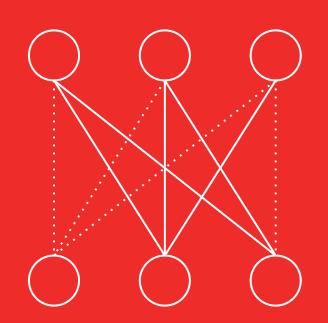
# What's Left After Distillation?

How Knowledge Transfer Impacts Fairness and Bias

# Yani Ioannou

Schulich Research Chair / Assistant Professor Schulich School of Engineering, University of Calgary



**Hanyang University** 

August 26th, 2025

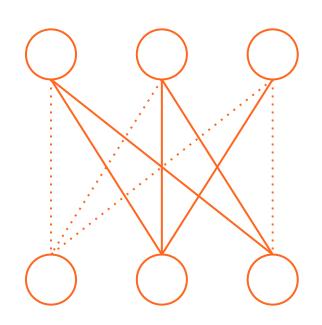




# **Short Biography**

# 1. Short Biography

- 2. Calgary ML Lab
- 3. Recent Research Highlights
- 4. Distillation and Fairness





# Biography: Yani Ioannou PhD, University of Cambridge, 2018

- Prof. Roberto Cipolla (Department of Engineering)
   Dr. Antonio Criminisi (Microsoft Research)
- Microsoft Research PhD Scholarship

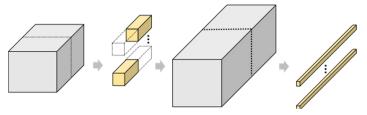
### Training CNNs with Low-Rank Filters for Efficient Image Classification.

Yani Ioannou, Duncan Robertson, Jamie Shotton, Roberto Cipolla, Antonio Criminisi. ICLR 2016

### **Deep Roots: Improving CNN Efficiency with Hierarchical Filter Groups**

Yani Ioannou, Duncan Robertson, Roberto Cipolla, Antonio Criminisi. CVPR 2017







2019 - 2020

**Google Brain (Toronto) Visiting Researcher** 

# **Gradient Flow in Sparse Neural Networks and How Lottery Tickets Win**

Utku Evci, Yani Ioannou, Cem Keskin, Yann Dauphin AAAI 2022 Oral Presentation



# **Industry & Applied Al Experience**





**Augmented Reality** 





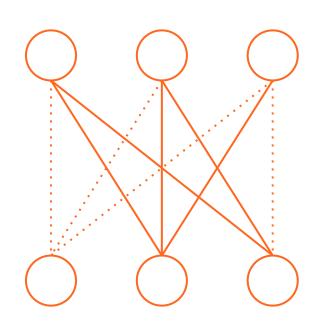


Blood Cell/Malaria Classification



# Recent Research Highlights

- 1. Short Biography
- 2. Calgary ML Lab
- 3. Recent Research Highlights
- 4. Distillation and Fairness





# DYNAMIC SPARSE TRAINING WITH STRUCTURED SPARSITY

Mike Lasby<sup>1</sup>, Anna Golubeva<sup>2,3</sup>, Utku Evci<sup>4</sup>, Mihai Nica<sup>5,6</sup>, Yani A. Ioannou<sup>1</sup>
<sup>1</sup>University of Calgary, <sup>2</sup>Massachusetts Institute of Technology, <sup>3</sup>IAIFI
<sup>4</sup>Google DeepMind, <sup>5</sup>University of Guelph, <sup>6</sup>Vector Institute for AI \*

### ABSTRACT

Dynamic Sparse Training (DST) methods achieve state-of-the-art results in sparse neural network training, matching the generalization of dense models while enabling sparse training and inference. Although the resulting models are highly sparse and theoretically less computationally expensive, achieving speedups with unstructured sparsity on real-world hardware is challenging. In this work, we propose a sparse-to-sparse DST method, Structured RigL (SRigL), to learn a variant of fine-grained structured N:M sparsity by imposing a constant fan-in constraint. Using our empirical analysis of existing DST methods at high sparsity, we additionally employ a neuron ablation method which enables SRigL to achieve state-of-the-art sparse-to-sparse structured DST performance on a variety of Neural Network (NN) architectures. Using a 90% sparse linear layer, we demonstrate a real-world acceleration of 3.4×/2.5× on CPU for online inference and 1.7×/13.0× on GPU for inference with a batch size of 256 when compared to equivalent dense/unstructured (CSR) sparse layers, respectively.

### 1 Introduction

Dynamic Sparse Training (DST) methods such as RigL (Evci et al., 2021) are the state-of-the-art in sparse training methods for Deep Neural Networks (DNNs). DST methods typically learn *unstructured* masks resulting in 85–95% fewer weights than dense models, while maintaining dense-like generalization and typically outperforming masks found via pruning. Furthermore, sparse-to-sparse DST algorithms are capable of employing sparsity both during training and inference, unlike pruning and dense-to-sparse DST methods such as SR-STE (Zhou et al., 2021) which only exploit sparsity at inference time.

While models trained with DST methods are highly sparse and enable a large reduction in Floating Point Operations (FLOPs) in theory, realizing these speedups on hardware is challenging when the sparsity pattern is unstructured. Even considering recent advances in accelerating unstructured Sparse Neural Networks (SNNs) (Gale et al., 2020; Elsen et al., 2020; Ji & Chen, 2022), structured sparsity realizes much stronger acceleration on real-world hardware. On the other hand, structured sparse pruning often removes salient weights, resulting in worse generalization than comparable unstructured SNNs for the same sparsity level (Fig. 1a). Our work presents a best-of-both-worlds approach: we exploit the DST framework to learn both a highly-sparse and structured representation while maintaining generalization performance. In summary, our work makes the following contributions:



Michael Lasby
PhD Student
(Currently interning at Cerebras)



Anna Golubeva
Postdoctoral Fellow, MIT/IAIFI



**Utku Evci**Research Scientist
Google DeepMind



**Mihai Nica** Associate Professor U. Guelph/Vector

 International Conference for Learning Representations (ICLR) 2024





### What's Left After Distillation? How Knowledge Transfer Impacts Fairness and Bias

Aida Mohammadshahi

aida.mohammadshahi@ucalgary.ca

Yani Ioannou

yani.ioannou@ucalgary.ca

Department of Electrical and Software Engineering Schulich School of Engineering, University of Calgary Calgary, AB, Canada

Reviewed on OpenReview: https://openreview.net/forum?id=xBbj46Y2fN

### Abstract

Knowledge Distillation is a commonly used Deep Neural Network (DNN) compression method, which often maintains overall generalization performance. However, we show that even for balanced image classification datasets, such as CIFAR-100, Tiny ImageNet and ImageNet, as many as 41% of the classes are statistically significantly affected by distillation when comparing class-wise accuracy (i.e. class bias) between a teacher/distilled student or distilled student/non-distilled student model. Changes in class bias are not necessarily an undesirable outcome when considered outside of the context of a model's usage. Using two common fairness metrics, Demographic Parity Difference (DPD) and Equalized Odds Difference (EOD) on models trained with the CelebA, Trifeature, and HateXplain datasets, our results suggest that increasing the distillation temperature improves the distilled student model's fairness and the distilled student fairness can even surpass the fairness of the teacher model at high temperatures, Additionally, we examine individual fairness, ensuring similar instances receive similar predictions. Our results confirm that higher temperatures also improve the distilled student model's individual fairness. This study highlights the uneven effects of distillation on certain classes and its potentially significant role in fairness, emphasizing that caution is warranted when using distilled models for sensitive application domains.

