

Disagreement Augmentation: A Socratic Approach to Knowledge Distillation

Anonymous ICCV submission

Paper ID 15298

Abstract

001 *Disagreement plays a fundamental role in the learning*
 002 *process, driving deeper understanding and improved gen-*
 003 *eralization. Inspired by the Socratic method, we intro-*
 004 *duce Disagreement Augmentation (DA), a novel approach*
 005 *to knowledge distillation that leverages disagreement be-*
 006 *tween teacher and student models as a learning signal. Tra-*
 007 *ditional distillation methods primarily focus on aligning the*
 008 *student with the teacher, minimizing divergence to transfer*
 009 *knowledge effectively. However, this approach may over-*
 010 *look critical underrepresented or ambiguous regions of the*
 011 *data distribution. Our method actively augments training*
 012 *samples to maximize disagreement between the student and*
 013 *teacher, encouraging the student to resolve conflicting pre-*
 014 *dictions and develop a more robust approximation of the*
 015 *teacher. We evaluate DA in both an image classification*
 016 *setting and a reinforcement learning setting, demonstrat-*
 017 *ing improved student model performance over typical base-*
 018 *lines. These results highlight the potential of disagreement*
 019 *as a powerful augmentation strategy in knowledge distil-*
 020 *lation. Code and implementation details are available on*
 021 *GitHub.*

022 1. Introduction

023 Disagreement is a catalyst for learning. This principle,
 024 rooted in the Socratic method, highlights the role of produc-
 025 tive conflict in refining ideas and uncovering deeper truths.
 026 Socrates, through his method of dialectical questioning, of-
 027 ten encouraged his students to confront contradictions in
 028 their beliefs, leading to a richer understanding of complex
 029 concepts. This pedagogical approach, centered on optimiz-
 030 ing the interplay between opposing perspectives, inspires a
 031 new direction in knowledge distillation. We propose that,
 032 much like in Socratic dialogue, fostering disagreement be-
 033 tween the teacher and student models can drive learning and
 034 enhance model performance.

035 Knowledge distillation traditionally aims to minimize
 036 the divergence between a large, well-trained teacher model
 037 and a smaller student model, transferring the teacher’s ex-

038 pertise to create a compact, deployable version of the origi-
 039 nal system [9]. This approach focuses on alignment, where
 040 the student learns to emulate the teacher’s soft predictions,
 041 thereby inheriting its generalization capabilities. However,
 042 this paradigm overlooks the potential benefits of disagree-
 043 ment—particularly as a mechanism to explore underrepre-
 044 sented or ambiguous aspects of the data distribution [21].
 045 By intentionally optimizing for areas where the student and
 046 teacher disagree, we aim to emulate the Socratic process,
 047 leveraging conflict as a driver of more robust learning.

048 In this work, we introduce a novel method of data aug-
 049 mentation rooted in disagreement. Our approach, Disagree-
 050 ment Augmentation (DA), augments training samples to
 051 maximize divergence between the student and teacher mod-
 052 els. These disagreement-optimized examples challenge the
 053 student to reconcile conflicting predictions, encouraging it
 054 to develop a more nuanced approximation of the teacher.
 055 This method of structured disagreement offers a comple-
 056 mentary perspective to traditional distillation methods.

057 Beyond its conceptual motivation, DA is designed with
 058 practical advantages. In contrast to standard data augmen-
 059 tation, which typically applies predefined transformations
 060 (e.g., cropping, rotation, or noise injection) [2], DA gen-
 061 erates task-specific augmentations tailored to expose the
 062 weaknesses of the student model. This targeted approach
 063 encourages the student to learn from its mistakes more ef-
 064 fectively, leading to improved generalization and robust-
 065 ness. Moreover, DA aligns with recent efforts in self-
 066 supervised learning and contrastive learning, where the in-
 067 troduction of difficult training examples has been shown to
 068 enhance representation learning [1, 6].

069 We demonstrate the effectiveness of this approach
 070 across multiple domains, showing that DA improves both
 071 generalization and robustness. Our results suggest that
 072 disagreement-driven augmentation can serve as a valuable
 073 tool in knowledge distillation, offering a novel perspective
 074 on how models can learn more efficiently from one an-
 075 other. Through this work, we aim to bridge the gap between
 076 classical pedagogical insights and modern machine learning
 077 methodologies, reinforcing the notion that structured con-
 078 flict—when properly harnessed—can be a powerful driver

079 of progress.

080 **2. Related Work**

081 Knowledge distillation, introduced by Hinton et al. [9], tra-
 082 ditionally aims to transfer knowledge from a large, well-
 083 trained teacher model to a smaller student model by min-
 084 imizing the divergence between their outputs. While this
 085 approach has been effective for model compression, recent
 086 work suggests that direct output matching may not always
 087 be optimal [10, 21]. The field has since evolved to recog-
 088 nize that valuable information exists not just in the teacher’s
 089 predictions but also in the underlying learning process. Ef-
 090 ficient knowledge transfer remains a key challenge, particu-
 091 larly under constraints of limited data or computational re-
 092 sources.

