\ 1.49_{\pm 0.01} \end{array}$	0.0650 $0.09_{\pm 0.01}$ $0.57_{\pm 0.00}$

Table 3. Scaling up the Forget set to 50% of the training sets: LoTUS outperforms state-of-the-art approaches in all metrics. Basic unlearning methods (Finetuning, NegGrad+, Rnd Labeling) are more efficient, but less effective than LoTUS.



Metric (↓)	Gold Std	Finetuning	NegGrad+	RndLbl	BadT	SCRUB	SSD	UNSIR	SalUn	LoTUS	synthetic D_u
Avg Gap JSD×1e4 Time (min.)	$\begin{array}{c c} 0.0000 \\ 0.00_{\pm 0.00} \\ 111.00_{\pm 1.99} \end{array}$	$\begin{array}{c} \underline{0.0075} \\ \textbf{0.01}_{\pm \textbf{0.00}} \\ 11.33_{\pm 0.03} \end{array}$	$\begin{array}{c} 0.0125 \\ 0.03_{\pm 0.00} \\ 12.61_{\pm 0.03} \end{array}$	$0.0125 \\ \frac{0.02_{\pm 0.01}}{12.78_{\pm 0.01}}$	$\begin{array}{c} 0.0375 \\ 0.12_{\pm 0.03} \\ 8.97_{\pm 0.02} \end{array}$	$\begin{array}{c} \textbf{0.0050} \\ \textbf{0.01}_{\pm 0.00} \\ 16.66_{\pm 0.02} \end{array}$	$\frac{0.0075}{13.65_{\pm 0.02}}$	$\begin{array}{c} 0.0100 \\ \textbf{0.01}_{\pm \textbf{0.01}} \\ 10.68_{\pm 0.02} \end{array}$	$\begin{array}{c} 0.0125 \\ \textbf{0.01}_{\pm \textbf{0.01}} \\ 37.97_{\pm 0.19} \end{array}$	$\begin{array}{c} \textbf{0.0050} \\ \textbf{0.01}_{\pm 0.00} \\ \hline 7.34_{\pm 0.19} \end{array}$	$\frac{0.0075}{0.01_{\pm 0.00}}$ 7.25 $_{\pm 0.06}$
$\stackrel{\infty}{\underset{\text{ZZ}}{\cong}} \begin{array}{l} \text{Avg Gap} \\ \text{JSD} \times 1e4 \\ \text{Time (min.)} \end{array}$	$\begin{array}{c} 0.0000 \\ 0.00_{\pm 0.00} \\ 5.32_{\pm 1.18} \end{array}$	$\begin{array}{c} 0.1375 \\ 1.03_{\pm 0.24} \\ 0.43_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.0975 \\ 1.06_{\pm 0.21} \\ 0.49_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.0925 \\ 1.00_{\pm 0.26} \\ 0.57_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.2650 \\ 2.39_{\pm 2.03} \\ 0.33_{\pm 0.00} \end{array}$	$\frac{0.0750}{0.41_{\pm 0.09}}\\ \frac{0.58_{\pm 0.00}}{0.58_{\pm 0.00}}$	$\begin{array}{c} 0.0825 \\ 0.82_{\pm 0.57} \\ 0.54_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1075 \\ 0.65_{\pm 0.05} \\ 0.45_{\pm 0.00} \end{array}$	$\begin{array}{c} 0.1800 \\ 1.09_{\pm 0.05} \\ 1.56_{\pm 0.01} \end{array}$	$\begin{array}{c} \underline{0.0350} \\ 0.32_{\pm 0.04} \\ 0.29_{\pm 0.00} \end{array}$	$\begin{array}{c} \textbf{0.0300} \\ \textbf{0.32}_{\pm 0.03} \\ \hline 0.30_{\pm 0.01} \end{array}$

Table 2. Unlearning 10% of CIFAR-10. LoTUS outperforms state-of-the-art approaches, both when the calibration/unseen set D_u consists of real data (•) and when it consists of synthetic data (•) from the CIFAKE dataset. We highlight the best and second-best scores.