### 1 Introduction

DNNs require significant computational resources, resulting in large overheads in compute, memory, and energy. Decreasing this computational overhead is necessary for many real-world applications where these costs would otherwise be prohibitive, or even make their application infeasible—e.g. the deployment of DNNs on mobile phones or edge devices with limited resources (Chen et al., 2016; Cheng et al., 2018; Gupta and Agrawal, 2022; Menghani, 2023). To address this challenge, DNN model compression methods have been developed that reduce the size and complexity of DNNs while maintaining their generalization performance (Cheng et al., 2017). One such widely used model compression method is Knowledge Distillation (distillation) (Hinton et al., 2015). Distillation has found extensive application in both industry and academia across various domains of artificial intelligence, encompassing areas such as Natural Language Processing (NLP) (Jiao et al., 2019; Fu et al., 2021; Liu et al., 2020), speech recognition (Ng et al., 2011; Nsecifically image classification (Zhu et al., 2019; Chen et al., 2019; cut al., 2019; cut

Distillation involves transferring knowledge from a complex model with superior performance (referred to as the teacher) to a simpler model (known as the student). In practice this allows the student model to achieve comparable or even better generalization than the teacher model, while using far fewer parameters (Hinton et al., 2015; Gou et al., 2021). Despite the widespread use of distillation, evaluation of the impact of distillation since its proposal by (Hinton et al., 2015) has overwhelmingly focused almost exclusively on the impact it has on generalization performance (Cho and Hariharan, 2019; Mirzadeh et al., 2020).



Aida Mohammadshahi MSc (Defended Jan. 2025)

ML Developer @ AltaML

- Presented at NeurIPS WiML
   Workshop in Dec. 2024
- Published in Transactions in Machine Learning Research (TMLR), March 2025



### Sparse Training from Random Initialization: Aligning Lottery Ticket Masks using Weight Symmetry

Mohammed Adnan \*12 Rohan Jain \*1 Ekansh Sharma 32 Rahul G, Krishnan 32 Yani Joannou 1

### Abstract

The Lottery Ticket Hypothesis (LTH) suggests there exists a sparse LTH mask and weights that achieve the same generalization performance as the dense model while using significantly fewer parameters. However, finding a LTH solution is computationally expensive, and a LTH's sparsity mask does not generalize to other random weight initializations. Recent work has suggested that neural networks trained from random initialization find solutions within the same basin modulo permutation, and proposes a method to align trained models within the same loss basin. We hypothesize that misalignment of basins is the reason why LTH masks do not generalize to new random initializations and propose permuting the LTH mask to align with the new optimization basin when performing sparse training from a different random init. We empirically show a significant increase in generalization when sparse training from random initialization with the permuted mask as compared to using the non-permuted LTH mask, on multiple datasets (CIFAR-10/100 & ImageNet) and models (VGG11 & ResNet20/50). Our codebase for reproducing the results is publicly available at here.

### 1. Introduction

In recent years, foundation models have achieved state-ofthe-art results for different tasks. However, the exponential increase in the size of state-of-the-art models requires a similarly exponential increase in the memory and computational costs required to train, store and use these models decreasing the accessibility of these models for researchers

Proceedings of the 42<sup>nd</sup> International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

and practitioners alike. To overcome this issue, different model compression methods, such as pruning, quantization and knowledge distillation, have been proposed to reduce the model size at different phases of training or inference. Post-training model pruning (Han et al., 2016) has been shown to be effective in compressing the model size, and seminal works have demonstrated that large models can be pruned after training with minimal loss in accuracy (Gale et al., 2019; Han et al., 2015). While model pruning makes inference more efficient, it does not reduce the computational cost of training the model.

Motivated by the goal of training a sparse model from a random initialization, Frankle & Carbin (2019) demonstrated that training with a highly sparse mask is possible and proposed the Lottery Ticket Hypothesis (LTH) to identify sparse subnetworks that, when trained, can match the performance of a dense model. The key caveat is that a dense model must first be trained to find the sparse mask, which can only be used with the same random initialization that was used to train the dense model. Despite LTH seeing significant interest in the research community, LTH masks cannot be used to train from a new random initialization. Furthermore, it has been observed empirically that the LTH is impractical for finding a diverse set of solutions (Evci et al., 2022).

This posits our main research questions: How can we train a LTH mask from a different random initialization while main-aining good generalization? Would doing so find a more diverse set of solutions than observed with the LTH itself?

In this work, we try to understand why the LTH does not work for different random initializations from a weight-space symmetry perspective. Our hypothesis is that to reuse the LTH winning ticket mask with a different random initialization, the winning ticket mask obtained needs to be permuted such that it aligns with the optimization basin associated with the new random initialization. We illustrate our hypothesis in Figure 1.

To empirically validate our hypothesis, we obtain a sparse mask using Iterative Magnitude Pruning (IMP) (Renda et al., 2020; Han et al., 2015) on model A (from Figure 1) and show that given a permutation that aligns the optimization basin of model A and a new random initialization, the mask can



Adnan Mohammad PhD Student



Rohan Jain MSc (Defended May, 2025)



**Ekansh Sharma**PhD Student (U. Toronto/Vector)



**Ekansh Sharma**Assistant Professor (U. Toronto/Vector)

 Presented at International Conference on Machine Learning (ICML), July 2025





<sup>\*</sup>Equal contribution <sup>1</sup> Schulich School of Engineering, University of Calgary <sup>2</sup>Vector Institute for AI <sup>3</sup>Dept. of Computer Science, University of Toronto. Correspondence to: Mohammed Adnan cadnan.ahmad@ucalgary.ca>, Yani Ioannou </rd>



# Calgary ML @ NeurIPS 2024

UNIVERSITY OF CALGARY

Tuesday, December 10th, 2024

### Muslims in ML (MusiML) Workshop

2:30-3:00 p.m. (Lightning Talk) \* 4:30-5:00 p.m. (Internal Poster Session), 6:30-8:00 p.m (Joint Poster Session for Affinity Groups)

### \* A Closer Look at Sparse Training in Deep Reinforcement Learning.

Muhammad Athar Ganaie, Vincent Michalski, Samira Ebrahimi Kahou, Yani loannou.

This paper explores sparse training in DRL, highlighting methods to improve dynamic sparse training performance at high sparsity, underscoring the need for DRL-specific strategies.

### Long-Tail Learning with Language Model Guided Curricula.

**Mohammed Adnan**, Rahul Krishnan, **Yani loannou**.

Improving performance on long-tail classes by leveraging LLMs to build curricula.