093 Recent studies have explored alternative approaches to
 094 distillation by considering additional aspects of model be-
 095 havior beyond simple output alignment [19]. For instance,
 096 Maroto et al. [14] demonstrated that knowledge distilla-
 097 tion can improve adversarial robustness, while Goldblum
 098 et al. [4] showed that adversarially robust teachers yield
 099 more resilient student networks. Of particular relevance to
 100 our work is the use of decision boundary information to en-
 101 hance distillation [8], which conceptually aligns with our
 102 disagreement-based augmentation method.

103 Distillation has also been investigated in the context of
 104 adversarial defense. Methods such as Adversarial Diffusion
 105 Distillation [18] and Adversarially Robust Distillation [4]
 106 demonstrate how transferring robustness properties from a
 107 teacher to a student can significantly improve model reli-
 108 ability. These approaches highlight the importance of de-
 109 cision boundaries in distillation, reinforcing the idea that
 110 structured exploration of disagreement can lead to better
 111 knowledge transfer.

112 In parallel, advancements in data-free model extraction
 113 have shown that student models can be trained without
 114 requiring access to the original training data. Truong et
 115 al. [20] introduced Data-Free Model Extraction (DFME),
 116 a technique that synthesizes queries to extract knowledge
 117 from black-box models. This method builds on data-free
 118 knowledge distillation by leveraging generative models to
 119 construct inputs that maximize disagreement between the
 120 teacher (victim) and student (stolen) models. Similarly,
 121 Fang et al. [3] proposed Data-Free Adversarial Distillation,
 122 which employs adversarial techniques to generate informa-
 123 tive samples for distillation without original data. Both ap-
 124 proaches align with our work by demonstrating how dis-
 125 agreement can guide knowledge transfer, reinforcing the
 126 role of structured model divergence in improving student
 127 learning.

128 Our work builds directly on the Committee Disagree-
 129 ment Sampling approach introduced by Goldfeder et al. [5].
 130 Their method identifies regions of the input space where

131 knowledge transfer is most needed by analyzing disagree-
 132 ment between multiple student models. While their work
 133 focused on exact parameter reconstruction, we adapt this
 134 technique for the more flexible problem of knowledge dis-
 135 tillation. This is closely related to adversarial sample gen-
 136 eration, where model disagreement often highlights decision
 137 boundary regions susceptible to adversarial attacks [13].
 138 By combining insights from disagreement-based learning
 139 and adversarial robustness, our method introduces a novel
 140 framework for enhancing knowledge transfer through struc-
 141 tured exploration of model differences.

142 **3. Methodology**

143 **3.1. Classification Experimental Setup**

144 We conducted our image classification experiments on the
 145 CIFAR-100 dataset [11], with three configurations of stu-
 146 dent/teacher pairs: Resnet8x4/Resnet32x4, VGG8/VGG13,
 147 and ShuffleNet-V2/Resnet32x4 [7, 19]. We used the origi-
 148 nal knowledge distillation method proposed by Hinton et
 149 al. [9], though our augmentation should be compatible with
 150 more modern techniques as well. The student model was
 151 trained to minimize the weighted sum of the knowledge dis-
 152 tillation loss and cross-entropy loss. Typical image augmen-
 153 tations were performed in both the baseline and DA exper-
 154 iments, such as random cropping and horizontal flipping.
 155 Experiments were run on a NVIDIA RTX 4090. All train-
 156 ing runs used an SGD optimizer, a batch size of 64, 240
 157 training epochs, an initial learning rate of 0.05, and learn-
 158 ing rate decay at epochs 150, 180 and 210. The learning rate
 159 here refers to the typical student learning rate, not the DA
 160 learning rate α . Both DA experiments and baselines with-
 161 out DA were run 5 times each to ensure statistical reliabil-
 162 ity, with results reported as the mean and standard deviation
 163 across these runs.

164 **3.2. Disagreement Augmentation Algorithm**

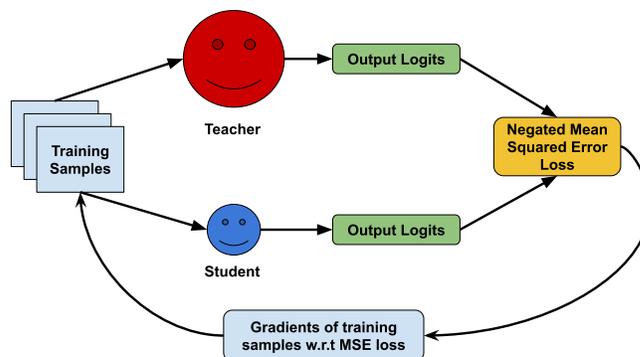


Figure 1. Schematic of the recursive DA algorithm. In practice only one epoch of augmentation occurs per batch.

