# SCAFFOLDING A STUDENT TO INSTILL KNOWLEDGE

ANIL KAG [1], DURMUS ALP EMRE ACAR [1], ADITYA GANGRADE [2], VENKATESH SALIGRAMA [1]

ICLR 2023

## Problem: Teaching a Tiny Model

➤ Knowledge Distillation Setup



$$-\frac{1}{N}\sum_{i=1}^{N}\sum_{y} t_y(x_i)\log\frac{s_y(x_i)}{t_y(x_i)}$$

Teacher: T — $l_T(x)$ — $\tau$ — $t(x)$

Student: S — $l_S(x)$ — $\tau$ — $s(x)$

Logits — Softened Probabilities

KL ( S || T )

➤ Tiny Models have limited capacity

Panda — Easy

Teacher: T — Bear — Hard — Student: S

Distillation focuses equally on easy / hard examples

➤ Our Proposal : Scaffold hard-to-learn inputs for student by exploiting teacher

Panda — Please learn — Scaffold / Help

Teacher: T — Bear — Ignore, if you cannot learn — Student: S

Teacher provides extra help on hard examples

## DiSK: Distilling Scaffolded Knowledge

➤ DiSK Setup



Teacher: T — t(x)

Help: g — Scaffold

Student: S — s(x) — Softened Probabilities

❑ g is computationally cheap
❑ g is used only in training

$$-\frac{1}{N}\sum_{i=1}^{N}\sum_{y} t_y(x_i)\log\frac{s_y(x_i)+g(x_i)}{t_y(x_i)}$$

KL ( S + g || T ) + Budget(g)

Avoid trivial solution (help on all examples)

$$\max\left(0, -\frac{1}{N}\sum_{i=1}^{N} g(x_i)\log s_{y_i}(x_i) - \delta\right)$$

➤ DiSK smoothens the student's loss-landscape, often eliminating suboptimal minima

Example: 1D Intervals

Teacher = 2 intervals
Student = 1 interval

KD Loss Surface — 83%, 67%, 87%
DISK Loss Surface — 87%

Teacher Minima (100% Accuracy)
Student Minima (67% Accuracy)
Student Minima (83% Accuracy)
Student Minima (87% Accuracy)
Data Distribution
0 1 2 3 4 5 6 7 8 9

Example: 2D Gaussians

Teacher = 3 layer NN
Student = 2 layer NN

KD Loss Surface — DISK Loss Surface

Data | S Min-A (Acc=0.7) | S Min-B (Acc=0.8)
T (Acc=1.0) | S Min-C (Acc=0.9) | S Min-D (Acc=1.0)

## Empirical Evaluation

➤ CIFAR-100 (Tiny Students): Up to 4% higher accuracy compared to KD

| Teacher | Teacher MACs | Student | Student MACs | CE | KD | DiSK |
|---|---|---|---|---|---|---|
| ResNet10-l | 64M | Resnet10-s | 4M | 52.16 | 54.92 | **58.14** |
| 71.99% | | Resnet10-m | 16M | 65.24 | 66.96 | **70.03** |
| ResNet18 | 555M | Resnet10-s | 4M | 52.16 | 55.76 | **58.11** |
| 76.56 | | Resnet10-m | 16M | 65.24 | 68.09 | **69.86** |

➤ CIFAR-100 (Standard Students): Up to 2.5% higher accuracy compared to KD

| Teacher | Teacher MACs | Student | Student MACs | CE | KD | DiSK |
|---|---|---|---|---|---|---|
| ResNet32x4 | 1083M | ShuffleNetV2 | 45M | 73.74 | 79.13 | **80.23** |
| 81.45% | | MobileNetV2x2 | 22M | 69.24 | 76.05 | **77.24** |
| Wide-ResNet | 327M | ShuffleNetV2 | 45M | 73.74 | 75.81 | **78.33** |
| 78.41 | | MobileNetV2x2 | 22M | 69.24 | 73.92 | **76.32** |

➤ ImageNet-1K: More than 1% higher accuracy compared to KD

| Teacher | Teacher MACs | Student | Student MACs | CE | KD | DiSK |
|---|---|---|---|---|---|---|
| ResNet50 | 4.12B | ResNet18 | 1.82B | 69.73 | 71.29 | **72.35** |
| ViT-Large | 59.65B | ViT-Tiny | 1.07B | 75.45 | 76.61 | **77.86** |
| ViT-Large | 59.65B | DeiT-Tiny | 1.07B | 72.2 | 74.5 | **75.59** |

➤ DiSK can be integrated with other procedures such as feature matching

| Teacher | Teacher MACs | Student | Student MACs | FitNet | SimKD | SimKD + DiSK |
|---|---|---|---|---|---|---|
| Wide-ResNet | 327M | ResNet8x4 | 177M | 75.02 | 76.75 | **77.13** |
| 78.41 | | Wide-ResNet-40-1 | 83M | 74.17 | 75.56 | **76.21** |

➤ https://github.com/anilkagak2/DiSK_Distilling_Scaffolded_Knowledge