### Women in Machine Learning (WiML) Workshop

6:30 p.m - 8:00 p.m (Joint Poster Session for Affinity Groups)

### Learning to Reweight Examples in Backdoor Defense.

Yufan Feng, Benjamin Tan, Yani loannou.

We extend the online sample reweighting method from robust learning to the context of backdoor defense.

### What's Left After Distillation? How Knowledge Transfer Impacts Fairness and Bias.

Aida Mohammadshahi, Yani loannou.

We explore the impact of knowledge distillation temperature on fairness for language and image classification models.

Friday, December 13th, 2024

### Main Conference

4:30-7:30 p.m. (Poster Session)

### Navigating Extremes: Dynamic Sparsity in Large Output Spaces.

Nasibullah Nasibullah, Erik Schultheis, **Mike Lasby, Yani loannou**, Rohit Babbar. Investigates Dynamic Sparse Training for large output spaces. Leveraging semi-structured sparsity, intermediate layers, and auxiliary loss, it enables end-to-end training with millions of labels.

Poster Location: Fast Exhibit Hall A-C #2004

Saturday, December 14th, 2024

### UniReps: Unifying Representations in Neural Models

4:30-7:30 p.m. (Poster Session)

### Winning Tickets from Random Initialization: Aligning Masks for Sparse Training.

Rohan Jain, Mohammed Adnan, Ekansh Sharma, Yani Ioannou.

Lottery Tickets can't be trained from random init. We show that permuting the mask to align with the new initialization's optimization basin results in a mask that better approaches LTH generalization.



- Neural Information Processing
   Systems (NeurIPS) 2024 in
   Vancouver
- 5 different works being presented by 6 CML students across main conference and workshops



# Calgary ML Lab @ ICML 2025



Monday, July 14th, 2025

### 4th Muslims in ML (MusiML) Workshop

West Meeting Room 211-214, 3:00 - 4:00 p.m. (Poster Session)

<u>Does Compression Exacerbate Large Language Models' Social Bias?</u>

Muhammad Athar Ganaie - Mohammed Adnan - Arfa Raja - Shaina Raza - Yani laannou

### Women in Machine Learning (WiML) Symposium

West Meeting Room 211-214, 3:00 - 4:00 p.m. (Poster Session)

Backdooring VLMs via Concept-Driven Triggers

Yufan Feng · Weimin Lyu · Yuxin Wang · Benjamin Tan · Yani loannou

Tuesday, July 15th, 2025

### Main Conference

East Exhibition Hall A-B, 11:00 a.m. - 1:30 p.m. (Poster Session)

Sparse Training from Random Initialization:

Aligning Lottery Ticket Masks using Weight Symmetry

Aligning Lottery Ticket Wasks using Weight Symmetry

Mohammed Adnan · Rohan Jain · Ekansh Sharma · Rahul G. Krishnan · Yani loannou

Poster Location: East Exhibition Hall A-B #E-2106

Friday, July 18th, 2025

### 3rd Workshop on High-dimensional Learning Dynamics (HiLD)

West Meeting Room 118-120, 4:45 - 5:30 p.m. (Poster Session)

**Understanding Normalization Layers for Sparse Training** 

Mohammed Adnan - Ekansh Sharma - Rahul G. Krishnan - Yani Joannou

Saturday, July 19th, 2025

### 3rd Workshop on Efficient Systems for Foundation Models (ES-FoMo III)

East Exhibition Hall A, 1:00 - 2:30 p.m (Poster Session)

SD2: Self-Distilled Sparse Drafters

Mike Lasby · Nish Sinnadurai · Valavan Manohararajah · Sean Lie · Yani loannou · Vithursan Thangarasa

### Data in Generative Models Workshop: The Bad, the Ugly, and the Greats (DIG-BUGS)

West Ballroom A, 3:00 - 3:45 p.m (Poster Session)

Backdooring VLMs via Concept-Driven Triggers

Yufan Feng · Weimin Lyu · Yuxin Wang · Benjamin Tan · Yani loannou

### Workshop on Technical Al Governance

West Meeting Room 109-110, 3:00 - 4:00 p.m (Poster Session)

**Exploring Functional Similarities of Backdoored Models** 

Yufan Feng · Benjamn Tan · Yani loannou

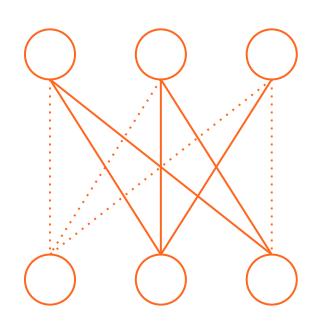


- International Machine
   Learning Conference (ICML)
   2025 in Vancouver in July
- Five students from CML Lab presenting 6 different works across 6 workshops and the main conference



# Recent Research Highlights

- 1. Short Biography
- 2. Calgary ML Lab
- 3. Recent Research Highlights
- 4. Distillation and Fairness

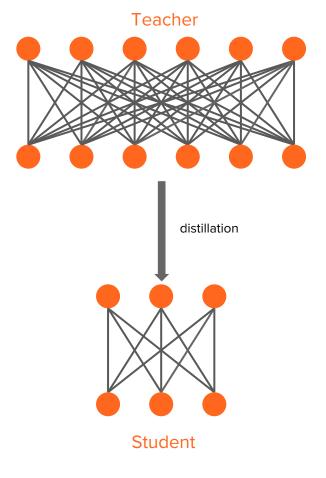




# Knowledge Distillation

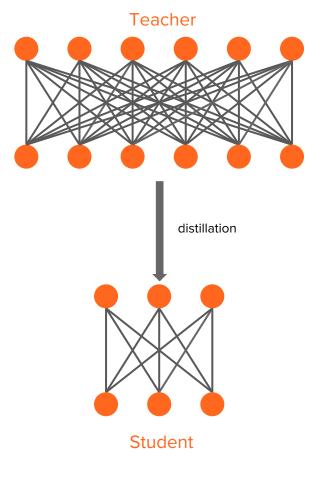
# What is Knowledge Distillation?

- A method of transferring "knowledge" from a larger model (or models) to a smaller model
- e.g. ensemble of models → single model
- Preserves generalization (test accuracy)
- Commonly used to compress large models
  - Large model → small model (Student)



# What is Knowledge Distillation?

- Commonly used to compress large models
  - Large model → small model (Student)
- Used extensively in industry to make models smaller for applications
  - Smaller models = cheaper compute costs
  - Smaller models enable mobile applications



# What is Knowledge Distillation?

- DeepSeek R1 (671B MoE Model)
  - Distilled smaller (1.5 70B) models, e.g. Llama
  - These smaller models are the models easier to use in practice

Model	AIME 2024		MATH-500	GPQA Diamond	LiveCode Bench	CodeForces	
	pass@1	cons@64	pass@1	pass@1	pass@1	rating	
GPT-40-0513	9.3	13.4	74.6	49.9	32.9	759	
Claude-3.5-Sonnet-1022	16.0	26.7	78.3	65.0	38.9	717	
OpenAI-o1-mini	63.6	80.0	90.0	60.0	53.8	1820	
QwQ-32B-Preview	50.0	60.0	90.6	54.5	41.9	1316	
DeepSeek-R1-Distill-Qwen-1.5B	28.9	52.7	83.9	33.8	16.9	954	
DeepSeek-R1-Distill-Qwen-7B	55.5	83.3	92.8	49.1	37.6	1189	
DeepSeek-R1-Distill-Qwen-14B	69.7	80.0	93.9	59.1	53.1	1481	
DeepSeek-R1-Distill-Qwen-32B	72.6	83.3	94.3	62.1	57.2	1691	
DeepSeek-R1-Distill-Llama-8B	50.4	80.0	89.1	49.0	39.6	1205	
DeepSeek-R1-Distill-Llama-70B	70.0	86.7	94.5	65.2	57.5	1633	

Table 5 | Comparison of DeepSeek-R1 distilled models and other comparable models on reasoning-related benchmarks.