165 The Disagreement Augmentation algorithm is designed to
 166 optimize input data by emphasizing areas of disagreement
 167 between a teacher model and a student model. The process
 168 begins by freezing the weights of both the teacher (T) and
 169 student (S) models to ensure that the augmentation process
 170 only modifies the input batch (I).

171 For each iteration of augmentation, the input batch is
 172 forward-propagated through the teacher and student mod-
 173 els to compute their respective output logits, denoted as L_T
 174 and L_S . These logits are then normalized to ensure they
 175 are on a comparable scale. The algorithm computes a dis-
 176 agreement loss l as the negative Mean Squared Error (MSE)
 177 between the normalized logits of the teacher and student:
 178 $l = -\text{MSE}(L_S, L_T)$. This loss function incentivizes maxi-
 179 mizing the discrepancy between the models’ predictions.

180 The disagreement loss is backpropagated to compute
 181 gradients with respect to the input batch I . These gradi-
 182 ents are then used to update I directly, employing a fixed
 183 learning rate α . This process is repeated for a predefined
 184 number of epochs e , iteratively refining the input batch to
 185 amplify disagreement between the models.

186 Once the iterations are complete, the optimized input
 187 batch I is returned as the final output of the algorithm, and
 188 used to train the student in typical knowledge distillation
 189 fashion. This approach ensures that the augmented data
 190 emphasizes areas where the teacher and student models di-
 191 verge, challenging the student model to learn more robust
 192 and generalizable features.

Algorithm 1 Disagreement Augmentation Algorithm

Require: Student S , teacher T , input batch I , learning rate
 α , epochs e
procedure DA(I, S, T, α, e)
 Freeze weights of S and T
for each epoch i in 1 to e **do**
 Forward-propagate I through S and T
 Compute logits: $L_S = S(I), L_T = T(I)$
 Normalize logits: $L_S \leftarrow \text{Normalize}(L_S), L_T \leftarrow$
 Normalize(L_T)
 Compute disagreement loss: $l =$
 $-\text{MSE}(L_S, L_T)$
 Back-propagate l and compute gradient w.r.t. I
 Update I using α
end for
 Return I
end procedure

193 **3.3. Policy Distillation Setup**

194 To extend our method to reinforcement learning environ-
 195 ments, we modified the original single-environment policy
 196 distillation methodology introduced by Rusu et al. [17].
 197 Our setup consists of three main stages: online data col-

lection, disagreement augmentation, and policy distillation. 198
 Both the student and teacher models are 4 layer deep Q- 199
 networks (DQNs), with 3 convolutional layers and one feed- 200
 forward layer [15]. The teacher DQN consists of 1.6 mil- 201
 lion paramaters, while the student has only 1% of that with 202
 roughly 16,000 (varies slightly across environments due to 203
 different action spaces). Students are trained for 500 epochs 204
 with a batch size of 32, a learning rate of 0.0001, and an 205
 SGD optimizer. 206

3.3.1. Online Data Collection 207

In this stage, a teacher pre-trained on an Atari environment 208
 is used to collect environment states. We used pre-trained 209
 teachers from RL Baselines3 Zoo [16]. The teacher inter- 210
 acts with the environment to generate trajectories, which 211
 are stored in a replay memory buffer. These stored states 212
 serve as the foundation for training the student. During 213
 each epoch of distillation, 54,000 environment states are 214
 generated, each consisting of 4 contiguous frames of the 215
 Atari game. The replay buffer has a capacity of 540,000 216
 states, which it maintains by removing excess states when 217
 new ones are added in a first-in-first-out manner. 218

To ensure diversity in the collected training data, we 219
 introduced a 5% exploration rate during data collection. 220
 Specifically, for each action taken by the teacher model, 221
 there is a 5% probability of selecting a random action in- 222
 stead of the teacher’s optimal policy decision. This con- 223
 trolled exploration helps capture a broader range of environ- 224
 ment states, including suboptimal transitions that can im- 225
 prove the robustness of the student model. 226

3.3.2. Disagreement Augmentation 227

Once a batch of environment states is retrieved from the re- 228
 play memory, we apply DA to emphasize areas in the state 229
 space where the student diverges from the teacher. The 230
 only difference between this instance of DA and what we 231
 used in the image classification setting is that we maxi- 232
 mize the Kullback-Leibler divergence (KLD) between the 233
 student and teacher Q-values rather than the mean squared 234
 error. We did this to remain consistent with the policy dis- 235
 tillation setting, as using KLD as the distillation loss is shown 236
 to improve performance over MSE [17]. For each batch of 237
 states, with probability $p = 0.3$ they undergo 1 epoch of 238
 DA with $\alpha = 0.001$ for Ms. Pacman and Space Invaders 239
 and $\alpha = 0.00001$ for Beam Rider. 240