Contributions

Large-Scale Unlearning Benchmark
 Simulates Real-World Conditions
 (ImageNet1k & Limited Data Access)

✓ Retrain-Free Jensen-Shannon Divergence (RF-JSD) Metric:

Ensures Evaluation without the Gold Standard Model

✓ LoTUS

Increases Entropy In Predictions up to an Information-Theoretical Bound



Figure 1. Machine Unlearning via smoothing prediction probabilities: LoTUS eliminates sample-specific information (e.g., unique fur patterns in cat images) that the DNN memorized and exposed through overconfident predictions. Then, the DNN responds to unlearned samples as if they were never part of the training set.





DNNs memorize sample-specific information from training data → Over-confident predictions →
Exploited by Privacy Attacks
Increase Model's Uncertainty / Entropy in Predictions



How Much Uncertainty?

Information-Theoretical Framework → Use Unseen (during training) Images for Calibration

Upper-Bounding Uncertainty





X : General Population

 X_s : Filtered Population

Data Processing Inequality: $X \rightarrow X_s \rightarrow f_w(X_s)$

$$\underbrace{I(f_w(X_s); X_s)}_{\text{total information captured by}} = \underbrace{I(f_w(X_s); X)}_{\text{global}} + \underbrace{I(f_w(X_s); X_s | X)}_{\text{subset-specific}}_{\text{information}}$$

$$I(f_{\text{orig}}(X_s), X_s \in D_f) = I(f_{\text{orig}}(X_s), X) + I(f_{\text{orig}}(X_s); X_s | X)$$

$$I(f_{\text{gold}}(X_s), X_s \in D_f) = I(f_{\text{gold}}(X_s), X) + I(f_{\text{gold}}(X_s); X_s | X)$$

$$I(f_{\text{un}}(X_s); X_s) \stackrel{\Delta}{=} I(f_{\text{un}}(X_s); X) \stackrel{\Delta}{=}^* I(f_{\text{orig}}(X_s); X)$$

2

*instance-wise unlearning

Upper-Bounding Uncertainty





X: General Population

tł

 X_s : Filtered Population

Data Processing Inequality: $X \rightarrow X_s \rightarrow f_w(X_s)$

$$\underbrace{I(f_w(X_s); X_s)}_{\text{otal information captured by}} = \underbrace{I(f_w(X_s); X)}_{\text{global}} + \underbrace{I(f_w(X_s); X_s \mid X)}_{\text{subset-specific}}$$

$$I(f_{\text{orig}}(X_s), X_s \in D_f) = I(f_{\text{orig}}(X_s), X) + I(f_{\text{orig}}(X_s); X_s \mid X)$$

$$I(f_{\text{gold}}(X_s), X_s \in D_f) = I(f_{\text{gold}}(X_s), X) + I(f_{\text{gold}}(X_s); X_s \mid X)$$

$$I(f_{un}(X_s); X_s) \stackrel{\Delta}{=} I(f_{un}(X_s); X) \stackrel{\Delta}{=}^* I(f_{\text{orig}}(X_s); X)$$

 $I(f_{\text{orig}}(X_s), X_s \in D_u) = I(f_{\text{orig}}(X_s), X) + I(f_{\text{orig}}(X_s), X_s \mid X)$

 $I(f_{\mathrm{un}}(X_s), X_s \in D_f) = I(f_{\mathrm{orig}}(X_s), X_s \in D_u) \stackrel{**}{\Rightarrow}$ $\underbrace{H(X_s \in D_f) - H(X_s \in D_f | f_{\mathrm{un}}(X_s)) =}_{H(X_s \in D_u) - H(X_s \in D_u | f_{\mathrm{orig}}(X_s)) \Rightarrow}$ $H(X_s \in D_f | f_{\mathrm{un}}(X_s)) = H(X_s \in D_u | f_{\mathrm{orig}}(X_s))$