### What's Left After Distillation? How Knowledge Transfer Impacts Fairness and Bias

Aida Mohammadshahi

aida.mohammadshahi@ucalgary.ca

Yani Ioannou
Department of Electrical and Software Engineering
Schulich School of Engineering, University of Calgary

yani.io annou@ucalgary.ca

Calgary, AB, Canada

 ${\bf Reviewed\ on\ OpenReview:\ https://openreview.\ net/forum?\ id=xBbj46Y2fN}$ 

### Abstract

Knowledge Distillation is a commonly used Deep Neural Network (DNN) compression method, which often maintains overall generalization performance. However, we show that even for balanced image classification datasets, such as CIFAR-100, Tiny ImageNet and ImageNet, as many as 41% of the classes are statistically significantly affected by distillation when comparing class-wise accuracy (i.e. class bias) between a teacher/distilled student or distilled student/non-distilled student model. Changes in class bias are not necessarily an undesirable outcome when considered outside of the context of a model's usage. Using two common fairness metrics, Demographic Parity Difference (DPD) and Equalized Odds Difference (EOD) on models trained with the CelebA, Trifeature, and HateXplain datasets, our results suggest that increasing the distillation temperature improves the distilled student model's fairness and the distilled student fairness can even surpass the fairness of the teacher model at high temperatures, Additionally, we examine individual fairness, ensuring similar instances receive similar predictions. Our results confirm that higher temperatures also improve the distilled student model's individual fairness. This study highlights the uneven effects of distillation on certain classes and its potentially significant role in fairness, emphasizing that caution is warranted when using distilled models for sensitive application domains.

### 1 Introduction

DNNs require significant computational resources, resulting in large overheads in compute, memory, and energy. Decreasing this computational overhead is necessary for many real-world applications where these costs would otherwise be prohibitive, or even make their application infeasible—e.g. the deployment of DNNs on mobile phones or edge devices with limited resources (Chen et al., 2016; Cheng et al., 2018; Gupta and Agrawal, 2022; Menghani, 2023). To address this challenge, DNN model compression methods have been developed that reduce the size and complexity of DNNs while maintaining their generalization performance (Cheng et al., 2017). One such widely used model compression method is Knowledge Distillation (distillation) (Hinton et al., 2015). Distillation has found extensive application in both industry and academia across various domains of artificial intelligence, encompassing areas such as Natural Language Processing (NLP) (Jiao et al., 2019; Firet al., 2012; Liu et al., 2012, openits of the compassion of the compa

Distillation involves transferring knowledge from a complex model with superior performance (referred to as the teacher) to a simpler model (known as the student). In practice this allows the student model to achieve comparable or even better generalization than the teacher model, while using far fewer parameters (Hinton et al., 2015; Gou et al., 2021). Despite the widespread use of distillation, evaluation of the impact of distillation since its proposal by (Hinton et al., 2015) has overwhelmingly focused almost exclusively on the impact it has on generalization performance (Cho and Hariharan, 2019; Mirzadeh et al., 2020).



Aida Mohammadshahi MSc (Defended Jan 2025)

- Presented at NeurlPS WiML
   Workshop in Dec. 2024
- Accepted at TMLR March 2025



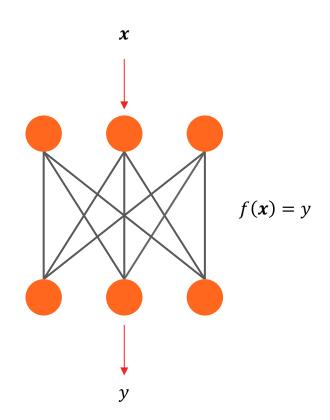




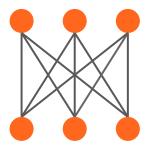
# Knowledge Distillation Introduction

# **Neural Networks as Functions**

- Neural Networks are function approximators
- A neural network learns a function
   f mapping an input x to an output y
- In practice, NNs for classification learn to predict a probability distribution p, from which the "hard" classification of a class y is made



# "Dark Knowledge"



$$f(\mathbf{x}) = \{0.4, 0.5, 0.1\}$$







- Trained models learn more than just how to predict labels
- They learn a function with rich knowledge of the domain
- An ImageNet model knows that a cat and dog are more similar to each other than an airplane



# Temperature Softmax

$$p_i = \frac{\exp\left(\frac{Z_i}{T}\right)}{\sum_j \exp\left(\frac{Z_j}{T}\right)}$$

$$f(x, T = 1) = \{0.09, 0.9, 0.01\}$$

$$f(x, T = 10) = \{0.4, 0.5, 0.1\}$$



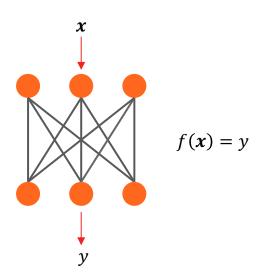
- A softmax p(z) gives us a probability output from logits z
- Distillation adds "temperature"
   T to softmax
- The typical softmax (T=1) gives very highly confident outputs for the target class, i.e. a "hard distribution
- Larger temp T gives "softer" distributions



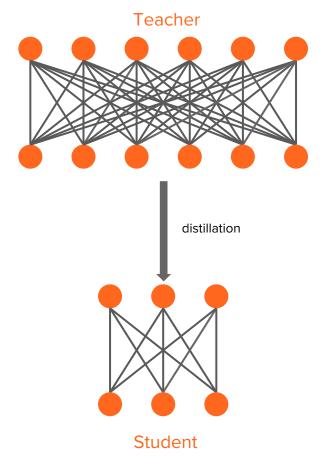
Knowledge
Distillation and
Fairness

# **Recall:** Neural Networks as Functions

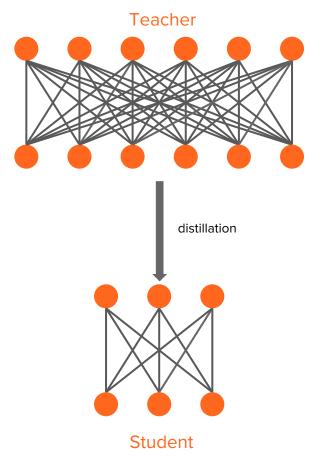
- NNs are function approximators
- A neural network learns a function
   f mapping an input x to an output y



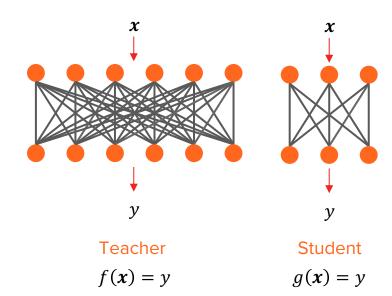
- When we distill a large teacher model to a small student, we often see generalization performance (test accuracy) maintained
- Does this mean that the Teacher and Student have learned similar functions?