3.3.3. Policy Distillation 241

Following the DA step, the policy distillation process takes 242
 place. The student policy is trained on the disagreement- 243
 augmented states, minimizing the KL divergence loss be- 244
 tween its action distribution and that of the teacher. Unlike 245
 Rusu et al., who used RMSProp, we instead optimize the 246
 student’s policy using SGD. This results in a training pro- 247
 cedure that is more sensitive to the nuanced differences be- 248

249 tween the teacher and student policies, ensuring improved
250 convergence dynamics.

251 **4. Experimental Results**

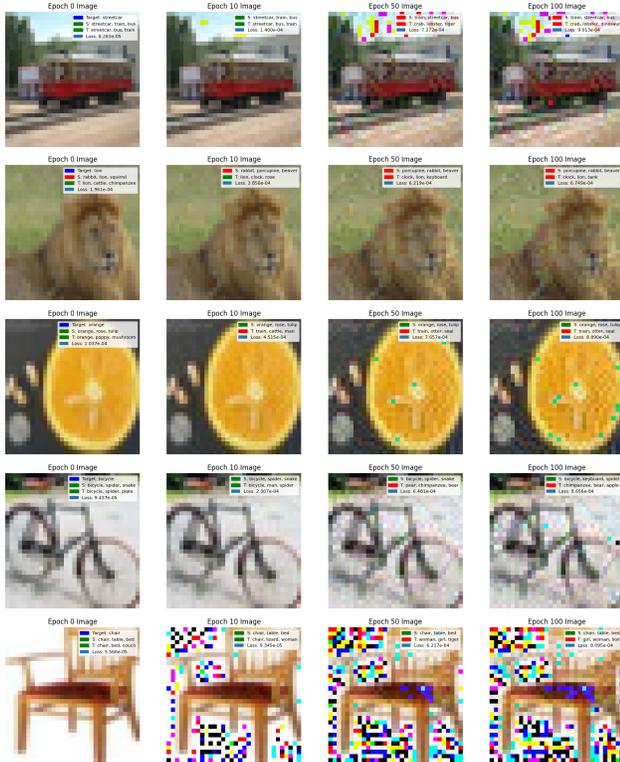


Figure 3. Examples of CIFAR-100 images undergoing various epochs of DA. The legend shows the ground truth target label, the Resnet8x4 student model’s top 3 predictions, the Resnet32x4 teacher model’s top 3 predictions, and the MSE loss between the student and teacher logits.

252 **4.1. Hyperparameter Search for Classification**

253 To optimize DA for our classification, we conducted a
254 Bayesian hyperparameter search with Hyperband early
255 stopping [12] over three hyperparameters: the number of
256 epochs of augmentation per batch e , the learning rate of
257 augmentation α , and the probability of augmentation per
258 batch p . The search was conducted using a Resnet32x4
259 teacher and a Resnet8x4 student, with the goal of maximiz-
260 ing student validation accuracy. It found the ideal paramet-
261 ers to be $e = 1$, $\alpha = 0.01778$, and $p = 0.7374$. These are
262 the parameters used in all classification experiments.

263 **4.2. Classification Results**

264 The results in Table 1 demonstrate that DA consistently
265 improves student model performance across different archi-
266 tectures. For the Resnet32x4 to Resnet8x4 transfer, DA
267 increases validation accuracy from $73.66\% \pm 0.26$ to

Table 1. Validation accuracy of baseline student models and student models trained with DA.

Teacher	Student	KD (%)	DA (%)
Resnet32x4	Resnet8x4	73.66 ± 0.26	74.59 ± 0.24
VGG13	VGG8	73.33 ± 0.25	73.76 ± 0.29
Resnet32x4	ShuffleNet-V2	71.67 ± 0.34	73.70 ± 0.19

268 $74.59\% \pm 0.24$, showing a clear performance gain. Simi-
269 larly, for the VGG13 to VGG8 transfer, DA yields a mod-
270 est improvement from $73.33\% \pm 0.25$ to $73.76\% \pm 0.29$.
271 The most significant relative improvement occurs in the
272 Resnet32x4 to ShuffleNet-V2 distillation, where DA raises
273 accuracy from $71.67\% \pm 0.34$ to $73.70\% \pm 0.19$, suggest-
274 ing that DA is particularly beneficial when distilling into
275 more compact, efficiency-oriented architectures. These re-
276 sults highlight that disagreement-driven augmentation pro-
277 vides a complementary boost to standard knowledge distil-
278 lation by encouraging more informative training dynamics.
279 The improvements observed across all tested student mod-
280 els suggest that DA is a robust and effective augmentation
281 strategy for classification tasks.