**distributional similarity

*instance-wise unlearning

Upper-Bounding Uncertainty





X: General Population

 X_s : Filtered Population

Data Processing Inequality: $X \rightarrow X_s \rightarrow f_w(X_s)$

$$\underbrace{I(f_w(X_s); X_s)}_{\text{total information captured by}} = \underbrace{I(f_w(X_s); X)}_{\text{global}} + \underbrace{I(f_w(X_s); X_s \mid X)}_{\text{subset-specific}}$$

$$I(f_{\text{orig}}(X_s), X_s \in D_f) = I(f_{\text{orig}}(X_s), X) + I(f_{\text{orig}}(X_s); X_s \mid X)$$

$$I(f_{\text{gold}}(X_s), X_s \in D_f) = I(f_{\text{gold}}(X_s), X) + I(f_{\text{gold}}(X_s); X_s \mid X)$$

$$I(f_{\text{un}}(X_s); X_s) \triangleq I(f_{\text{un}}(X_s); X) \triangleq^* I(f_{\text{orig}}(X_s); X)$$



$$I(f_{un}(X_s), X_s \in D_f) = I(f_{orig}(X_s), X_s \in D_u) \stackrel{**}{\Rightarrow}$$
$$\underbrace{H(X_s \in D_f) - H(X_s \in D_f | f_{un}(X_s)) =}_{H(X_s \in D_u) - H(X_s \in D_u | f_{orig}(X_s)) \Rightarrow}$$
$$H(X_s \in D_f | f_{un}(X_s)) = H(X_s \in D_u | f_{orig}(X_s))$$

Fano's Inequality: Lower Prediction Error Probability Implies Lower $H(X_s | f_w(X_s))$

$$P_e \ge \frac{H(X_s \mid f_w(X_s)) - 1}{\log |A(X_s)|}$$

$$\operatorname{Acc}(f_{\operatorname{un}}, D_f) = \operatorname{Acc}(f_{\operatorname{orig}}, D_u)$$



**distributional similarity

*instance-wise unlearning



Temperature Parameter (Dynamically Adjusted) $\tau_d = \exp\left(\alpha \cdot \left(\operatorname{Acc}(f_{\operatorname{un}}, D_f) - \operatorname{Acc}(f_{\operatorname{orig}}, D_u)\right)\right)$

Gumbel – Softmax Activation Function

$$p_{i} = gs(\pi, \tau) = \frac{\exp\left(\left(\log \pi_{i} + g_{i}\right)/\tau\right)}{\sum_{j=1}^{k} \exp\left(\left(\log \pi_{j} + g_{j}\right)/\tau\right)}, \ i = 1, \dots, k$$

Loss in the Teacher – Student Framework

- T: Pre-trained Model + Output Perturbation
- S: Unlearned Model initialized with T's weights





Figure 2. Contribution of the dynamically adjusted temperature τ_d to convergence toward the objective $\Delta Acc = Acc(f_{un}, D_f) - Acc(f_{orig}, D_u) = 0$. The steps are denoted as follows: ①: sharp step towards objective, ②: smaller step (proportional to ΔAcc), ③: drastic accuracy restoration when over-unlearning.

Why Gumbel – Softmax?

			Vision Transf	former		ResNet18				
_		TinyImageNet	CIFAR-100	CIFAR-10	MUFAC	TinyImageNet	CIFAR-100	CIFAR-10	MUFAC	
Avg Gap	Gumbel- Softmax	0.0150	0.125	0.0050	0.0200	0.1675	0.1200	0.0350	0.1250	
	Softmax with Temperature	0.0675	0.0225	0.0050	0.0200	0.1850	0.1075	0.0675	0.1175	
(1e4	Gumbel- Softmax	0.03	0.04	0.01	0.05	0.62	1.67	0.32	6.90	
JSD	Softmax with Temperature	0.15	0.04	0.01	0.08	0.65	1.36	0.41	7.33	

Table 10. Contribution of Gumbel noise into the activation function. Ablation analysis using Gumbel-Softmax and Softmax with Temperature as activation functions. LoTUS performs better with Gumbel-Softmax in the majority of the benchmarks.

Conclusion & Thank you!

✓ LoTUS:

- 1st Entropy-based Unlearning Method that Meticulously Increases Uncertainty up to an Information-Theoretical Bound
- Outperforms SoTA in Balancing Forgetting-Retention, Unlearning Effectiveness & Resilience to the Streisand Effect, and Efficiency
- ✓ Retrain-Free Jensen-Shannon Divergence (RF-JSD) Metric:
 - Ensures Evaluation without the Gold Standard Model
 - Strong Correlations with the Established JSD Metric
 - More Interpretable & Efficient than the existing ZRF Metric
- ✓ Large-Scale Unlearning Benchmark Simulating Real-World Conditions