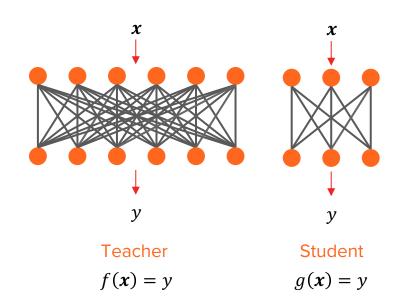
- When we distill a large teacher model to a small student, we often see generalization performance maintained
- Does this mean that the Teacher and Student have learned similar functions?
  - Not necessarily: accuracy is aggregate measure over many samples in test set



- When we distill a large teacher model to a small student, we often see generalization performance maintained
- However, student can learn different function than teacher
- Why does this matter?

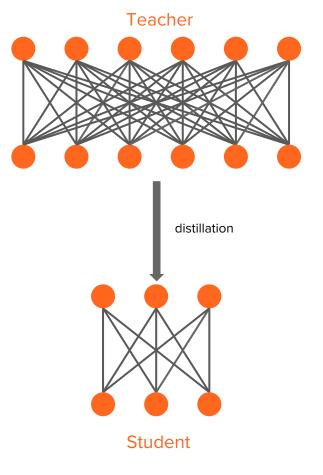


- When we distill a large teacher model to a small student, we often see generalization performance maintained
- However, student can learn different function than teacher
- Why does this matter?
- Student may learn different algorithmic biases than Teacher!



## **Research Questions**

- Q: What classes are significantly affected by distillation?
- Q: What is the impact of increase temperature T on the model's class biases?
- Q: How does distillation temperature affect group fairness?
- Q: How does distillation temperature affect individual fairness?



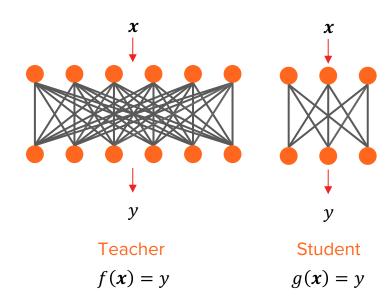
# Class-wise Bias: Analysis

- Q: What classes are significantly affected by distillation?
- Disagreement of the models f, g on predictions for  $x_n$ :

$$CMP(f(\mathbf{x}_n), g(\mathbf{x}_n)) = \begin{cases} 0 & \text{if } f(\mathbf{x}_n) = g(\mathbf{x}_n) \\ 1 & \text{if } f(\mathbf{x}_n) \neq g(\mathbf{x}_n) \end{cases}$$

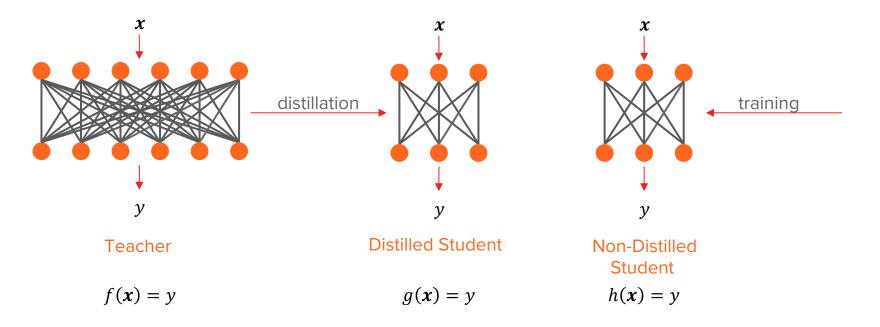
Compare the teacher f and distilled student g model's disagreement for each class c:

$$CMP(f(\mathbf{x}_n), g(\mathbf{x}_n))$$
 where  $(\mathbf{x}_n, y_n \mid y_n = c)$ 



# Class-wise Bias: Analysis

- Compare the teacher f and distilled student g model's disagreement for each class c:
- We use a non-distilled student h (trained from scratch) as a baseline



# Class-wise Bias: Models/Datasets

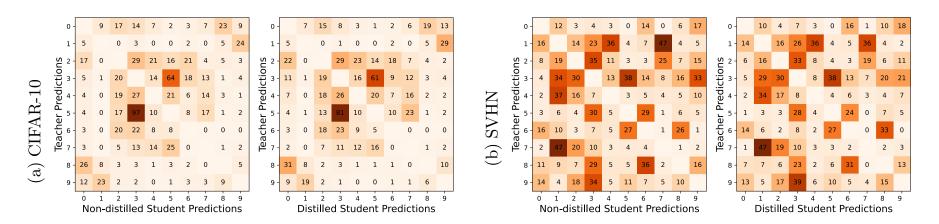


Figure 2: Class-wise Disagreement. Disagreement between a ResNet-56 teacher and ResNet-20 (left) non-distilled/(right) distilled student for (a) CIFAR-10 using T=9 and (b) SVHN using T=7. The diagonals are excluded since here both models predict the same class without any disagreement.

Dataset	Teacher (#param)	Student (#param)			
CIFAR-10/100, SVHN	ResNet56 (0.85M)	ResNet20 (0.27M)			

# Class-wise Bias: Analysis

- Q: What is the impact of increase temperature T on the model's class biases?
- TC = Teacher vs. Distilled Student, SC = Trained Student vs. Distilled Student

Table 1: Class-wise Bias and Distillation. The number of statistically significantly affected classes comparing the class-wise accuracy of teacher vs. Distilled Student (DS) models, denoted #TC, and Non-Distilled Student (NDS) vs. distilled student models, denoted #SC.