282 **4.3. Robustness to Disagreement Augmented Sam-
283 ples**

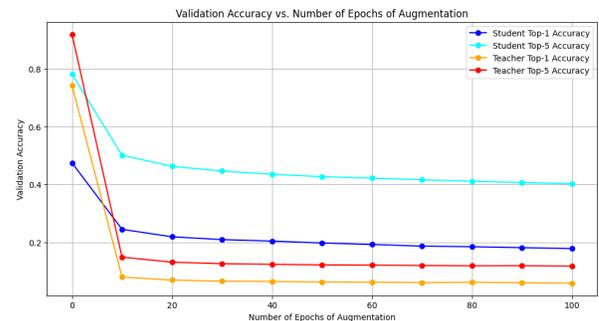


Figure 4. Validation accuracy on a DA augmented validation set vs. number of epochs of augmentation.

284 We hypothesized that training a student model with dis-
285 agreement augmented samples would result in a more ro-
286 bust model. To investigate, we evaluated the validation
287 accuracy of a pre-trained Resnet32x8 teacher and a DA-
288 trained Resnet8x4 student under varying levels of augmen-
289 tation intensity, measured by the number of augmentation
290 epochs. Here, augmentation occurs on the validation set
291 to ensure that the evaluation reflects whether training with
292 disagreement-augmented samples leads to improved robus-
293 tness against such perturbations. Additionally, we compared
294 the performance of the student model with its teacher to as-
295 sess whether the knowledge distillation process, combined
296 with disagreement-based augmentation, enables the student

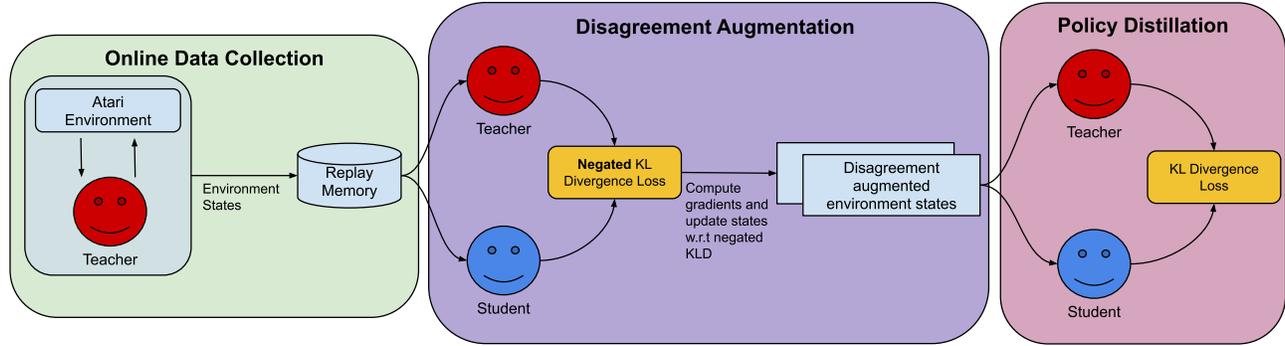


Figure 2. Policy distillation with Disagreement Augmentation.

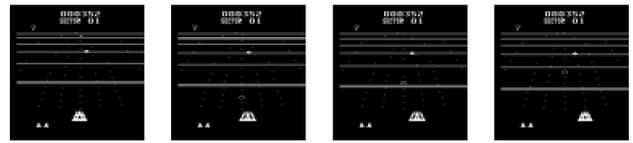
297 to achieve similar or superior resilience in handling these
 298 adversarial-like inputs. This approach allowed us to vali-
 299 date the hypothesis that disagreement-driven training fos-
 300 ters a more adaptable and robust student model.

301 **4.4. Policy Distillation Results**

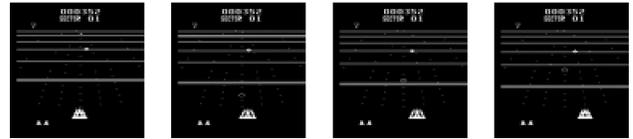
302 We evaluated DA in the policy distillation setting across
 303 three Atari environments: Beam Rider, Ms. Pacman, and
 304 Space Invaders. The goal of this experiment was to assess
 305 whether DA could improve the performance of a distilled
 306 student policy compared to standard behavior cloning using
 307 policy distillation.

308 The teacher policy in each case was a DQN model
 309 trained on the respective environment, while the student was
 310 a significantly smaller DQN model. The results are summa-
 311 rized in Table 2.

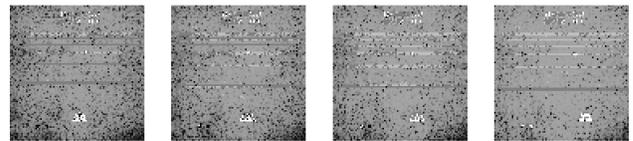
312 The results demonstrate that incorporating DA consis-
 313 tently improves student policy performance across all tested
 314 Atari environments. In Beam Rider, DA increases the
 315 student’s average score from 2105.46 to 2663.53, raising
 316 its relative performance from 51.62% to 65.30% of the
 317 teacher’s score. Similarly, in Ms. Pacman, DA enhances the
 318 student’s performance from 3094.92 to 3341.92, yielding a
 319 relative improvement from 117.44% to 126.81%. The most
 320 notable increase occurs in Space Invaders, where DA raises
 321 the student’s performance from 585.24 to 637.02, boost-
 322 ing the relative score from 112.24% to 122.17%. These
 323 results suggest that disagreement-driven augmentation en-
 324 hances policy distillation by exposing the student to more
 325 informative training samples, leading to improved gener-
 326 alization and robustness.



4 frame stack from Beam Rider after typical preprocessing for Atari environments (downsized to 84×84 , grayscale).



The same 4 frame stack after one iteration of DA, with $\alpha = 0.00001$.



Absolute difference between original and augmented frames, normalized to the range $[0, 255]$ for visualization. This is equivalent to the (normalized) absolute value of the gradient of the input frames with respect to the negated KLD loss.

Figure 5. Example of how DA affects environment states. This was generated with a pre-trained teacher and a student trained with DA.

5. Discussion

5.1. Interpretation of Results

The results of our experiments demonstrate that incorporat-
 ing DA into the knowledge distillation process significantly
 improves the generalization and robustness of student mod-
 els. Across all tested configurations, models trained with
 DA consistently outperformed their baseline counterparts in
 both classification and reinforcement learning tasks. This
 suggests that the structured introduction of disagreement

Table 2. Policy distillation results across different Atari environments. Scores are calculated as average scores over 1000 episodes. The relative score represents student performance as a percentage of the teacher’s performance.

Env	Method	Teacher Score	Student Score	Relative Score
Beam Rider	PD	4078.726	2105.460	51.62%
Beam Rider	PD + DA	4078.726	2663.526	65.30%
Ms. Pacman	PD	2635.370	3094.920	117.44%
Ms. Pacman	PD + DA	2635.370	3341.920	126.81%
Space Invaders	PD	521.405	585.235	112.24%
Space Invaders	PD + DA	521.405	637.020	122.17%

336 during training helps the student model better learn nuanced
337 representations of the teacher’s decision boundaries.

338 In the classification experiments, DA led to improved
339 validation accuracy across all student architectures, with
340 particularly strong gains in compact models such as
341 ShuffleNet-V2. These improvements indicate that DA is
342 especially beneficial for lightweight models, where stan-
343 dard distillation may struggle to fully capture the teacher’s
344 knowledge.

345 The reinforcement learning experiments further high-
346 light the effectiveness of DA in policy distillation. In
347 Atari environments, DA consistently improved student pol-
348 icy performance across all tested games, increasing rela-
349 tive student scores compared to standard policy distillation.
350 Notably, the largest improvements were observed in Beam
351 Rider and Space Invaders, where DA enhanced the student’s
352 ability to generalize across diverse game states. These re-
353 sults suggest that DA is not only beneficial in supervised
354 learning but also in reinforcement learning settings, where
355 effectively transferring policy knowledge remains a major
356 challenge.

357 Furthermore, the robustness evaluation confirmed our
358 hypothesis that disagreement-driven training fosters re-
359 siliance to adversarial-like inputs. By augmenting the val-
360 idation set to contain disagreement-optimized samples, we
361 observed that DA-trained students were better equipped to
362 reconcile these challenging inputs, achieving performance
363 levels comparable to or surpassing their teachers. The im-
364 provements observed in both classification and reinforce-
365 ment learning settings demonstrate that DA provides a gen-
366 eralizable mechanism for improving knowledge transfer.

367 5.2. Comparison with Previous Studies

368 Our findings align with and expand upon prior work that
369 has explored the role of adversarial robustness in knowledge
370 distillation. Although earlier studies, such as Goldblum et
371 al. [4], demonstrated the benefits of robust teachers for im-
372 proving student resilience, our method extends this concept
373 by actively incorporating disagreement between models as
374 a training signal. Compared to approaches like adversari-

ally robust distillation, DA introduces a more generalizable
framework that does not rely on predefined attack methods
but instead leverages natural divergences between teacher
and student predictions. This positions DA as a comple-
mentary and scalable strategy for enhancing robustness in
distillation tasks.