		CIFAR-100						${\bf ImageNet}$					
Teacher/Student		ResNet56/ResNet20			DenseNet169/DenseNet121		ResNet50/ResNet18			ViT-Base/TinyViT			
Model	Temp	Test Acc. (%)	#SC	#TC	Test Acc. (%)	#SC	#TC	Test Top-1 Acc. (%)	#SC	#TC	Test Top-1 Acc. (%)	#SC	#TC
Teacher	-	$70.87 \pm 0.21$	-	-	$72.43 \pm 0.15$	_	-	$76.1 \pm 0.13$	-	-	$81.02 \pm 0.07$	-	-
NDS	-	$68.39 \pm 0.17$	-	-	$70.17\pm0.16$	-	-	$68.64 \pm 0.21$	-	-	$78.68 \pm 0.19$	-	-
$\overline{\mathrm{DS}}$	2	$68.63 \pm 0.24$	5	15	$70.93 \pm 0.21$	4	12	$68.93 \pm 0.23$	77	314	$78.79 \pm 0.21$	83	397
DS	3	$68.92 \pm 0.21$	7	12	$71.08 \pm 0.17$	4	11	$69.12 \pm 0.18$	113	265	$78.94 \pm 0.14$	137	318
DS	4	$69.18 \pm 0.19$	8	9	$71.16 \pm 0.23$	5	9	$69.57 \pm 0.26$	169	237	$79.12 \pm 0.23$	186	253
DS	5	$69.77 \pm 0.22$	9	8	$71.42 \pm 0.18$	8	9	$69.85 \pm 0.19$	190	218	$79.51 \pm 0.17$	215	206
DS	6	$69.81 \pm 0.15$	9	8	$71.39 \pm 0.22$	8	8	$69.71 \pm 0.13$	212	193	$80.03 \pm 0.19$	268	184
DS	7	$69.38 \pm 0.18$	10	6	$71.34 \pm 0.16$	9	7	$70.05 \pm 0.18$	295	174	$79.62 \pm 0.23$	329	161
DS	8	$69.12 \pm 0.21$	13	6	$71.29 \pm 0.13$	11	7	$70.28 \pm 0.27$	346	138	$79.93 \pm 0.12$	365	127
DS	9	$69.35 \pm 0.27$	18	9	$71.51 \pm 0.23$	12	9	$70.52 \pm 0.09$	371	101	$80.16 \pm 0.17$	397	96
DS	10	$69.24 \pm 0.19$	22	11	$71.16 \pm 0.21$	14	10	$70.83 \pm 0.15$	408	86	$79.98 \pm 0.12$	426	78

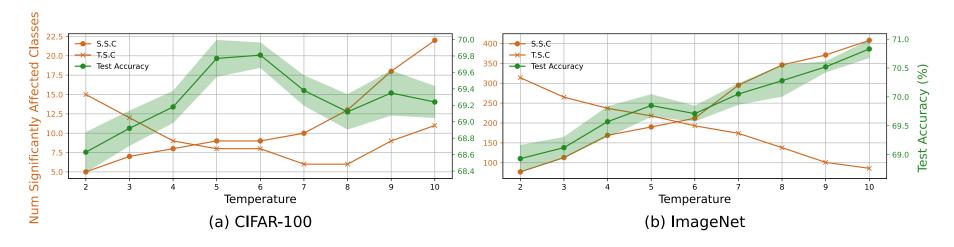


Figure 3: **Temperature vs. Test Accuracy/Class Bias.** Number of non-distilled vs. distilled student significantly affected classes (S.S.C.) and the number of teacher vs. distilled student significantly affected classes (T.S.C.) by distillation in (a) CIFAR-100 (ResNet-56/ResNet-20) and (b) ImageNet datasets (ResNet-50/ResNet-18), with 100 and 1000 total classes respectively. As the temperature used for distillation increases up to T=10, the S.S.C. rises for both datasets. For ImageNet, T.S.C. decreases, while for CIFAR-100, it first decreases and then slightly increases. The changes in the distilled student's test accuracy over all classes are also depicted in the figure.

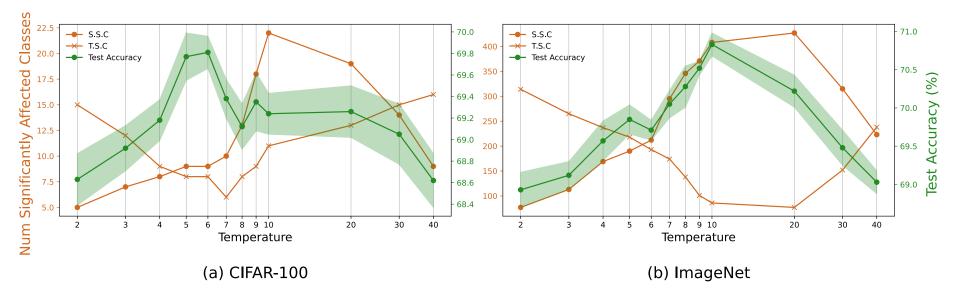
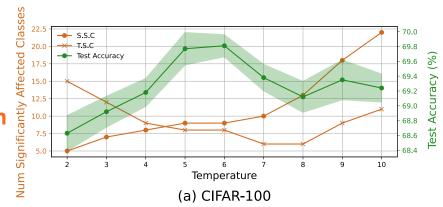
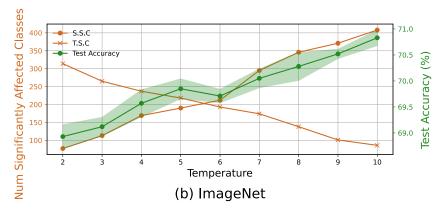


Figure 8: **Temperature vs. Test Accuracy/Class Bias.** Number of non-distilled vs. distilled student significantly affected classes (S.S.C.) and the number of teacher vs. distilled student significantly affected classes (T.S.C.) by distillation in (a) CIFAR-100 (ResNet-56/ResNet-20) and (b) ImageNet datasets (ResNet-50/ResNet-18), with 100 and 1000 total classes respectively. As the temperature used for distillation increases, the S.S.C. rises for both datasets up to a certain T, after which it decreases. Meanwhile, T.S.C. decreases first and then increases. The changes in the distilled student Test Accuracy over all classes are also depicted in the figure.

# Distillation and Class Bias

- When we distill a large teacher model to a small student, clearly the learned function is different
- Distillation does not affect class-wise accuracy uniformly
- However, a change in class bias alone is not meaningful (bad or good) in itself
- How can we judge if this is good or bad for applications?

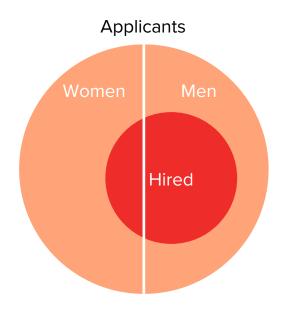




# Group Fairness

# Group Fairness

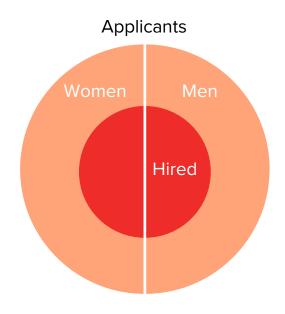
- A change in class bias alone is not meaningful (bad or good) in itself...
- What is clearly bad are unfair outcomes, i.e. a model not treating individuals from different groups equitably
- An example is a hiring system that accepts more men than women



### **Group Fairness:** Demographic Parity

- We want individuals belongs to different groups to have equal probability of a positive outcome
  - e.g. we want men and women to have equal odds of being hired
- Let A be the sensitive attribute (gender), and  $\hat{Y} = 1$  be the outcome (i.e. hired), we want:

$$P(\hat{Y} = 1 | A = a) = P(\hat{Y} = 1 | A = b)$$



### Group Fairness Metrics: Demographic Parity Difference

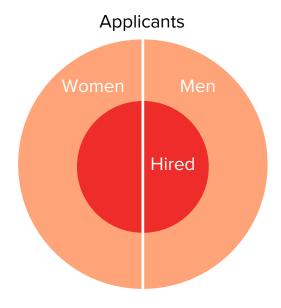
Demographic Parity:

$$P(Y = 1 | A = a) = P(Y = 1 | A = b)$$

 A metric based on demographic parity is the Demographic Parity Difference (DPD):

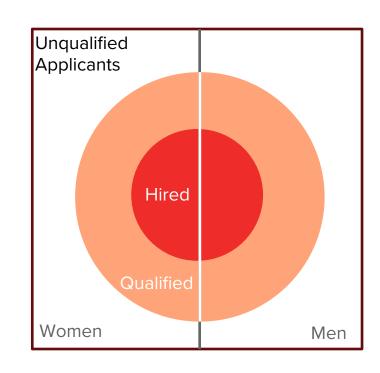
$$DPD = \max_{a \in A} P(Y = 1 | A = a) - \min_{a \in A} P(Y = 1 | A = a)$$

 DPD = 0 means perfectly fair in demographic parity fairness



# **Group Fairness:** Equalized Odds

- We want individuals to have equal probability of a positive or negative outcome given a condition is true
  - i.e. want groups to have equal probability of outcomes AND to have same TPR and FPR rates
  - e.g. we want men and women to have equal odds of being hired/not, if they are qualified
- Let A be the sensitive attribute,  $\hat{Y}$  be the outcome, and Y be the true label, we want:

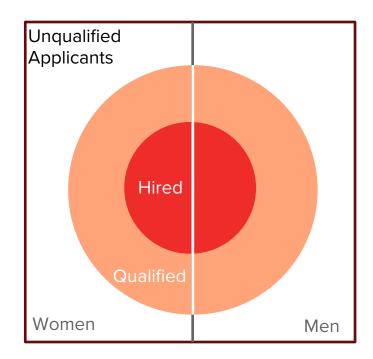


### Group Fairness Metrics: Demographic Parity Difference

Demographic Parity:

$$P(\hat{Y} = 1 | Y = y, A = a) = P(\hat{Y} = 1, Y = y | A = b)$$

- We use a metric based on equalized odds:
   Equalized Odds Difference (EOD)
- EOD=0 means perfectly fair in equalized odds fairness



#### CelebA Dataset



- CelebA is a dataset of celebrity photos
- CelebA has protected attributes, such as gender and age
- Also has independent attributes such as "smiling" or "glasses"
- Often used in fairness, but is also a deeply problematic dataset...

Deep Learning Face Attributes in the Wild. Liu, Ziwei and Luo, Ping and Wang, Xiaogang and Tang, Xiaoou. Proceedings of International Conference on Computer Vision (ICCV), 2015.



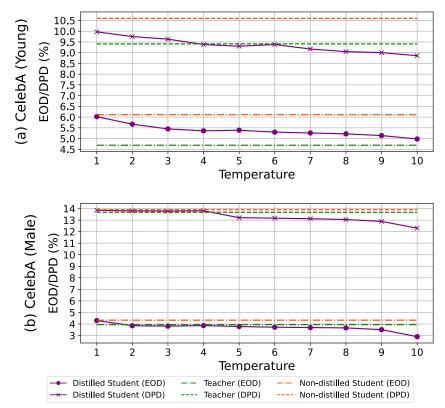


Figure 4: Evaluation of Fairness Metrics for Distilled Students in Computer Vision (CV). Equalized Odds Difference (EOD) and Demographic Parity Difference (DPD) are reported in % and lower values indicate improved fairness. (a) illustrates fairness metrics for the CelebA dataset with 'smiling' label concerning the 'Young' demographic attribute and (b) concerning the 'Male' demographic attribute. (c) presents fairness metrics for the Trifeature dataset with 'shape' label with regard to the 'color' attribute and (d) with regard to the 'texture' attribute. It is notable that the models are fairer for the Trifeature dataset compared to the CelebA dataset with lower values in metrics. The explanation lies in the fact that the Trifeature dataset maintains a balanced distribution of demographic attributes, while the CelebA dataset contains biases that mirror real-world disparities. As seen in the second column, the downward trend does not continue at very high temperatures (T=20,30,40), as the teacher model generates nearly uniform softmax outputs.

- ResNet50 (24M) → ResNet18
   (11.4M) distillation with CelebA dataset
- Protected attribute is Age (top) and Gender (bottom)
- Evaluated on "smiling" classification
- Fairness improves (i.e. EOD/DPD decreases) with higher T

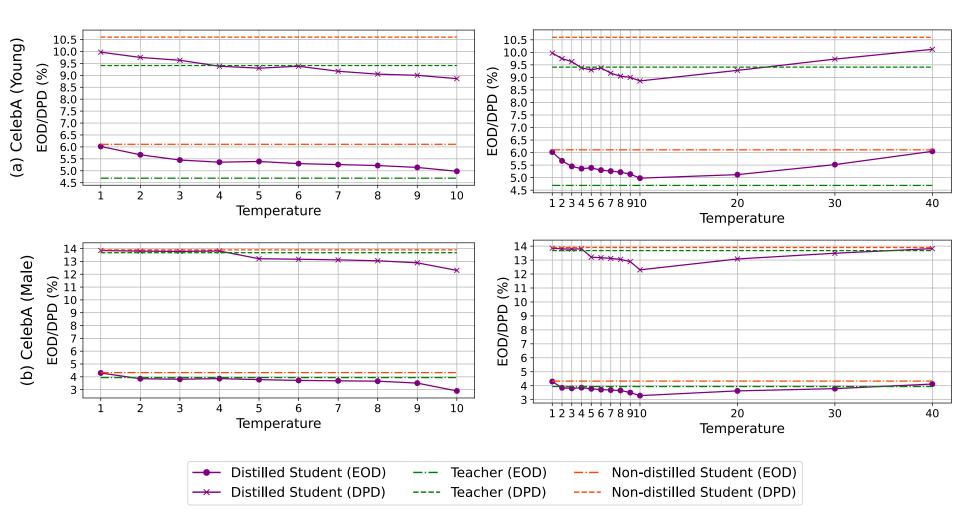


		CelebA (smiling)			
Teacher/Student:		ResNet-50 / ResNet-18			
Model	Temp	Test Acc. (%) ↑	EOD↓	DPD ↓	
Teacher	_	$93.09 \pm 0.08$	$\textbf{4.69} \pm \textbf{0.06}$	$9.41 \pm 0.11$	
NDS	_	$92.03 \pm 0.03$	$6.11 \pm 0.05$	$10.60\pm0.08$	
DS	1	$92.12 \pm 0.06$	$6.02 \pm 0.11$	$9.97 \pm 0.08$	
DS	2	$92.14 \pm 0.11$	$5.67 \pm 0.08$	$9.75 \pm 0.09$	
DS	3	$92.53 \pm 0.13$	$5.45 \pm 0.05$	$9.63 \pm 0.06$	
DS	4	$92.17 \pm 0.10$	$5.36 \pm 0.02$	$9.38 \pm 0.03$	
DS	5	$92.29 \pm 0.05$	$5.39 \pm 0.04$	$9.30 \pm 0.05$	
DS	6	$92.26 \pm 0.08$	$5.30 \pm 0.01$	$9.38 \pm 0.07$	
DS	7	$92.12 \pm 0.08$	$5.26 \pm 0.05$	$9.17 \pm 0.10$	
DS	8	$92.66 \pm 0.12$	$5.22 \pm 0.02$	$9.05 \pm 0.04$	
DS	9	$93.18 \pm 0.15$	$5.14 \pm 0.04$	$9.01 \pm 0.08$	
DS	10	$92.57 \pm 0.11$	$4.98 \pm 0.03$	$\textbf{8.86} \pm \textbf{0.04}$	