In reinforcement learning, previous work on policy dis-
tillation has primarily focused on directly matching teacher
policies [17], often struggling to capture uncertainty or out-
of-distribution states effectively. Our results suggest that
introducing structured disagreement in policy distillation
improves knowledge transfer, potentially helping student
policies generalize beyond trajectories demonstrated by the
teacher. This complements recent studies on uncertainty-
aware policy distillation, reinforcing the idea that controlled
divergence can be a useful signal in both supervised and re-
inforcement learning settings.

5.3. Challenges and Limitations

Despite its promising results, DA is not without chal-
lenges. One limitation is the additional computational cost
incurred during the augmentation process, as optimizing
input batches over multiple epochs introduces overhead.
While this cost was manageable in our experiments with
CIFAR-100, scaling to larger datasets or models may re-
quire further optimization of the augmentation procedure.
Similarly, in reinforcement learning environments, gener-
ating disagreement-optimized samples requires additional
exploration, which can slow down training if not carefully
managed.

Another limitation is the reliance on hyperparameter tun-
ing to achieve optimal performance. As shown in our hyper-
parameter search, the number of augmentation epochs (e),
learning rate (α), and probability of augmentation (p) are
critical to the success of DA. While our method performed
well across multiple settings, the need for manual tuning
may limit accessibility. Automating or simplifying this tun-
ing process could improve the scalability of the method,
particularly for reinforcement learning applications where
hyperparameter sensitivity is often high.

414 **5.4. Future Directions**

415 Future work could address the computational challenges of
416 DA by exploring methods to reduce augmentation overhead,
417 such as adaptive augmentation strategies that selectively ap-
418 ply disagreement-based modifications based on confidence
419 thresholds. Additionally, while our reinforcement learning
420 experiments demonstrated the benefits of DA in Atari envi-
421 ronments, further research is needed to assess its effective-
422 ness in more complex RL tasks, such as continuous control
423 or multi-agent settings.

424 Another promising direction is extending DA to other
425 domains, such as natural language processing (NLP) or
426 self-supervised learning, where structured disagreement
427 could help improve representation learning. For example,
428 disagreement-based augmentation could be applied to NLP
429 models by modifying token embeddings to create diverse
430 training sequences, potentially leading to better generaliza-
431 tion in text classification and translation tasks.

432 Finally, investigating the theoretical underpinnings of
433 disagreement as a learning signal, particularly in the con-
434 text of decision boundary exploration, could further refine
435 and justify the approach. A deeper understanding of why
436 and when DA is most effective could help develop more
437 principled augmentation strategies that generalize across a
438 broader range of learning tasks.

439 **6. Conclusion**

440 **6.1. Summary of Findings**

441 This work introduced Disagreement Augmentation (DA), a
442 novel method for improving knowledge distillation by in-
443 tentively optimizing the input to maximize disagreement
444 between teacher and student models. Inspired by the So-
445 cratic method, DA leverages structured conflict to challenge
446 the student model, encouraging it to develop more robust
447 and generalizable representations.

448 Experimental results on CIFAR-100 demonstrated
449 that DA-trained students consistently outperformed base-
450 line models in validation accuracy and robustness to
451 disagreement-augmented samples. Furthermore, extending
452 DA to reinforcement learning environments showed that
453 disagreement-driven augmentation significantly enhances
454 policy distillation. In Atari games, DA improved student
455 policies across all tested environments, increasing their
456 ability to generalize beyond the teacher’s demonstrated tra-
457 jectories. These results suggest that DA is a versatile aug-
458 mentation strategy applicable to both supervised and rein-
459 forcement learning tasks.

460 **6.2. Contributions**

461 Our primary contributions are as follows:

- 462 • The introduction of Disagreement Augmentation as a
- 463 generalizable data augmentation strategy for knowledge

- distillation across both classification and reinforcement 464
learning. 465
- Empirical validation of DA’s effectiveness, demonstrating 466
improved generalization and robustness across multiple 467
teacher-student configurations in classification tasks and 468
enhanced policy transfer in reinforcement learning. 469
- A conceptual shift in knowledge distillation, emphasizing 470
the role of structured disagreement as a catalyst for 471
learning. 472

Acknowledgment 473

This work was built on [this codebase](#) [22] 474
[23]. 475

References 476

[1] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Ge- 477
offrey Hinton. A simple framework for contrastive learning 478
of visual representations. In *International conference on ma- 479
chine learning*, pages 1597–1607. PmLR, 2020. 1 480

[2] Ekin D. Cubuk, Barret Zoph, Dandelion Mane, Vijay Va- 481
sudevan, and Quoc V. Le. Autoaugment: Learning augmen- 482
tation strategies from data. In *Proceedings of the IEEE/CVF 483
Conference on Computer Vision and Pattern Recognition 484
(CVPR)*, 2019. 1 485

[3] Gongfan Fang, Jie Song, Chengchao Shen, Xinchao Wang, 486
Da Chen, and Mingli Song. Data-free adversarial distillation. 487
arXiv preprint arXiv:1912.11006, 2019. 2 488

[4] Micah Goldblum, Liam Fowl, Soheil Feizi, and Tom Gold- 489
stein. Adversarially robust distillation. In *Proceedings of the 490
AAAI conference on artificial intelligence*, pages 3996–4003, 491
2020. 2, 6 492

[5] Judah Goldfeder, Quinten Roets, Gabe Guo, John Wright, 493
and Hod Lipson. Sequencing the neurome: Towards scal- 494
able exact parameter reconstruction of black-box neural net- 495
works. *arXiv preprint arXiv:2409.19138*, 2024. 2 496

[6] Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin 497
Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, 498
Bernardo Avila Pires, Zhaohan Guo, Mohammad Ghesh- 499
laghi Azar, et al. Bootstrap your own latent—a new approach 500
to self-supervised learning. *Advances in neural information 501
processing systems*, 33:21271–21284, 2020. 1 502

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 503
Deep residual learning for image recognition. In *Proceed- 504
ings of the IEEE conference on computer vision and pattern 505
recognition*, pages 770–778, 2016. 2 506

[8] Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young 507
Choi. Knowledge distillation with adversarial samples sup- 508
porting decision boundary. In *Proceedings of the AAAI con- 509
ference on artificial intelligence*, pages 3771–3778, 2019. 2 510

[9] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distill- 511
ing the knowledge in a neural network. *arXiv preprint 512
arXiv:1503.02531*, 2015. 1, 2 513

[10] Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan 514
Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, 515
Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! 516

- 517 outperforming larger language models with less training data
518 and smaller model sizes. *arXiv preprint arXiv:2305.02301*,
519 2023. 2
- 520 [11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple
521 layers of features from tiny images. 2009. 2
- 522 [12] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Ros-
523 tamizadeh, and Ameet Talwalkar. Hyperband: A novel
524 bandit-based approach to hyperparameter optimization.
525 *Journal of Machine Learning Research*, 18(185):1–52, 2018.
526 4
- 527 [13] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt,
528 Dimitris Tsipras, and Adrian Vladu. Towards deep learn-
529 ing models resistant to adversarial attacks. *arXiv preprint*
530 *arXiv:1706.06083*, 2017. 2
- 531 [14] Javier Maroto, Guillermo Ortiz-Jiménez, and Pascal
532 Frossard. On the benefits of knowledge distillation for adver-
533 sarial robustness. *arXiv preprint arXiv:2203.07159*, 2022. 2
- 534 [15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, An-
535 drei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves,
536 Martin Riedmiller, Andreas K Fidfjeland, Georg Ostrovski,
537 et al. Human-level control through deep reinforcement learn-
538 ing. *nature*, 518(7540):529–533, 2015. 3
- 539 [16] Antonin Raffin. RL baselines3 zoo. [https://github.](https://github.com/DLR-RM/rl-baselines3-zoo)
540 [com/DLR-RM/rl-baselines3-zoo](https://github.com/DLR-RM/rl-baselines3-zoo), 2020. 3
- 541 [17] Andrei A Rusu, Sergio Gomez Colmenarejo, Caglar Gul-
542 cehre, Guillaume Desjardins, James Kirkpatrick, Raz-
543 van Pascanu, Volodymyr Mnih, Koray Kavukcuoglu, and
544 Raia Hadsell. Policy distillation. *arXiv preprint*
545 *arXiv:1511.06295*, 2015. 3, 6
- 546 [18] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin
547 Rombach. Adversarial diffusion distillation. In *European*
548 *Conference on Computer Vision*, pages 87–103. Springer,
549 2024. 2
- 550 [19] Karen Simonyan and Andrew Zisserman. Very deep convo-
551 lutional networks for large-scale image recognition. *arXiv*
552 *preprint arXiv:1409.1556*, 2014. 2
- 553 [20] Jean-Baptiste Truong, Pratyush Maini, Robert J Walls, and
554 Nicolas Papernot. Data-free model extraction. In *Proceed-*
555 *ings of the IEEE/CVF conference on computer vision and*
556 *pattern recognition*, pages 4771–4780, 2021. 2
- 557 [21] Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold
558 Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou.
559 A survey on knowledge distillation of large language models.
560 *arXiv preprint arXiv:2402.13116*, 2024. 1, 2
- 561 [22] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun
562 Liang. Decoupled knowledge distillation. In *Proceedings of*
563 *the IEEE/CVF Conference on computer vision and pattern*
564 *recognition*, pages 11953–11962, 2022. 7
- 565 [23] Borui Zhao, Quan Cui, Renjie Song, and Jiajun Liang.
566 Dot: A distillation-oriented trainer. *arXiv preprint*
567 *arXiv:2307.08436*, 2023. 7