Table 2: Fairness Metrics and Distillation. The performance of teacher, Non-Distilled Student (NDS), and Distilled Student (DS) models with a range of temperatures T on the Trifeature and CelebA datasets. Fairness metrics are presented for Trifeature with regard to color attribute and for CelebA with regard to the Young demographic attribute. With increasing temperature, EOD and DPD have a downward trend signifying enhanced fairness. Mean and std. dev. are over five random inits.

- ResNet50 (24M) → ResNet18
   (11.4M) distillation with CelebA dataset
- Protected attribute is Age (top) and Gender (bottom)
- Evaluated on "smiling" classification
- Fairness improves (i.e. EOD/DPD decreases) with higher T





### HateExplain Dataset

Target groups	Categories
Race	African, Arabs, Asians, Caucasian, Hispanic
Religion	Buddhism, Christian, Hindu, Islam, Jewish
Gender	Men, Women
Sexual Orientation	Heterosexual, Gay
Miscellaneous	Indigenous, Refugee/Immigrant, None, Others

Table 3: Target groups considered for the annotation.

	Twitter	Gab	Total
Hateful	708	5,227	5,935
Offensive	2,328	3,152	5,480
Normal	5,770	2,044	7,814
Undecided	249	670	919
Total	9,055	11,093	20,148

Table 4: Dataset details. "Undecided" refers to the cases where all the three annotators chose a different class.

- HateExplain is a dataset used for detecting hate speech in online discourse
- Covers a range of protected groups (we use target groups aggregated, e.g. religion)
- We combine hateful/offensive to make task binary classification ("toxic" v.s. "normal")



HateXplain: A Benchmark Dataset for Explainable Hate Speech Detection. Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, Animesh Mukherjee. AAAI 2021.

#### BERT-Base (110M) → DistilBERT (66M) distillation

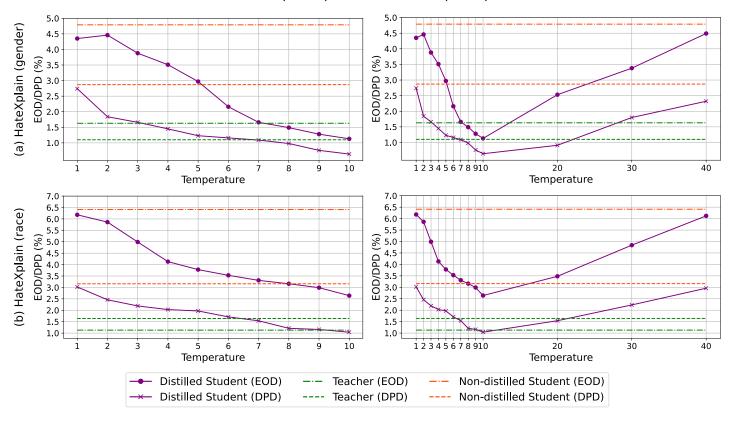


Figure 5: Evaluation of Fairness Metrics for Distilled Students in Natural Language Processing (NLP). Equalized Odds Difference (EOD) and Demographic Parity Difference (DPD) are reported in % and lower values indicate improved fairness. (a) illustrates fairness metrics for the HateXplain dataset concerning the 'gender' demographic attribute, and (b) with regard to the 'race' attribute. The teacher employed the BERT architecture, while the student used the DistilBERT architecture.

#### **Individual Fairness Metrics**

- Individual fairness metrics are very different
  - Group Fairness: individuals with different protected attributes should see similar outcomes
  - Individual Fairness: similar individuals should see similar outcomes
- Metric captures whether a model provides consistent predictions for semantically similar inputs, ensuring fairness at an individual level
- Lipschitz condition proposed by Dwork et al. (2012),
   smaller values = more fair

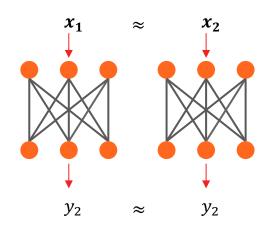


Table 4: Individual Fairness Metrics Across Datasets. Individual fairness scores for Teacher, Non-Distilled Student (NDS), and Distilled Student (DS) models across CelebA, Trifeature, and HateXplain datasets. Scores for DS models are reported for varying temperature values T.

		Individual Fairness $\downarrow$			
		CelebA	Trifeature	HateXplain	
Model	Temp	ResNet-50 / ResNet-18	ResNet-20 / LeNet-5	BERT-Base / DistilBERT	
Teacher	_	0.0407	0.016	0.0320	
NDS	_	0.124	0.0462	0.1078	
$\overline{\mathrm{DS}}$	1	0.113	0.0422	0.0994	
$_{ m DS}$	2	0.104	0.0407	0.0985	
$_{ m DS}$	3	0.0908	0.0393	0.0927	
DS	4	0.0906	0.0387	0.0882	
DS	5	0.0886	0.0384	0.0823	
DS	6	0.0799	0.0377	0.0768	
DS	7	0.0753	0.0356	0.0727	
DS	8	0.0712	0.0349	0.0689	
DS	9	0.0701	0.0341	0.0681	
DS	10	0.0697	0.0338	0.0654	

Clear increase in individual fairness with increased distillation temp

#### Conclusion

- Knowledge Distillation is pervasive in its use, you are likely affected by the decisions of a distilled model daily
- And yet the effect of distillation temperature on model fairness has not been looked at previously!
- We find across models, datasets and both vision and language modalities that distillation temperature affects the bias and fairness of models
- We also consistently find that higher distillation temperatures leads to more fair models
- In some cases, distilled models (with high T) can be fairer than even the (much larger) teacher model!



#### **Future Directions**

- Can distillation be an effective method of improving model fairness?
- Are there any trade offs to using large temperatures, less typically used with distillation in practice?
- Does distillation have a similar effect on LLMs, e.g. DeepSeek?



# Questions?

yani.ioannou@ucalgary.ca



amazon science



Aida Mohammadshahi MSc (Defended Jan 2025)

Research Engineer @ AltaML

What's Left After Distillation?
How Knowledge Transfer Impacts Fairness and Bias.
Aida Mohammadshahi, Yani Ioannou
Transactions in Machine Learning Research (TMLR), March 2025



